

An LALR Parser Generator Supporting Conflict Resolution

Leonardo Teixeira Passos

(Federal University of Minas Gerais, Belo Horizonte, Brazil
leonardo@dcc.ufmg.br)

Mariza A. S. Bigonha

(Federal University of Minas Gerais, Belo Horizonte, Brazil
mariza@dcc.ufmg.br)

Roberto S. Bigonha

(Federal University of Minas Gerais, Belo Horizonte, Brazil
bigonha@dcc.ufmg.br)

Abstract: Despite all the advance brought by LALR parsing method by DeRemer in the late 60's, conflicts continue to be removed in a non-productive way, by means of analysis of a huge amount of textual and low level data dumped by the parser generator tool. For the purpose of changing this scenario, we present a parser generator capable of automatically removing some types of conflicts, along with a supported methodology that guides the process of manual removal. We also discuss the internal algorithms and how the created parsers are compact in terms of memory usage.

Key Words: lalr parsing, automatic conflict removal, table compression, methodology

Category: D.3.4, F.4.2

1 Introduction

The usual way to build a bottom-up syntax analyzer is by writing syntax specifications that are readable by parser generator tools, such as YACC, Bison, CUP, etc. These tools can only work with LALR(1) grammars, otherwise conflicts are reported, i.e., non-determinism points resulted from the grammar. The number of conflicts in a programming language can easily reach hundreds, if not more than a thousand conflicts. To illustrate how frequent conflicts are, the original grammars of Algol-60, Scheme and Notus[Tirelo and Bigonha 2006] programming languages result in 61, 78 and 575 conflicts, respectively. In the first two grammars, the average density is one conflict for each two productions; in the latter there are two per production.

To remove the conflicts in a syntax specification, one must inspect each conflict. Parser generators assist users in conflict removal by writing log files. These files consist of pure text data such as the list of conflicts and the LALR(1) automaton. Log files, however, tend to be too extensive for analysis. For the Notus

programming language, Bison creates a log file of 54 Kb in size, having 6,244 words and 2,257 lines. Besides its size, log files contain a low level of abstraction in their content, requiring expertise in LALR inner working. All together, these factors contribute to a low productivity when building LALR parsers.

In order to change the current scenario, we present in this article an LALR parser generator that automatically removes conflicts and supports a methodology to guide the process in cases of manual removal. Such methodology is the result of new techniques that extend the original work in [Passos et al. 2007]. In addition, we also discuss how the proposed tool generates parsers with little memory requirements.

This article is organized as follows: Section 2 discusses the types of conflicts, Section 3 presents the proposed parser generator, with its algorithms explained in Section 4. Section 5 presents some experimental results and Section 6 concludes the article.

2 Types of Conflicts

A conflict, either a shift/reduce or reduce/reduce, is reported by the parser generator by two reasons: lack of right context and ambiguity.

Conflicts caused by lack of right context indicate that the number of inspected lookaheads is not enough to decide which action to execute – shift or reduce within a set of possible reductions. Part of these conflicts can be removed if the value of k is increased accordingly. This results in an LALR(k) parser, where $k = \max\{k', \text{such that } k' \text{ is the number of lookaheads needed to solve a given conflict among all conflicts reported by the parser generator and } k \text{ is not infinite}\}$. This solution, although correct, is unfeasible in most cases. If we take, for example, the LALR(1) parsing table of Visual Basic .NET¹ and make it LALR(3), the result matrix would contain 1,899,085,824 entries. The other part of this set of conflicts can only be removed by rewriting the grammar, assuming the language in question is indeed LALR. Such cases result from an infinite amount of lookaheads necessary to solve a given conflict.

Conflicts caused by ambiguity must be removed by rewriting the grammar. Again, this can only be performed if the language in question is LALR. Some parser generators deal with ambiguity by means of precedence and associativity or simply by performing a shift in a shift/reduce conflict.

A discussion of all types of conflicts can be found in [Passos et al. 2007].

¹ This test was performed using the grammar provided in <http://www.devincook.com/goldparser/grammars/index.html>.

3 Proposed Tool

The proposed tool, named SAIDE², is an integrated development environment with an internal parser generator. Figure 1 shows the overall appearance of SAIDE's graphical interface. The text editor window is located in the upper left corner, loaded with a syntax specification. After asking for the validation of the grammar, the user starts the main cycle of the methodology supported by the tool. The main cycle is divided in two major steps: automatic and manual conflict removal.

In the automatic removal, SAIDE tries to remove all conflicts without user intervention. The non-removed conflicts are then listed to the user. This listing is performed using a heuristic that sorts all conflicts considering the order in which they must be removed. A conflict must be listed before those that appear as a consequence of the existence of the first. In order to calculate such removal priority, SAIDE needs to know the whole set of conflicts. In Figure 1, the listing of conflicts is shown below the editor.

After the listing, the manual removal step starts. According to the methodology, to manually remove a conflict one must go through four phases: (i) understanding; (ii) classification; (iii) editing and (iv) testing.

In the understanding phase, the user tries to deduce the cause of the conflict using derivation trees. This has the advantage of manipulating a more intuitive and higher level structure compared to the low level data available in log files. Derivation trees are presented after the user clicks on the *Debug conflict* option, shown as a hyperlink below the conflict's item set. For expert users, low level content is still available, as can be seen by the LALR automaton shown next to the editor window, in Figure 1.

In the classification phase the user defines the category in which the conflict belongs, i.e., determine whether a conflict is due to lack of right context or ambiguity. In the latter case, a catalog of some well known ambiguity constructions, along with their solutions, is available for consultancy and can be extended with user defined entries. At this phase, the user must define a strategy to rewrite the grammar so the given conflict can be removed. The identification of the conflict's category adds confidence, as we expect that a strategy used in removing a past conflict can be applied many times to other conflicts in the same category.

Next, the user edits the grammar in order to apply the strategy defined in the last phase and submits the specification to be validated. The main cycle of the methodology is then restarted and continues until no conflicts are reported.

² SAIDE /said/: Syntax Analyzer Integrated Development Environment.

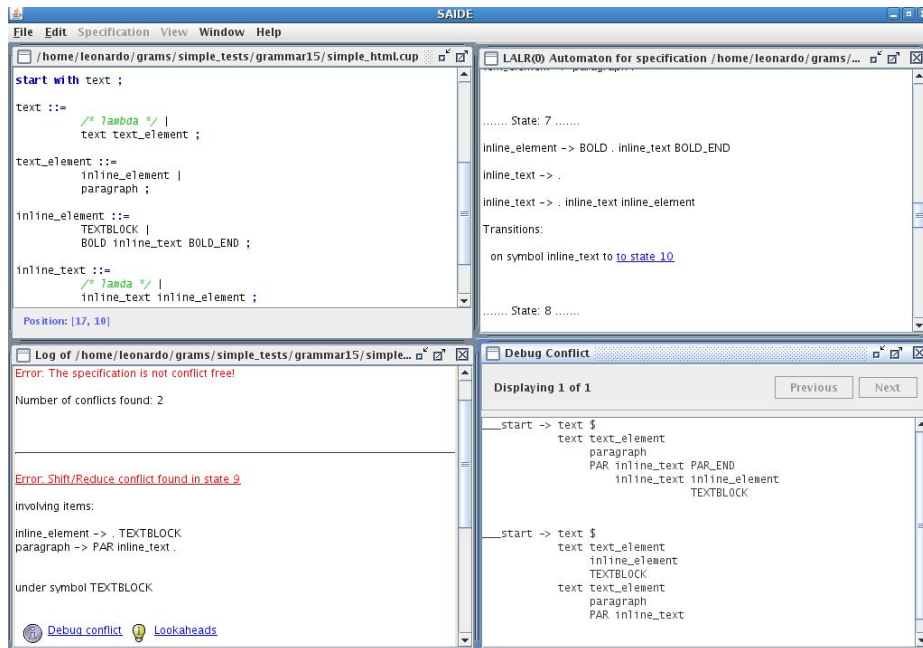


Figure 1: SAIDE's main window.

4 Automatic Conflict Resolution

Deviating from conflict resolution based on nondeterminism like Generalized LR Parsing [Tomita 1991], we address an automatic and deterministic resolution approach. In this section we give an overview of the technique proposed by Charles [Charles 1991], capable of removing some conflicts caused by lack of right context.

Charles suggests LALR(k) generation as a mechanism to remove conflicts. To perform this while decreasing storage needs, he proposes the extension of the number of inspected lookaheads only when necessary to uniquely define which parse action to execute. Thus, the parser inspects up to k tokens before choosing a parse action. The value of k is limited by a constant k_{max} . This approach has the advantage of reducing the number of entries in the LALR parsing table when compared to complete LALR(k) tables. These parsers will be referred as LALR(k_v), with k_v denoting the variable length characteristic of k .

To build the LALR(k_v) parser, each entry (p, a) in the action table with at least one conflict will become the start state of a deterministic finite automaton (DFA). From p following the conflict symbol a , another state, represented as q , is

reached. For all tokens t_1, \dots, t_n that may follow a , the parsing actions in (q, t_i) are properly set. If at this point, all (q, t_i) contain cardinality equal to one, then all conflicts have been removed. Otherwise, if there is a least one entry (q, t_i) whose cardinality is greater than one, another level of lookaheads is calculated. These lookaheads consist of the tokens that may follow t_i , given a conflictual entry (q, t_i) . This process continues until the conflicts are removed or the depth of the DFA reaches k_{max} . The described scheme is depicted in terms of an example in Figure 2. Suppose, that the original LALR parsing table entry (q, a) has the following actions: $\{S2, R5, R6\}$. In this entry, there are three conflicts: two shift/reduces $(\{S2, R5\}, \{S2, R6\})$ and a reduce/reduce $(\{R5, R6\})$. If k_{max} is taken as three, q is made the start state of the DFA. From state q following a , the destination q_1 is reached. If b or c follow a , then a unique parsing action is determined. At this point, the shift/reduce conflicts vanish. The entry, (q_1, d) still reports $\{R5, R6\}$, and the automaton is extended in another level of lookaheads. The tokens that can be found after d are $\{b, c\}$. Extending d with these tokens uniquely determines each reduction. Note that although k_{max} was three, in q_1 only two lookaheads were used.

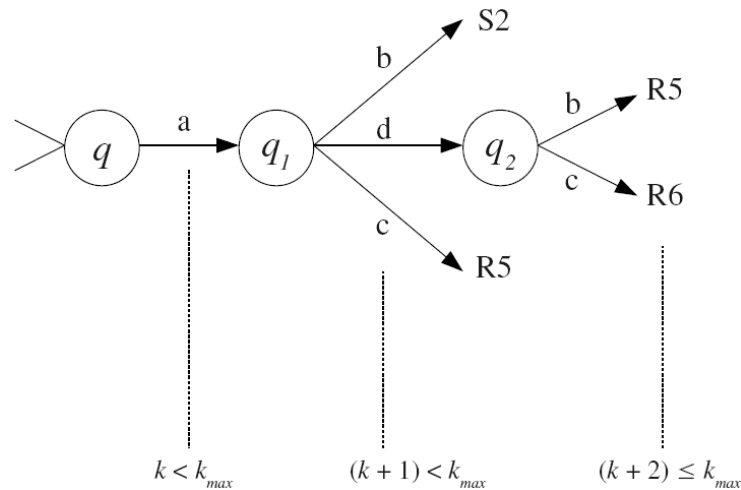


Figure 2: DFA to solve a given conflict by extending the number of lookaheads respected the limit $k_{max} = 3$.

The action parsing table of an LALR(k_v) parser is encoded as follows: each action that points to a DFA becomes a *lookahead action* – Ln , where n is the first line available in the table. The rest of the DFA’s transition table is appended starting from the n -th line of the action parsing table. Note that this coding

scheme preserves the original layout of the LALR(1) action parsing table, since the columns are still indexed by a single token, in contrast with LALR(k) encoding method. To illustrate the produced table, consider the following grammar:

$$\begin{aligned} bnf &\rightarrow rlist \\ rlist &\rightarrow rlist \text{ rule} \\ &\mid \lambda \\ rule &\rightarrow \mathbf{s} \text{ “}\rightarrow\text{” } slist \\ slist &\rightarrow slist \mathbf{s} \\ &\mid \lambda \end{aligned}$$

This grammar is not LALR(1); it is LALR(2). The following action parsing table results from the LALR(k_v) approach:

Action	s	“→”	“\$”
	0 R3		R3
	1 S3		R1
	2		ACC
	3	S5	
	4 R2		R2
	5 R6		R6
	6 L8		R4
	7 R5		R5
Lookaheads	8 S7	R4	S7

An LALR(k_v) parser works almost as an LALR(k) parser. The main difference occurs when dealing with lookahead actions. Given an entry Ln in position (q, a) in the action table, the parser switches to the appropriate DFA table. By consulting the tokens that may follow a in the input, the parser tries to find a path from the current DFA state that leads to a non-error parse action (different from blank). This possibly implies in making other lookahead actions, that differ from conventional shift actions in the sense that the tokens consulted in the input are not consumed in any way.

Charles' algorithms to calculate the action tables for LALR(k_v) parsers restrict the input grammars to the set of grammars that: (a) do not have circularity, i.e., given a nonterminal A , A derives A in one or more steps; (b) do not imply in the existence of strong connected components (SCCs) in the *reads* relation graph. The *reads* graph represents a relation between transitions³ [DeRemer and Pennello 1982]:

$$(p, A) \text{ reads } (q, B) \text{ iff } \text{GOTO}_0(p, A) = q \text{ e } B \xrightarrow{*} \lambda$$

³ (p, A) denotes a transition under A in state p in the LR(0) automaton.

where $GOTO_0$ is the transition function of the LR(0) automaton, defined over $(M_0 \times V) \rightarrow M_0$. The set M_0 contains all LR(0) states and V is the vocabulary of the grammar, made of the union of terminal (Σ) and nonterminal symbols (N). In case a grammar does not fit in such characteristics the algorithm stops (such grammars will be hereafter referred as *NLALR grammars*). If continuation is performed, Charles states that there might be a chance of entering infinite loop.

5 SAIDE's Internal Algorithms

This section presents the core algorithms we propose to support the methodology discussed in Section 3. For those already discussed in [Passos et al. 2007], no presentation is made. The algorithms are organized in three parts: conflict removal, conflict listing and table compression.

5.1 Conflict Removal

Charles calculates the lookaheads necessary to extend a given token by simulating the steps of the LR(0) automaton. His algorithms stop as soon as the grammar is found to be NLALR, ignoring the possibility of having non-inspected automatic solvable conflict entries in the action table. If such entries exist, given the LALR(k_v) approach, users face a number of conflicts that is greater than the real quantity. In addition, they may spend time manually removing conflicts that were supposed to disappear automatically. In this scenario, SAIDE's methodology as originally proposed becomes inapplicable.

SAIDE overcomes this by establishing the genuine set of conflicts that cannot be solved. It uses six algorithms that are the result of modifications over the ones originally proposed by Charles. Due to a proper infinite loop control, all proposed algorithms are guaranteed to terminate, even in the presence of NLALR grammars.

The whole process of solving conflicts starts with SWEEP, shown in Figure 3. The SWEEP algorithm must be called for each state p that contains a conflict. This procedure inspects the number of entries in the pairs (p, a) ⁴ of the action matrix, where a is a terminal symbol. If (p, a) has cardinality greater than two, it determines all possible stacks (sources), given the initial stack $[p]$, that result in reading a . If (p, a) contains a shift action, then $[p]$ is a valid stack. Reduce actions must also be considered. If a reduction by $A \rightarrow \omega$ belongs to (p, a) , then the start stack is one made of a predecessor state of p under ω . Function PRED returns such predecessor states. In both cases, the stacks are stored in the *sources* dictionary⁵ as part of pairs (stk, w) , where w corresponds to the string of

⁴ The functions FST and SND will be used to retrieve the first and the second component of a pair.

⁵ The *sources* dictionary is indexed by parsing actions.

```

SWEEP( $p$ )
1  for  $a \in \Sigma$ 
2  do if  $|Action[p, a]| > 0 \wedge a \neq \$$ 
3      then  $sources \leftarrow \emptyset$ 
4          for  $act \in Action[p, a]$ 
5              do if  $act$  is a shift action
6                  then  $sources[act] \leftarrow \{([p], a)\}$ 
7                  else ASSERT( $act = \text{reduce by rule } A \rightarrow \alpha$ )
8                      for  $p_0 \in PRED(p, \alpha)$ 
9                          do FOLLOW-SOURCES $_A(sources[act],$ 
10                              $[p_0], A, a, \lambda)$ 
11                             RESOLVE-CONFLICTS( $p, a, sources, 2$ )

```

Figure 3: Main procedure to automatically remove conflicts.

tokens that would have been read by the LR(0) automaton having the stack stk . Calling FOLLOW-SOURCES $_A$ completes the initial set of pairs (stk, w) . After the definition of all sources, RESOLVE-CONFLICTS is executed in an attempt to remove the conflicts in (p, a) .

The procedure FOLLOW-SOURCES $_A$, shown in Figure 4, is a façade procedure to FOLLOW-SOURCES $_B$. Five arguments are received: the set of pairs $(stack, w)$ – stored in $sources$, the current stack, the current symbol $X \in V$ denoting the transition that must be performed from the top state of the current stack, a marking the terminal symbol that has to be eventually found as a transition symbol, and w as the current processed string. Before FOLLOW-SOURCES $_B$ starts simulating the LR(0) steps in order to find all stacks that will lead to the reading of a , FOLLOW-SOURCES $_A$ puts the current stack in a tree format. The idea of using a tree structure comes from [Kristensen and Madsen 1981] and is used as a way to prevent infinite loop. Each node in the tree stores, besides its list of children, a pair $(state, z)$ as its value. Given a node n , the string z stands as the string processed by the LR(0) automaton given the states from the root of the tree to the node n . In the built tree, each node has a unique numeric identifier. Instantiating nodes using NODE controls such uniqueness. The value and the identifier of a node can be retrieved at any time by calling VALUE and ID, respectively.

FOLLOW-SOURCES $_B$, presented in Figure 5, aims to find from a given initial stack, encoded as a path in the tree of states, the set of stacks that will lead to the reading of a . Eight arguments must be received: the set of sources, as in FOLLOW-SOURCES $_A$, the current transition to be performed as part of the LR(0) simulation – structured as a triple $(source\text{-}state, transition\text{-}symbol, desti-$

nation-state), the terminal a , the current processed string w , the root of the tree, the current node tree and the sets *visited* and *roots*. When a appears as a transition symbol, the procedure stores in *sources* the pair $(stack, wa)$, where *stack* is given by the states in the values of the nodes in the path in the tree from *root* to *node*. The first part of infinite loop control in FOLLOW-SOURCES_B is located in lines 2 – 8. At that point, each time the procedure is activated, FOLLOW-SOURCES_B checks if the current stack and transition were visited in a past time. Instead of storing the whole sequence of states on the stack, it is enough to know the current node's identifier, for it also defines a unique stack. For optimization purposes, the pairs $(id, transition)$ are stored in two distinct sets: *roots* and *visited*. The *roots* set is only used for stacks whose size is one. The size of the current stack is obtained by calling HEIGHT, which returns the height from the root node to the current node. Note that if the size is one, storing the node's identifier is not necessary, since the root state corresponds to the source state in the *transition* triple. For stacks greater than one in size, *visited* is used. Lines 10 – 25 inspect all edges leaving q , where q is the destination state of the *transition* parameter. Two possibilities arise: a transition over a nullable symbol or a terminal symbol. From the first we must continue the simulation by pushing q onto the stack. To do this, we add a new child node to the current node's children list. The created node stores a pair (q, w) as its value. Next, FOLLOW-SOURCES_B is recursively called. The second situation occurs when there is a transition symbol from q that matches a . In this case, the current stack is retrieved by getting the inverse order of the states from the current node to the root of the tree, and the pair $(stack, wa)$ is added to *sources*. Again, infinite loop control is performed in the 11th line. To aid this control, we use the GET-FROM function. Given a node n_x and a value v , GET-FROM returns a non-nullable reference to a node n_y in the path from n_x (inclusively) to the root of the tree whose value coincide with v . If so, pushing q onto the stack results in a cycle: $([p_1 p_2 \dots p_n q], w) \vdash^* ([p_1 p_2 \dots p_n q q_1 q_2 \dots q_n q], w)$. The lines 27 – 40 are responsible for making reductions. Reductions while simulating LR(0) may lead to underflow situations, i.e., the act of popping more states than currently available on the stack. If $|\gamma_1 \gamma_2|$ states should be popped, but only $|\gamma_2|$ are available on the stack, being $\gamma_1 \gamma_2$ the right hand side of a production, the predecessor state of the γ_1 is retrieved and put as the top element of an unitary stack, used as a parameter to a recursive call.

The procedure RESOLVE-CONFLICTS, shown in Figure 6, checks if a conflict is removed. If not, it extends the DFA in another level of lookaheads, respected k_{max} . Four arguments are mandatory: a state q containing conflicts under t , also a received parameter, the *sources* dictionary (as in SWEEP and FOLLOW-SOURCES_A) and n , the number of lookaheads used so far. Its execution starts checking if n is greater than k_{max} or t is the EOF marker. In

```

FOLLOW-SOURCESA(sources, stack, X, a, w)
1  ASSERT(stack = [p1...pn])
2  root ← node ← NODE((p1, λ))
3  for 2 ≤ i ≤ n
4  do node2 ← NODE((pi, λ))
5     ADD-CHILD(node, node2)
6     node ← node2
7  SND(VALUE(node)) ← w
8  FOLLOW-SOURCESB(sources, (pn, X, GOTO0(pn, X)), a, w,
9  root, node, ∅, ∅)

```

Figure 4: Façade procedure FOLLOW-SOURCES_A.

either case, the extension of lookaheads should not go further and the procedure returns. Otherwise, a new line p is allocated in the action table, after its last line, and each entry in p is given the empty set. Later, the conflict entry (q, t) points to p by a lookahead action $-Lp$. Next, for each action indexing $sources$, each source $(stack, w)$ is inspected. Each token than can follow t , given $stack$, is calculated by calling NEXT-LOOKAHEADS_A. For each returned token a , the appropriate actions are put into (p, a) in the action table. After determining the values in p 's entries, for the entries whose cardinality is greater than one, the conflict removal process continues by calling FOLLOW-SOURCES_A. This is necessary, since NEXT-LOOKAHEADS_A returns the tokens that can extend t , but not the context in which they were obtained (source stacks).

Analogous to FOLLOW-SOURCES_A, NEXT-LOOKAHEADS_A, shown in Figure 7, is a façade procedure to NEXT-LOOKAHEADS_B. It structures the received stack in a tree format.

Having such tree, the procedure NEXT-LOOKAHEADS_B, presented in Figure 8, fast calculates the tokens that can be found given a stack and a transition. To achieve this, it uses the external functions READ₁ and FOLLOW₁, both defined in [DeRemer and Pennello 1982]. From a state p and a symbol X , READ₁ returns the tokens that can be read from GOTO₀(p, X) either directly or under nullable transitions; FOLLOW₁ returns the tokens either in READ₁(p, X) or in FOLLOW₁(p_0, C), as long as $C \rightarrow \alpha \bullet X\beta \in p$, $\beta \xrightarrow{*} \lambda$ and $p_0 \in \text{PRED}(p, \alpha)$ ⁶. Given a stack and a transition, NEXT-LOOKAHEADS_B grabs all tokens returned by READ₁, i.e., the tokens that can be read given the current context. Reductions are treated as in FOLLOW-SOURCES_B, except that in cases of un-

⁶ Originally, READ₁ and FOLLOW₁ were defined over nonterminal transitions - in NEXT-LOOKAHEADS_B, READ₁, as defined by Charles, was generalized to terminal and nonterminal transitions.

derflow, the simulation does not go further. Instead, the algorithm retrieves the desired tokens by calling FOLLOW₁. To reduce under non-underflow cases, the procedure pops states by calling UP; given a node *n* and a value *k*, UP returns the *k*-th ancestor of *n*. By using READ₁ and FOLLOW₁, which can be precomputed, NEXT-LOOKAHEADS_B does not have to search for source stacks when looking for lookaheads. The procedure's infinite loop control is achieved just like FOLLOW-SOURCES_B, i.e., by storing all visited transitions.

The presented algorithms do not have the limitation of stopping when dealing with NLALR grammars; termination is guaranteed to be reached under any circumstances.

5.2 Conflict Listing

When using LALR(*k_v*) action tables, the parser generator must not miscalculate the number of remaining conflicts. To illustrate this, consider the following table when *k_{max}* = 1:

	a	b	c	d	e	f	\$
8					S11, R8	R3, R8	R3
...							

Using *k_{max}* = 2, the table is given by:

	a	b	c	d	e	f	\$
8					L13	L14	R3
...							
13						S11, R8	
14					R3, R8	R3, R8	R3, R8

Simply examining the number of conflictual entries in the latter table allows identifying four conflicts, instead of the original two. When using *k_{max}* ≥ 2, SAIDE performs a depth first search from the lookahead action in an entry (*p*, *a*), and retrieves the actions in the entries that still contain conflicts. The obtained set of actions *acts* is a subset of the original set of conflicts. When performing the search, SAIDE also keeps track of all traversed edges of the DFA, and thus obtain the strings of length up to *k_{max}* for which the conflict is not removed. The number of conflicts for (*p*, *a*) is given by:

$$(|shift| \times |reds|) + (\lambda x. \text{if } x \geq 2 \text{ then } 1 \text{ else } 0) |reds|$$

where

$$\begin{aligned}
 shift &= \{s \mid s \in acts \wedge s \text{ is a shift action}\} \\
 reds &= \{r \mid r \in acts \wedge r \text{ is a reduce action}\}
 \end{aligned}$$

```

FOLLOW-SOURCESB(sources, transition, a, w, root, node, visited, roots)
1  stackSize ← HEIGHT(node, root)
2  if stackSize = 1
3    then if transition ∈ roots
4          then return
5          else roots ← roots ∪ {transition}
6    else if (ID(node), transition) ∈ visited
7          then return
8          else visited ← visited ∪ {(ID(node), transition)}
9  ASSERT(transition = (ts, X, q))
10 for Y ∈ V | GOTO0(q, Y) is defined
11 do if Y  $\xrightarrow{*}$  λ ∧ GET-FROM(node, (q, w)) = nil
12   then node2 ← NODE((q, w))
13        ADD-CHILD(node, node2)
14        FOLLOW-SOURCESB
15        (sources, (q, Y, GOTO0(q, Y)), a, w,
16        root, node2, visited, roots)
17   else if Y = a
18         then node2 ← node
19              list ← [FST(VALUE(node2))]
20              while (node2 ← PARENT(node2)) ≠ nil
21                do list ← list + [FST(VALUE(node2))]
22
23              ASSERT(list = [pn...p1])
24              stack ← [p1...pnq]
25              sources ← sources ∪ {(stack, wa)}
26 bottom ← FST(VALUE(root))
27 for C → γ• ∈ ts | C ≠ S
28 do if |γ| + 1 < stackSize
29   then node2 ← UP(node, |γ|)
30        SND(VALUE(node2)) ← w
31        ts2 ← FST(VALUE(node2))
32        FOLLOW-SOURCESB
33        (sources, (ts2, C, GOTO0(ts2, C)),
34        a, w, root, node2, visited, roots)
35   else ASSERT(γ = γ1γ2), where |γ2| = stackSize - 1
36        for p0 ∈ PRED(bottom, γ1)
37          do root2 ← NODE((p0, w))
38              FOLLOW-SOURCESB
39              (sources, (p0, C, GOTO0(p0, C)), a, w, root2,
40              root2, visited, roots)

```

Figure 5: Procedure FOLLOW-SOURCES_B.

For the conflicts that could not be automatically removed, SAIDE lists them in a heuristic manner. The heuristic here discussed builds a conflict graph, whose vertexes represent LALR states with at least one non-solved conflict. A directed edge connects p to q if there is a path from p to q in the LALR automaton, i.e., p propagates lookaheads to q . From this graph, a second one is built, formed by the SCCs of the first. In this graph, a directed edge between two vertexes exists

```

RESOLVE-CONFLICTS( $q, t, sources, n$ )
1  if  $t = \$ \vee n > k_{max}$ 
2    then return
3  allocate a new line  $p$  in the action table
4  for  $a \in \Sigma$ 
5    do  $Action[p, a] \leftarrow \emptyset$ 
6   $Action[q, t] \leftarrow \{Lp\}$ 
7  for act indexing sources
8    do for  $src \in sources[act]$ 
9      do ASSERT( $src = (stack, w)$ )
10      $la \leftarrow NEXT-LOOKAHEADS_A(stack, t)$ 
11     for  $a \in la$ 
12       do  $Action[p, a] \leftarrow Action[p, a] \cup \{act\}$ 
13 for  $a \in (\Sigma - \{\$\}) \mid |Action[p, a]| > 1$ 
14 do for  $act \in Action[p, a]$ 
15   do  $nSources \leftarrow \emptyset$ 
16     for  $src \in sources[act]$ 
17       do ASSERT( $src = (stack, w)$ )
18         FOLLOW-SOURCES $_A(nSources[act], stack, t, a, w)$ 
19
20   RESOLVE-CONFLICTS( $p, a, nSources, n + 1$ )

```

Figure 6: Procedure RESOLVE-CONFLICTS.

if at least one vertex in the first SCC connects to another vertex in the second SCC. The SCCs graph is then topologically sorted. From the obtained graph, the conflicts are listed according to the order of the SCCs, from left to right. Given an SCC c , all conflicts in state $p \in c$ are put on the listing.

5.3 Table Compression

An important problem using LALR(k) parsers is due to space requirements, for the size of the action table substantially grows when the value of k increases.

In LALR(k_v) parsers, the number of inspected lookaheads is minimized, but the size of action parsing tables are still considerable. We performed tests with nine LALR(1) compression schemes ⁷:

- the compression scheme proposed in [Aho et al. 1986] (ACS);

⁷ The authors of ACS and BCS didn't name their methods. For reference purposes, the adopted names reflect the corresponding bibliography entries in this work.

```

NEXT-LOOKAHEADSA(stack, t)
1  la ← ∅
2  ASSERT(stack = [p1...pn])
3  root ← node ← NODE(p1)
4  for 2 ≤ i ≤ n
5  do node2 ← NODE(pi)
6     ADD-CHILD(node, node2)
7     node ← node2
8  NEXT-LOOKAHEADSB(la, (pn, t, GOTO0(pn, t)), root, node, ∅)
9  return la

```

Figure 7: Façade procedure NEXT-LOOKAHEADS_A.

```

NEXT-LOOKAHEADSB(la, transition, root, node, visited)
1  ASSERT(transition = (ts, X, q))
2  bottom ← VALUE(root)
3  if (ID(node), transition) ∈ visited
4  then return
5  la ← la ∪ READ1(ts, X)
6  stackSize ← HEIGHT(node, root)
7  nStacks ← ∅
8  for C → γ • Xδ ∈ ts | δ* ⇒ λ ∧ C ≠ S
9  do if |γ| + 1 < stackSize
10     then node2 ← UP(node, |γ|)
11         nStacks ← nStacks ∪ {(node2, C)}
12
13     else ASSERT(γ = γ1γ2), where |γ2| = stackSize - 1
14         for p0 ∈ PRED(bottom, γ1)
15         do la ← la ∪ FOLLOW1(p0, C)
16 for (n, C) ∈ nStacks
17 do ts ← VALUE(n)
18     NEXT-LOOKAHEADSB(la, (ts, C, GOTO0(ts, C)), root, n, visited)

```

Figure 8: Procedure NEXT-LOOKAHEADS_B.

- the compression scheme proposed in [Bigonha et al. 1983] (BCS);
- row displacement scheme (RDS), reported in [Dencker et al. 1984];
- significance distance scheme (SDS), proposed in [Beach 1974];
- row column scheme (RCS), proposed in [Tewarson 1967];
- suppression of zeros scheme (SZS), reported in [Dencker et al. 1984];
- graph coloring scheme (GCS), proposed in [Schmitt 1979];
- line elimination scheme (LES), proposed in [Bell, 1974];
- row merging scheme (RMS), reported in [Passos 2007].

Each method was tested using the grammars of C, C#, HTML, Java and VB programming languages. The result of this study is directly generalized to LALR(k_v) parsers, due to the same layout of LALR(1) and LALR(k_v) parsing tables.

From the experiment, BCS and the combination of GCS, LES and RMS presented the highest compression rates, respectively 95% and 87%, in average. For BCS, there is the price of the overhead caused by the substitution of a direct access by a linear access to an interval of an array containing syntactic actions. In this scheme, the generated parser might execute unnecessary reductions, although correctness is preserved. The combination of GCS, LES and RMS preserves $O(1)$ access time and guarantees execution of the same number of actions as the non-compressed parser.

6 Experimental results

We performed tests in order to evaluate the automatic conflict mechanism and how the propose infinite loop control impacts in execution time. Table 1 shows the results obtained for the Algol-60, Scheme, Oberon-2 and Notus original grammars.

From the experiment, the number of conflicts was reduced in 51% and 97% in Scheme and Oberon-2 when using $k_{max} = 2$ in comparison with one token ahead. In Algol-60 the number of conflicts is not affected at any time; in Notus, there is a reduction in 6% when using at most two lookaheads. From $k_{max} \geq 3$, there's practically no overall change in the number of conflicts. The discrepancy between Scheme and Oberon-2 when compared to Notus and Algol-60 is that these are more human readable than closer to LALR conformance. The opposite situation occurs with Scheme and Oberon-2.

Grammar	$k_{max} = 1$		$k_{max} = 2$		$k_{max} = 3$		$k_{max} = 4$	
	Confs.	T. (ms)	Confs.	T. (ms)	Confs.	T. (ms)	Confs.	T. (ms)
Algol-60	61	670	61	2,480	61	4,837	61	18,017
Scheme	78	839	38	1,320	38	2,485	38	4,957
Oberon-2	32	898	1	1,464	1	1,567	1	1,493
Notus	575	1,068	541	34,828	539	38,921	539	67,467

Table 1: Evaluation of increasing k in automatic conflict removal.

Our tool was also used in building a parser for the Machina programming language [Bigonha et al. 2007]. The writing of its syntax specification was incrementally performed. An increment is the result of adding new rules to the previous increment when the set of conflicts of the latter becomes empty. In the experiment, $k_{max} = 2$ was used. At each reported conflict, the four phases for manual removal were applied. The application of the four phases defines a step. Figure 9 shows the dot graph for the number of conflicts obtained in each step. For reading purposes, the dots were connected to better identify the increasing and decreasing of conflicts. In the presented graph, the frontier between increments is marked by a dotted line. With exception of increment four, all steps inside other increments showed a decrease in the number of conflicts. The increase of conflicts occurs in the border of two increments, which is expected due to the adding of new rules. The produced parsing table was compressed with BCS and the combination of GCS, LES and RMS, resulting in 98% and 94% of compression rate.

7 Conclusion

This article presented an LALR parser generator supporting conflict resolution. Among the contributions of our work, we highlight the following:

- the process of conflict removal is eased by automatic conflict removal. In particular, the present algorithms remove some conflicts caused by lack of right context;
- for the cases in which manual removal is required, the tool assists users through a well defined methodology;
- modifications in Charles' proposal in order to accept any context free grammar. This makes the methodology independent of grammar characteristics;
- generalization of LALR(1) compression schemes to LALR(k_v) approach, turning it viable in the sense that memory requirements are minimized.

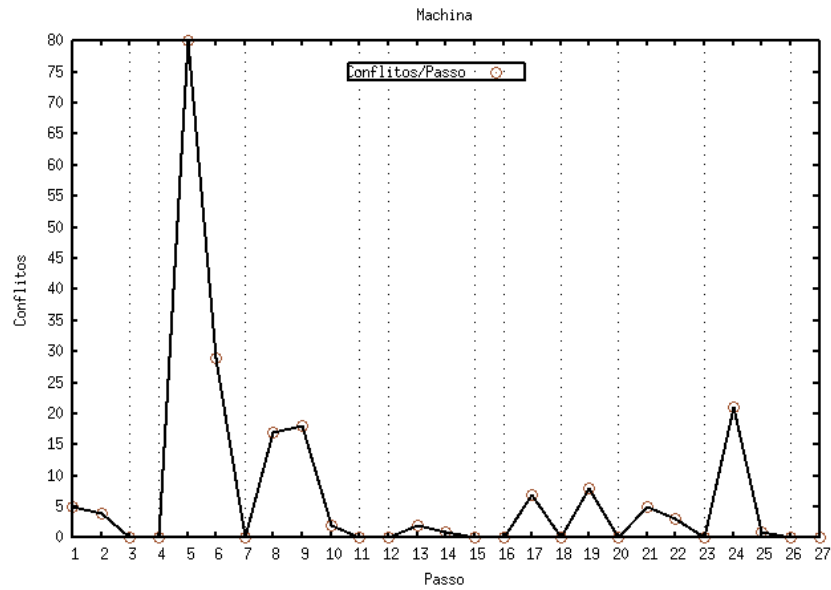


Figure 9: Conflict removal graph for the Machina programming language.

All obtained results are based on empirical data, determined by using established and also new programming languages, such as Notus and Machina. The presented results indicate that the application of the methodology contributes to the constant decrease on the number of conflicts in a grammar and that, in general, the time spent in calculating lookaheads in any context free grammar does not incur in a major overhead given today CPUs clocks.

References

- [Aho et al. 1986] Aho, A. V., Sethi, R., Ullman, J. D.: “Compilers, Principles, Techniques, and Tools”; Addison-Wesley, 1986.
- [Beach 1974] Beach, B.: “Storage Organization for a Type of Sparse Matrix”; Proc. 5th Southeastern Conference on Combinatorics, Graph Theory and Computing, 1974, 245-252.
- [Bell, 1974] Bell, J. R.: “A Compression Method for Compiler Precedence Tables.”; Proc. IFIP Congress, Elsevier North-Holland, New York (1974), 359–362.
- [Bigonha et al. 1983] Bigonha, R. S. and Bigonha, M. A. S.: “A Method for Efficient Compactation of LALR(1) Parsing Tables”; Technical Report, Federal University of Minas Gerais, Programming Languages Laboratory (1983).
- [Bigonha et al. 2007] Bigonha, R. S., Tirelo, F., Iorio, V. O., Bigonha, M. A. S.: “A Linguagem de Especificação Formal Machina 2.0”; Technical Report, Federal University of Minas Gerais, Programming Languages Laboratory (2007).
- [Charles 1991] Charles, P.: “A Practical Method for Constructing Efficient LALR(*k*) Parsers with Automatic Error Recovery”; Doctoral Thesis, New York University (1991).

- [Dencker et al. 1984] Dencker, P., Dürre, K., Heuft, H.: “Optimization of Parser Tables for Portable Compilers”; (TOPLAS) ACM Transactions on Programming Languages and Systems, 6, 4 (1984), 546-572.
- [DeRemer and Pennello 1982] DeRemer, F., Pennello, F.: “Efficient Computation of LALR(1) Look-Ahead Sets”; (TOPLAS) ACM Transactions on Programming Languages and Systems, 4, 4 (1982), 615-649.
- [Kristensen and Madsen 1981] Kristensen, B. B., Madsen, O. L.: “Methods for Computing LALR(k) Lookahead”; (TOPLAS) ACM Transactions on Programming Languages and Systems, 3, 1 (1981), 60-82.
- [Passos et al. 2007] Passos, L. T., Bigonha, Mariza A. S., Bigonha, R. S.: “A Methodology for Removing LALR(k) Conflicts”; J.UCS (Journal for Universal Computer Science), 13, 6 (2007), 737-752.
- [Passos 2007] Passos, L. T.: “Gerador LALR com Suporte a Resolução de Conflitos”; Master’s Thesis, Federal University of Minas Gerais, Programming Languages Laboratory (2007).
- [Schmitt 1979] Schmitt, A.: “Minimizing Storage Space of Sparse Matrices by Graph Coloring Algorithms”; Graphs, Data Structures, Algorithms, 1979, 157-168.
- [Tewarson 1967] Tewarson, R. P.: “Row Column Permutation of Sparse Matrices”; Computer Journal; 10, 3 (1967), 300-305.
- [Tirelo and Bigonha 2006] Tirelo, F., Bigonha, R.: “Notus”; Technical Report, Federal University of Minas Gerais, Programming Languages Laboratory (2006).
- [Tomita 1991] Tomita, M.: “Generalized L.R. Parsing”; Kluwer Academic Publishers, Norwell (1991).