

بسم الله الرحمن الرحيم

دانشگاه صنعتی اصفهان – دانشکده مهندسی برق و کامپیوتر
(نیم سال تحصیلی ۴۰۲۲)

کامپایلر

حسین فلسفین

Constructing LALR Parsing Tables

We now introduce our last parser construction method, the LALR (lookahead-LR) technique. This method is often used in practice, because the tables obtained by it **are considerably smaller than the canonical LR tables**, yet most common syntactic constructs of programming languages **can be expressed conveniently** by an LALR grammar. The same is almost true for SLR grammars, but there are a few constructs that cannot be conveniently handled by SLR techniques (see Example 4.48, for example).

For a comparison of parser size, the SLR and LALR tables for a grammar always have the same number of states, and this number is typically **several hundred** states for a language like C. The canonical LR table would typically have **several thousand** states for the same-size language. Thus, it is **much easier and more economical** to construct SLR and LALR tables than the canonical LR tables.

نکاتی درباره چیهستی و چراییی تجزیه LALR(1)

👉 LR(1) *parsers would be impractical in that the space required for their deterministic automata would be prohibitive. A modest grammar might already require hundreds of thousands or even millions of states.*

LR(1) *parsing tables can be very large, with many states. A smaller table can be made by merging any two states whose items are identical except for lookahead sets. The result parser is called an LALR(1) parser, for lookahead LR(1). The surprising thing is that this procedure preserves almost all the original lookahead power and still saves an enormous amount of memory.*

👉 LALR(1) parsers are powerful, almost as powerful as LR(1) parsers, they have fairly modest memory requirements, only slightly inferior to (= larger than) those of LR(0) parsers, and they are time-efficient. LALR(1) parsing may very well be **the most-used parsing method in the world today**. Probably the most famous LALR(1) parser generators are **yacc** and its GNU version **bison**.

👉 Even though LALR(1) grammars are not as powerful as LR(1) grammars, they are sufficiently powerful to describe most programming languages. This, together with their small (relative to LR) table size, makes the LALR(1) family of grammars an excellent candidate for the automatic generation of parsers. Stephen C. Johnson's **YACC**, for “Yet Another Compiler-Compiler”, based on LALR(1) techniques, was probably the first practical bottom-up parser generator. GNU has developed an open-source version called **Bison**.

👉 LALR(1) grammars make for parsers that are almost as powerful as LR(1) grammars but result in **much more space-efficient parsing tables**. This goes some way in explaining the popularity of parser generators such as **YACC and Bison**.

👉 An LR(1) parsing table for a typical programming language such as Java can have **thousands of states**, and so thousands of rows. One could argue that, given the inexpensive memory nowadays, this is not a problem. On the other hand, smaller programs and data make for faster running programs so it would be advantageous if we might be able to reduce the number of states. LALR(1) is a parsing method that does just this.

Example 4.60: Again consider grammar (4.55) whose GOTO graph was shown in Fig. 4.41. As we mentioned, there are three pairs of sets of items that can be merged. I_3 and I_6 are replaced by their union:

$$\begin{aligned} I_{36}: \quad & C \rightarrow c \cdot C, \ c/d/\$ \\ & C \rightarrow \cdot cC, \ c/d/\$ \\ & C \rightarrow \cdot d, \ c/d/\$ \end{aligned}$$

I_4 and I_7 are replaced by their union:

$$I_{47}: \quad C \rightarrow d \cdot, \ c/d/\$$$

and I_8 and I_9 are replaced by their union:

$$I_{89}: \quad C \rightarrow cC \cdot, \ c/d/\$$$

The LALR action and goto functions for the condensed sets of items are shown in Fig. 4.43.

STATE	ACTION			GOTO	
	<i>c</i>	<i>d</i>	\$	<i>S</i>	<i>C</i>
0	s36	s47		1	2
1			acc		
2	s36	s47			5
36	s36	s47			89
47	r3	r3	r3		
5			r1		
89	r2	r2	r2		

Figure 4.43: LALR parsing table for the grammar of Example 4.54

To see how the GOTO's are computed, consider $\text{GOTO}(I_{36}, C)$. In the original set of LR(1) items, $\text{GOTO}(I_3, C) = I_8$, and I_8 is now part of I_{89} , so we make $\text{GOTO}(I_{36}, C)$ be I_{89} . We could have arrived at the same conclusion if we considered I_6 , the other part of I_{36} . That is, $\text{GOTO}(I_6, C) = I_9$, and I_9 is now part of I_{89} . For another example, consider $\text{GOTO}(I_2, c)$, an entry that is exercised after the shift action of I_2 on input c . In the original sets of LR(1) items, $\text{GOTO}(I_2, c) = I_6$. Since I_6 is now part of I_{36} , $\text{GOTO}(I_2, c)$ becomes I_{36} . Thus, the entry in Fig. 4.43 for state 2 and input c is made s36, meaning shift and push state 36 onto the stack. \square

ساخت جدول LALR(1)

Algorithm 4.59: An easy, but space-consuming LALR table construction.

INPUT: An augmented grammar G' .

OUTPUT: The LALR parsing-table functions ACTION and GOTO for G' .

METHOD:

1. Construct $C = \{I_0, I_1, \dots, I_n\}$, the collection of sets of LR(1) items.
2. For each **core** present among the set of LR(1) items, find all sets having that **core**, and replace these sets by their union.
3. Let $C' = \{J_0, J_1, \dots, J_m\}$ be the resulting sets of LR(1) items. The parsing actions for state i are constructed from J_i in the same manner as in Algorithm 4.56. If there is a parsing action conflict, the algorithm fails to produce a parser, and the grammar is said not to be LALR(1).
4. The GOTO table is constructed as follows. If J is the union of one or more sets of LR(1) items, that is, $J = I_1 \cup I_2 \cup \dots \cup I_k$, then the cores of $\text{GOTO}(I_1, X)$, $\text{GOTO}(I_2, X)$, \dots , $\text{GOTO}(I_k, X)$ are the same, since I_1, I_2, \dots, I_k all have the same core. Let K be the union of all sets of items having the same core as $\text{GOTO}(I_1, X)$. Then $\text{GOTO}(J, X) = K$.

There are algorithms for directly constructing the tables for LALR(1) parsing without first constructing the tables for LR(1) parsing, but we do not consider them here.

With LALR (lookahead LR) parsing, we attempt to reduce the number of states in an LR(1) parser by merging similar states. This reduces the number of states to the same as SLR(1), but still retains *some* of the power of the LR(1) lookaheads.

$S' \rightarrow S$
 $S \rightarrow XX$
 $X \rightarrow aX$
 $X \rightarrow b$

$I_0:$ $S' \rightarrow \bullet S, \$$
 $S \rightarrow \bullet XX, \$$
 $X \rightarrow \bullet aX, a/b$
 $X \rightarrow \bullet b, a/b$

$I_1:$ $S' \rightarrow S \bullet, \$$

$I_2:$ $S \rightarrow X \bullet X, \$$
 $X \rightarrow \bullet aX, \$$
 $X \rightarrow \bullet b, \$$

$I_3:$ $X \rightarrow a \bullet X, a/b$
 $X \rightarrow \bullet aX, a/b$
 $X \rightarrow \bullet b, a/b$

$I_4:$ $X \rightarrow b \bullet, a/b$

$I_5:$ $S \rightarrow XX \bullet, \$$

$I_6:$ $X \rightarrow a \bullet X, \$$
 $X \rightarrow \bullet aX, \$$
 $X \rightarrow \bullet b, \$$

$I_7:$ $X \rightarrow b \bullet, \$$

$I_8:$ $X \rightarrow aX \bullet, a/b$

$I_9:$ $X \rightarrow aX \bullet, \$$

I₃₆: $X \rightarrow a \bullet X, a/b/\$$
 $X \rightarrow \bullet aX, a/b/\$$
 $X \rightarrow \bullet b, a/b/\$$

I₄₇: $X \rightarrow b \bullet, a/b/\$$

I₈₉: $X \rightarrow aX \bullet, a/b/\$$

But isn't this just SLR(1) all over again? In the above example, yes, since after the merging we coincidentally end up with the complete follow sets as the lookahead.

This is not always the case however. Consider this example:

$I_0: \quad S' \rightarrow \bullet S, \$$

$S \rightarrow \bullet Bbb, \$$

$S' \rightarrow S$
 $S \rightarrow Bbb \mid aab \mid bBa$
 $B \rightarrow a$

$S \rightarrow \bullet aab, \$$

$S \rightarrow \bullet bBa, \$$

$B \rightarrow \bullet a, b$

$I_2: \quad S \rightarrow B \bullet bb, \$$

$I_3: \quad S \rightarrow a \bullet ab, \$$

$B \rightarrow a \bullet, b$

....

$I_1: \quad S' \rightarrow S \bullet, \$$

چرا اینجا پارسر LALR(1) از پارسر SLR(1) متمایز خواهد بود؟

In an SLR(1) parser there is a shift-reduce conflict in state 3 when the next input is anything in FOLLOW(B) which includes a and b . In LALR(1), state 3 will shift on a and reduce on b . Intuitively, this is because the LALR(1) state "remembers" that we arrived at state 3 after seeing an a . Thus we are trying to parse either Bbb or aab . In order for that first a to be a valid reduction to B , the next input has to be exactly b since that is the only symbol that can follow B in this particular context. Although elsewhere an expansion of B can be followed by an a , we consider only the subset of the follow set that can appear here, and thus avoid the conflict an SLR(1) parser would have.

یک پارسر LALR(1) برای گرامری که اتوماتون LR(1) آن ۲۲ استیت داشت

$$0. E' ::= E$$

$$1. E ::= E + T$$

$$2. E ::= T$$

$$3. T ::= T * F$$

$$4. T ::= F$$

$$5. F ::= (E)$$

$$6. F ::= \text{id}$$

یک پارسر LALR(1) برای گرامری که اتوماتون LR(1) آن ۲۲ استیت داشت

$s_0 =$

$\{[E' ::= \cdot E, \#],$	$\text{goto}(s_0, E) = s_1$
$[E ::= \cdot E + T, +/\#],$	$\text{goto}(s_0, T) = s_{2.9}$
$[E ::= \cdot T, +/\#],$	$\text{goto}(s_0, F) = s_{3.10}$
$[T ::= \cdot T * F, +/*/\#],$	$\text{goto}(s_0, () = s_{4.11}$
$[T ::= \cdot F, +/*/\#],$	$\text{goto}(s_0, \text{id}) = s_{5.12}$
$[F ::= \cdot (E), +/*/\#],$	
$[F ::= \cdot \text{id}, +/*/\#]$	

$s_1 =$

$\{[E' ::= E \cdot, \#],$	$\text{goto}(s_1, +) = s_{6.16}$
$[E ::= E \cdot + T, +/\#]$	

$s_{6.16} =$

$\{[E ::= E + \cdot T, +/\#],$	$\text{goto}(s_{6.16}, T) = s_{13.19}$
$[T ::= \cdot T * F, +/*/\#],$	$\text{goto}(s_{6.16}, F) = s_{3.10}$
$[T ::= \cdot F, +/*/\#],$	$\text{goto}(s_{6.16}, () = s_{4.11}$
$[F ::= \cdot (E), +/*/\#],$	$\text{goto}(s_{6.16}, \text{id}) = s_{5.12}$
$[F ::= \cdot \text{id}, +/*/\#]$	

$s_{7.17} =$

$\{[T ::= T * \cdot F, +/*/\#],$	$\text{goto}(s_{7.17}, F) = s_{14.20}$
$[F ::= \cdot (E), +/*/\#],$	$\text{goto}(s_{7.17}, () = s_{4.11}$
$[F ::= \cdot \text{id}, +/*/\#]$	$\text{goto}(s_{7.17}, \text{id}) = s_{5.12}$

یک پارسر LALR(1) برای گرامری که اتوماتون LR(1) آن ۲۲ استیت داشت

$s_{2.9} =$

$\{[E ::= T \cdot, +/)/\#], \text{ goto}(s_{2.9}, *) = s_{7.17}$
 $[T ::= T \cdot * F, +/)/\#]\}$

$s_{3.10} =$

$\{[T ::= F \cdot, +/)/\#]\}$

$s_{4.11} =$

$\{[F ::= (\cdot E), +/)/\#], \text{ goto}(s_{4.11}, E) = s_{8.18}$
 $[E ::= \cdot E + T, +/)], \text{ goto}(s_{4.11}, T) = s_{2.9}$
 $[E ::= \cdot T, +/)], \text{ goto}(s_{4.11}, F) = s_{3.10}$
 $[T ::= \cdot T * F, +/)/\#], \text{ goto}(s_{4.11}, () = s_{4.11}$
 $[T ::= \cdot F, +/)/\#], \text{ goto}(s_{4.11}, \text{id}) = s_{5.12}$
 $[F ::= \cdot (E), +/)/\#],$
 $[F ::= \cdot \text{id}, +/)/\#]\}$

$s_{5.12} = \{[F ::= \text{id} \cdot, +/)/\#]\}$

$s_{8.18} =$

$\{[F ::= (E \cdot), +/)/\#], \text{ goto}(s_{8.18},) = s_{15.21}$
 $[E ::= E \cdot + T, +/)]\} \text{ goto}(s_{8.18}, +) = s_{6.16}$

$s_{13.19} =$

$\{[E ::= E + T \cdot, +/)/\#], \text{ goto}(s_{13.19}, *) = s_{7.17}$
 $[T ::= T \cdot * F, +/)/\#]\}$

$s_{14.20} = \{[T ::= T * F \cdot, +/)/\#]\}$

$s_{15.21} = \{[F ::= (E) \cdot, +/)/\#]\}$

جدول LALR(1) بجای ۲۲ سطر، ۱۲ سطر دارد

	Action						Goto		
	+	*	()	id	#	E	T	F
0			s4		s5		1	2	3
1	s6					accept			
2.9	r2	s7.17		r2		r2			
3.10	r4	r4		r4		r4			
4.11			s11		s12		8.18	2.9	3.10
5.12	r6	r6		r6		r6			
6.16			s4		s5			13.19	3.10
7.17			s4		s5				14.20
8.18	s16			s15					
13.19	r1	s7				r1			
14.20	r3	r3				r3			
15.21	r5	r5				r5			

درباره ماهیت کانفلیکتهای در تجزیه LALR(1)

Suppose we have an LR(1) grammar, that is, one whose sets of LR(1) items produce no parsing-action conflicts. If we replace all states having the same core with their union, it is possible that the resulting union will have a conflict, **but it is unlikely** for the following reason: Suppose in the union there is a conflict on lookahead a because there is an item $[A \rightarrow \alpha \bullet, a]$ calling for a reduction by $A \rightarrow \alpha$, and there is another item $[B \rightarrow \beta \bullet a \gamma, b]$ calling for a shift. Then some set of items from which the union was formed has item $[A \rightarrow \alpha \bullet, a]$, and since the cores of all these states are the same, it must have an item $[B \rightarrow \beta \bullet a \gamma, c]$ for some c . But then this state has the same shift/reduce conflict on a , and the grammar was not LR(1) as we assumed. Thus, the merging of states with common cores **can never produce a shift/reduce conflict** that was not present in one of the original states, because shift actions depend only on the core, not the lookahead.

A shiftreduce conflict cannot exist in a merged set unless the conflict was already present in one of the original LR(1) configuring sets. When merging, the two sets must have the same core items. If the merged set has a configuration that shifts on a and another that reduces on a, both configurations must have been present in the original sets, and at least one of those sets had a conflict already.

It is possible, however, that a merger will produce a reduce/reduce conflict, as the following example shows.

Example 4.58: Consider the grammar

$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow a A d \mid b B d \mid a B e \mid b A e \\ A &\rightarrow c \\ B &\rightarrow c \end{aligned}$$

which generates the four strings acd , ace , bcd , and bce . The reader can check that the grammar is LR(1) by constructing the sets of items. Upon doing so, we find the set of items $\{[A \rightarrow c\cdot, d], [B \rightarrow c\cdot, e]\}$ valid for viable prefix ac and $\{[A \rightarrow c\cdot, e], [B \rightarrow c\cdot, d]\}$ valid for bc . Neither of these sets has a conflict, and their cores are the same. However, their union, which is

$$\begin{aligned} A &\rightarrow c\cdot, d/e \\ B &\rightarrow c\cdot, d/e \end{aligned}$$

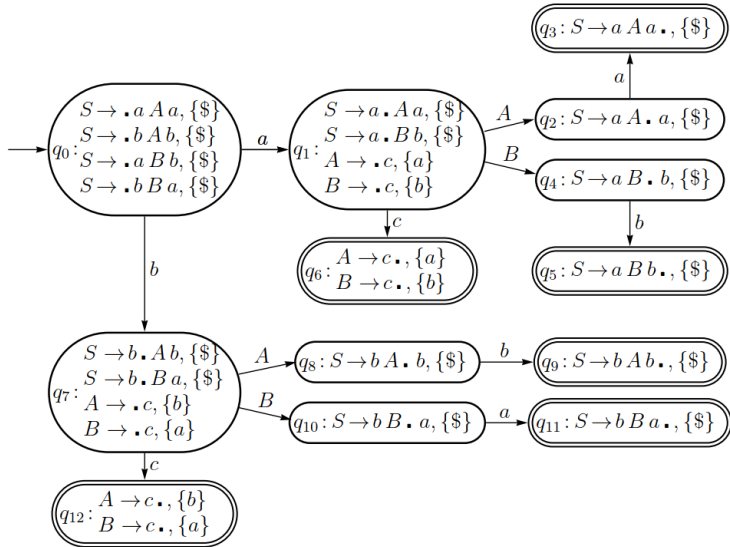
generates a reduce/reduce conflict, since reductions by both $A \rightarrow c$ and $B \rightarrow c$ are called for on inputs d and e . \square

Example (a non-LALR(1) grammar): Let us consider the grammar G with axiom S and the following productions:

1. $S \rightarrow a A a$
2. $S \rightarrow b A b$
3. $S \rightarrow a B b$
4. $S \rightarrow b B a$
5. $A \rightarrow c$
6. $B \rightarrow c$

Since the axiom S does not occur on the right hand side of any production, we can take the given grammar G to be the augmented grammar of the grammar G itself.

The finite automaton for the LR(1) parsing of the grammar G



LR(1) *Parsing table for the grammar G*

	<u>action</u>				<u>goto</u>	
	<i>a</i>	<i>b</i>	<i>c</i>	<i>\$</i>	<i>A</i>	<i>B</i>
<i>q</i> ₀	<i>sh q</i> ₁	<i>sh q</i> ₇				
<i>q</i> ₁			<i>sh q</i> ₆		<i>q</i> ₂	<i>q</i> ₄
<i>q</i> ₂	<i>sh q</i> ₃					
<i>q</i> ₃				<i>red 1</i>		
<i>q</i> ₄		<i>sh q</i> ₅				
<i>q</i> ₅				<i>red 3</i>		
<i>q</i> ₆	<i>red 5</i>	<i>red 6</i>				
<i>q</i> ₇			<i>sh q</i> ₁₂		<i>q</i> ₈	<i>q</i> ₁₀
<i>q</i> ₈		<i>sh q</i> ₉				
<i>q</i> ₉				<i>red 2</i>		
<i>q</i> ₁₀	<i>sh q</i> ₁₁					
<i>q</i> ₁₁				<i>red 4</i>		
<i>q</i> ₁₂	<i>red 6</i>	<i>red 5</i>				

- where:
1. $S \rightarrow a A a$
 2. $S \rightarrow b A b$
 3. $S \rightarrow a B b$
 4. $S \rightarrow b B a$
 5. $A \rightarrow c$
 6. $B \rightarrow c$

We fuse the states q_6 and q_{12} , thereby getting

$$q_{612}: \begin{array}{l} A \rightarrow c \cdot, \{a, b\} \\ B \rightarrow c \cdot, \{a, b\} \end{array}$$

the state: . Thus, in the LALR(1) parsing table we **will get two reduce-reduce conflicts**, because in the columns a and b of that table we have both the *red 5* action (for the production $A \rightarrow c$) and also the *red 6* action (for the production $B \rightarrow c$). **Hence, the given grammar G is not an LALR(1) grammar.**

جمع‌بندی و نکات نهایی بحث تجزیه پائین به بالا

*Theoretically, LR(1) grammars are the largest category of grammars that can be parsed deterministically while looking ahead just one token. Of course, LR(k) grammars for $k > 1$ are even more powerful, but one must look ahead k tokens, more importantly, the parsing tables must (in principle) keep track of all possible token strings of length k . So, in principle, the tables can grow **exponentially with k** .*

FACT [Relationships Between $LALR(k)$ Grammars and $LR(k)$ Grammars] The class of $LALR(0)$ grammars coincides with the class of $LR(0)$ grammars. For $k \geq 1$ the class of $LALR(k)$ grammars is properly contained in the class of $LR(k)$ grammars.

EXAMPLE The grammar with axiom S and productions:

$$S \rightarrow aAa \mid bAb \mid aBb \mid bBa$$

$$A \rightarrow c$$

$$B \rightarrow c$$

is an $LR(1)$ grammar which is *not* $LALR(k)$ for any $k \geq 0$.

FACT [Hierarchy of the $LR(k)$ Grammars] For all $k \geq 0$, the grammar with axiom S and the productions:

$$S \rightarrow ab^k c \mid Ab^k d$$

$$A \rightarrow a$$

is an $LR(k+1)$ grammar which is *not* an $LR(k)$ grammar.

A hierarchy of grammar classes

