

## Notes on SLR and LR parsing

*Handout written by Maggie Johnson and revised by me.*

### The problem with LR(0)

LR(0) is the simplest technique in the LR family, and although that makes it the easiest to learn, these parsers are too weak to be of practical use for anything but a very limited set of grammars. The examples given at the end of the LR(0) handout show how even small changes can introduce conflicts that make the grammar not LR(0). The fundamental limitation of LR(0) is the 0. It is a stifling constraint to have to make decisions based on only what has already been read, without even glancing at what comes next in the input. If we could peek at the next token coming up and use that as part of the decision-making, we will find that it allows for a much larger class of grammars to be parsed.

### SLR(1)

We will first consider SLR(1) where the S is for Simple. SLR(1) parser uses the same LR(0) configuring sets and has the same table structure and parser operation, so everything you've already learned about LR(0) applies here. The difference comes in assigning table actions, where we are going to use a token of lookahead to help arbitrate among the conflicts. If we think back to the kind of conflicts we encountered in LR(0) parsing, it was the reduce actions that seemed to cause us grief. An LR(0) grammar can have at most one reduce action in a particular state and cannot have both shift and reduce actions in the same state. Basically a reduce action in an LR(0) parser is a trump card, it beats everything. But is that always appropriate? If we were to look at the next upcoming token, it may tell us something that invalidates that reduction. If the string on top of the stack could be reduced to the non-terminal A, what tokens do we expect to find as the next input? What tokens would tell us that the reduction is not appropriate? Doesn't it seem like Follow(A) could be useful here?

The simple improvement that SLR(1) makes on the basic LR(0) parser is to only perform a reduction action if the next input token is a member of the follow set of the non-terminal being reduced. When filling in the table, we don't assume a reduce on all inputs as we did in LR(0), we selectively choose the reduction only when the next input symbols in a member of the follow set. To be more precise, here is the algorithm for SLR(1) table construction (note all steps are the same as for LR(0) table construction except for 2a)

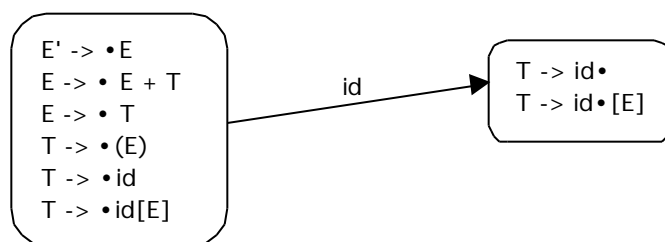
1. Construct  $F = \{I_0, I_1, \dots, I_n\}$ , the collection of LR(0) configuring sets for  $G'$ .
2. State  $i$  is determined from  $I_i$ . The parsing actions for the state are determined as follows:
  - a) If  $A \rightarrow \underline{u}\bullet$  is in  $I_i$  then set  $\text{Action}[i, a]$  to reduce  $A \rightarrow \underline{u}$  for all  $a$  in  $\text{Follow}(A)$  ( $A$  may not be  $S'$ ).
  - b) If  $S' \rightarrow S\bullet$  is in  $I_i$  then set  $\text{Action}[i, \$]$  to accept.
  - c) If  $A \rightarrow \underline{u}\bullet a\bar{v}$  is in  $I_i$  and  $\text{succ}(I_i, a) = I_j$ , then set  $\text{Action}[i, a]$  to shift  $j$  ( $a$  must be a terminal).
3. The goto transitions for state  $i$  are constructed for all non-terminals  $A$  using the rule: If  $\text{succ}(I_i, A) = I_j$ , then  $\text{Goto}[i, A] = j$ .
4. All entries not defined by rules 2 and 3 are errors.
5. The initial state is the one constructed from the configuring set containing  $S' \rightarrow \bullet S$ .

In SLR(1) parser, it is allowable for there to be both shift and reduce items in the same state as well as multiple reduce items. The SLR(1) parser will be able to determine what action to take as long as the necessary follow sets are disjoint.

Let's consider those changes at the end of the LR(0) handout to the simplified expression grammar that would have made it no longer LR(0). Here is the version with the addition of array access:

$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E + T \mid T \\ T &\rightarrow (E) \mid id \mid id[E] \end{aligned}$$

Here are the first two LR(0) configurating sets traced out after first seeing id in the input.

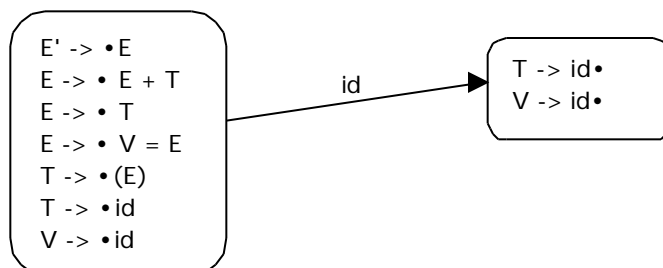


In an LR(0) parser, the set on the right has a shift-reduce conflict. However, an SLR(1) will compute  $\text{Follow}(T) = \{ + \} ] \$ \}$  and only enter the reduce action on those tokens. The input  $[$  will shift and there is no conflict. Thus this grammar is SLR(1).

Similarly, the simplified expression grammar with the assignment addition:

$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E + T \mid T \mid V = E \\ T &\rightarrow (E) \mid id \\ V &\rightarrow id \end{aligned}$$

Here are the first two LR(0) configurating sets traced out after first seeing id in the input.



In an LR(0) parser, the set on the right has a reduce-reduce conflict. However, an SLR(1) parser will compute  $\text{Follow}(T) = \{ + \} \$ \}$  and  $\text{Follow}(V) = \{ = \}$  and thus can distinguish which reduction to apply depending on the next input token. The modified grammar is SLR(1).

### SLR(1) grammars

A grammar is SLR(1) if the following two conditions hold for each configurating set:

1. For any item  $A \rightarrow \underline{u} \bullet x \underline{v}$  in the set, with terminal  $x$ , there is no complete item  $B \rightarrow \underline{w} \bullet$  in that set with  $x$  in  $\text{Follow}(B)$ . In the tables, this translates no shift-reduce conflict on any state. This

means the successor function for  $x$  from that set either shifts to a new state or reduces, but not both.

2. For any two complete items  $A \rightarrow \underline{u} \bullet$  and  $B \rightarrow \underline{v} \bullet$  in the same configuring set, the follow sets must be disjoint, e.g.  $\text{Follow}(A) \cap \text{Follow}(B)$  is empty. This translates to no reduce-reduce conflict on any state. If more than one non-terminal could be reduced from this set, it must be possible to uniquely determine which using only one token of lookahead.

All LR(0) grammars are SLR(1) but the reverse is not true, as the two extensions to our expression grammar demonstrated. The addition of just one token of lookahead and use of the follow set greatly expands the class of grammars that can be parsed without conflict.

### The limitations of SLR(1)

The SLR(1) technique still leaves something to be desired, because we are not using all the information that we have at our disposal. When we have a completed configuration (i.e. dot at the end) such as  $X \rightarrow \underline{u} \bullet$ , we know this corresponds to a situation in which we have  $\underline{u}$  as a handle on top of the stack which we then can reduce, i.e., replacing  $\underline{u}$  by  $X$ . We allow such a reduction whenever the next symbol is in  $\text{Follow}(X)$ . However, it may be that we should not reduce for every symbol in  $\text{Follow}(X)$ , because the symbols below  $\underline{u}$  on the stack preclude  $\underline{u}$  being a handle for reduction in this case. In other words, SLR(1) states only tell us about the sequence on top of the stack, not what is below it on the stack. We may need to divide an SLR(1) state into separate states to differentiate the possible means by which that sequence has appeared on the stack. By carrying more information in the state, it will allow us to rule out these invalid reductions.

Consider this example from the Aho/Sethi/Ullman that defines a small grammar for assignment statements, using the non-terminal  $L$  for l-value and  $R$  for r-value and  $*$  for contents-of.

$S' \rightarrow S$   
 $S \rightarrow L = R$   
 $S \rightarrow R$   
 $L \rightarrow *R$   
 $L \rightarrow \text{id}$   
 $R \rightarrow L$

$I_0:$      $S' \rightarrow \bullet S$   
            $S \rightarrow \bullet L = R$   
            $S \rightarrow \bullet R$   
            $L \rightarrow \bullet *R$   
            $L \rightarrow \bullet \text{id}$   
            $R \rightarrow \bullet L$

$I_5:$      $L \rightarrow \text{id} \bullet$   
  
 $I_6:$      $S \rightarrow L = \bullet R$   
            $R \rightarrow \bullet L$   
            $L \rightarrow \bullet *R$   
            $L \rightarrow \bullet \text{id}$

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

$I_7:$      $L \rightarrow *R \bullet$

$I_2:$      $S \rightarrow L \bullet = R$   
            $R \rightarrow L \bullet$

$I_8:$      $R \rightarrow L \bullet$

$I_9:$      $S \rightarrow L = R \bullet$

$I_3:$      $S \rightarrow R \bullet$

$I_4:$      $L \rightarrow * \bullet R$   
            $R \rightarrow \bullet L$   
            $L \rightarrow \bullet *R$   
            $L \rightarrow \bullet \text{id}$

Consider parsing the expression  $id = id$ . After working our way to configuring set  $I_2$  having reduced the first  $id$  to  $L$ , we have a choice upon seeing  $=$  coming up in the input. The first item in the set wants to set  $Action[2,=]$  be shift 6, which corresponds to moving on to find the rest of the assignment. However,  $=$  is also in  $Follow(R)$  because  $S \Rightarrow L=R \Rightarrow *R = R$ . Thus, the second configuration wants to reduce in that slot  $R \rightarrow L$ . This is a shift-reduce conflict but not because of any problem with the grammar. A SLR parser does not remember enough left context to decide what should happen when it encounters a  $=$  in the input having seen a string reducible to  $L$ . Although the sequence on top of the stack could be reduced to  $R$ , but we don't want to choose this reduction because there is no possible right sentential form that begins  $R = \dots$  (there is one beginning  $*R = \dots$  which is not the same). Thus, the correct choice is to shift.

It's not further lookahead that the SLR tables are missing, we don't need to see additional symbols coming up in the input, we have already seen the information that allows us to determine the correct choice. What we need is to retain a little more of the left context that brought us here. In this example grammar, the only time we should consider reducing by production  $R \rightarrow L$  is during a derivation that has already seen a  $*$  or an  $=$ . Just using the entire follow set is not discriminating enough as the guide for when to reduce. The follow set contains symbols that can follow  $R$  in any valid sentence but it does not precisely indicate which symbols follow  $R$  within a particular derivation. So we will augment our states to include information about exactly what portion of the follow set is appropriate given the path we have taken to that state.

We can be in state 2 for one of two reasons, we are trying to build a sentence from  $S \rightarrow L = R$  or from  $S \rightarrow R \rightarrow L$ . If the upcoming symbol is  $=$ , then that rules out the second choice and we must be building the first, which tells us to shift. The reduction should only be applied if the next input symbol is  $\$$ . Even though  $=$  is  $Follow(R)$  because of the other contexts that an  $R$  can appear, in this particular situation, it is not appropriate because when deriving a sentence  $S \rightarrow R \rightarrow L$ ,  $=$  cannot follow  $R$ .

### Constructing LR(1) parsing tables

LR or *canonical LR* parsing incorporates the required extra information into the state by redefining configurations to include a terminal symbol as an added component. LR(1) configurations have the general form:

$$A \rightarrow X_1 \dots X_i \bullet X_{i+1} \dots X_j, a$$

This means we have states corresponding to  $X_1 \dots X_i$  on the stack and we are looking to put states corresponding to  $X_{i+1} \dots X_j$  on the stack and then reduce, but only if the token following  $X_j$  is the terminal  $a$ .  $a$  is called the *lookahead* of the configuration. The lookahead only comes into play with LR(1) configurations with a dot at the right end:

$$A \rightarrow X_1 \dots X_j \bullet, a$$

This means we have states corresponding to  $X_1 \dots X_j$  on the stack but we may only reduce when the next symbol is  $a$ . The symbol  $a$  is either a terminal or  $\$$  (end of input marker). With SLR(1) parsing, we would reduce if the next token was any of those in  $Follow(A)$ . With LR(1) parsing, we reduce only if the next token is exactly  $a$ . We may have more than one symbol in the lookahead for the configuration, as a convenience, we list those symbols separated by a forward slash. Thus, the configuration  $A \rightarrow \underline{u} \bullet, a/b/c$  says that it is valid to reduce  $\underline{u}$  to  $A$  only if the next token is equal to  $a$ ,  $b$ , or  $c$ . The configuration lookahead will always be a subset of  $Follow(A)$ , but it could be a proper subset.

Recall the definition of a *viable prefix* from the SLR handout. Viable prefixes are those prefixes of right sentential forms that can appear on the stack of a shift-reduce parser. Formally we say that a configuration  $[A \rightarrow \underline{u} \bullet \underline{v}, a]$  is valid for a viable prefix if there is a rightmost derivation  $S \Rightarrow^* A \underline{w} \Rightarrow^* \underline{u} \underline{v} \underline{w}$  where  $\underline{u} = \underline{u}$  and either  $a$  is the first symbol of  $\underline{w}$  or  $\underline{w}$  is  $\$$  and  $a$  is  $\$$ . For example,

$$\begin{aligned} S &\rightarrow ZZ \\ Z &\rightarrow xZ \mid y \end{aligned}$$

There is a rightmost derivation  $S \Rightarrow^* xxZxy \Rightarrow^* xxxZxy$ . We see that configuration  $[Z \rightarrow \underline{x} \bullet Z, x]$  is valid for viable prefix  $\underline{u} = xxx$  by letting  $\underline{u} = xx$ ,  $A = Z$ ,  $\underline{w} = xy$ ,  $\underline{u} = x$  and  $\underline{v} = Z$ . Another example is from the rightmost derivation  $S \Rightarrow^* ZxZ \Rightarrow^* ZxxZ$ , making  $[Z \rightarrow \underline{x} \bullet Z, \$]$  valid for viable prefix  $\underline{u} = Zxx$ .

Often we have a number of LR(1) configurations that differ only in their lookahead components. The addition of a lookahead component to LR(1) configurations allows us to make parsing decisions beyond the capability of SLR(1) parsers. There is, however, a big price to be paid. There will be more distinct configurations and thus many more possible configuring sets. This increases the size of the goto and action tables considerably. In the past when memory was smaller, it was difficult to find storage-efficient ways of representing these tables, but now this is not as much of an issue. Still, it's a big job building LR tables for any substantial grammar by hand.

The method for constructing the configuring sets of LR(1) configurations is essentially the same as for SLR, but there are some changes in the closure and successor operations because we must respect the configuration lookahead. To compute the closure of an LR(1) configuring set  $I$ :

Repeat the following until no more configurations can be added to state  $I$ :

- For each configuration  $[A \rightarrow \underline{u} \bullet B \underline{v}, a]$  in  $I$ , for each production  $B \rightarrow \underline{w}$  in  $G'$ , and for each terminal  $b$  in  $\text{First}(\underline{va})$  such that  $[B \rightarrow \bullet \underline{w}, b]$  is not in  $I$ : add  $[B \rightarrow \bullet \underline{w}, b]$  to  $I$ .

What does this mean? We have a configuration with the dot before the non-terminal  $B$ . In LR(0), we computed the closure by adding all  $B$  productions with no indication of what was expected to follow them. In LR(1), we are a little more precise—we add each  $B$  production but insist that each have a lookahead of  $\underline{va}$ . The lookahead will be  $\text{First}(\underline{va})$  since this is what follows  $B$  in this production. Remember that we can compute first sets not just for a single non-terminal, but also a sequence of terminal and non-terminals.  $\text{First}(\underline{va})$  includes the first set of the first symbol of  $\underline{v}$  and then if that symbol is nullable, we include first set of the following symbol, and so on. If the entire sequence  $\underline{v}$  is nullable, we add the lookahead  $a$  already required by this configuration.

The successor function for the configuring set  $I$  and symbol  $X$  is computed as this:

Let  $J$  be the configuring set  $[A \rightarrow \underline{u} X \bullet \underline{v}, a]$  such that  $[A \rightarrow \underline{u} X \underline{v}, a]$  is in  $I$ .  
 $\text{successor}(I, X)$  is the closure of configuring set  $J$ .

We take each production in a configuring set, move the dot over a symbol and close on the resulting production. This is pretty much the same successor function as defined for LR(0), but we have to propagate the lookahead when computing the transitions.

We construct the complete family of all configuring sets  $F$  just as we did before.  $F$  is initialized to the set with the closure of  $[S' \rightarrow S, \$]$ . For each configuring set  $I$  and each grammar symbol  $X$  such that  $\text{successor}(I, X)$  is not empty and not in  $F$ , add  $\text{successor}(I, X)$  to  $F$  until no other configuring set can be added to  $F$ .

Let's consider an example. The augmented grammar below that recognizes the regular language  $a^*ba^*b$  (this example from pp. 231-236 Aho/Sethi/Ullman).

- 0)  $S' \rightarrow S$
- 1)  $S \rightarrow XX$
- 2)  $X \rightarrow aX$
- 3)  $X \rightarrow b$

Here is the family of LR configuration sets:

$I_0:$	$S' \rightarrow \bullet S, \$$ $S \rightarrow \bullet XX, \$$ $X \rightarrow \bullet aX, a/b$ $X \rightarrow \bullet b, a/b$	$I_4:$	$X \rightarrow b \bullet, a/b$
$I_1:$	$S' \rightarrow S \bullet, \$$	$I_5:$	$S \rightarrow XX \bullet, \$$
$I_2:$	$S \rightarrow X \bullet X, \$$ $X \rightarrow \bullet aX, \$$ $X \rightarrow \bullet b, \$$	$I_6:$	$X \rightarrow a \bullet X, \$$ $X \rightarrow \bullet aX, \$$ $X \rightarrow \bullet b, \$$
$I_3:$	$X \rightarrow a \bullet X, a/b$ $X \rightarrow \bullet aX, a/b$ $X \rightarrow \bullet b, a/b$	$I_7:$	$X \rightarrow b \bullet, \$$
		$I_8:$	$X \rightarrow aX \bullet, a/b$
		$I_9:$	$X \rightarrow aX \bullet, \$$

The above grammar would only have seven SLR states, but has ten in canonical LR. We end up with a lot more states because we have split states that have different lookaheads. For example, states 3 and 6 are the same except for lookahead, state 3 corresponds to the context where we are in the middle of parsing the first X, state 6 is the second X. Similarly, states 4 and 7 are completing the first and second X respectively. In SLR, those states are not distinguished, and if we were attempting to parse a single b by itself, we would allow that to be reduced to X, even though this will not lead to a valid sentence. The SLR parser will eventually notice the syntax error, but the LR parser figures it out a bit sooner.

To fill in the entries in the action and goto tables, we use the pretty much the same algorithm as we did for SLR(1). It is a bit easier to assign the reduce actions for an LR(1) table because the lookahead sets tell us exactly where to put each reductions.

1. Construct  $F = \{I_0, I_1, \dots, I_n\}$ , the collection of configuring sets for the augmented grammar  $G'$  (augmented by adding the special production  $S' \rightarrow S$ ).
2. State  $i$  is determined from  $I_i$ . The parsing actions for the state are determined as follows:
  - a) If  $[A \rightarrow \underline{u} \bullet, a]$  is in  $I_i$  then set  $\text{Action}[i, a]$  to reduce  $A \rightarrow \underline{u}$  (note  $A$  may not be  $S'$ ).
  - b) If  $[S' \rightarrow S \bullet, \$]$  is in  $I_i$  then set  $\text{Action}[i, \$]$  to accept.
  - c) If  $[A \rightarrow \underline{u} \bullet a \underline{v}, b]$  is in  $I_i$  and  $\text{succ}(I_i, a) = I_j$ , then set  $\text{Action}[i, a]$  to shift  $j$  (a must be a terminal).
3. The goto transitions for state  $i$  are constructed for all non-terminals  $A$  using the rule: If  $\text{succ}(I_i, A) = I_j$ , then  $\text{Goto}[i, A] = j$ .
4. All entries not defined by rules 2 and 3 are errors.
5. The initial state is the one constructed from the configuring set containing  $S' \rightarrow \bullet S$ .

Following the algorithm using the configuring sets given above, we construct this canonical LR parse table:

State on top of stack	Action			Goto	
	a	b	\$	S	X
0	s3	s4		1	2
1			acc		
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8	r2	r2			
9			r2		

Let's parse an example string baab. It is a valid in this language as shown by this leftmost derivation:

S → XX  
       bX  
      baX  
     baaX  
    baab

STACK STACK	REMAINING INPUT	PARSER ACTION
S <sub>0</sub>	baab\$	Shift S <sub>4</sub>
S <sub>0</sub> S <sub>4</sub>	aab\$	Reduce 3) X → b, goto S <sub>2</sub>
S <sub>0</sub> S <sub>2</sub>	aab\$	Shift S <sub>6</sub>
S <sub>0</sub> S <sub>2</sub> S <sub>6</sub>	ab\$	Shift S <sub>6</sub>
S <sub>0</sub> S <sub>2</sub> S <sub>6</sub> S <sub>6</sub>	b\$	Shift S <sub>7</sub>
S <sub>0</sub> S <sub>2</sub> S <sub>6</sub> S <sub>6</sub> S <sub>7</sub>	\$	Reduce 3) X → b, goto S <sub>9</sub>
S <sub>0</sub> S <sub>2</sub> S <sub>6</sub> S <sub>6</sub> S <sub>9</sub>	\$	Reduce 2) X → aX, goto S <sub>9</sub>
S <sub>0</sub> S <sub>2</sub> S <sub>6</sub> S <sub>9</sub>	\$	Reduce 2) X → aX, goto S <sub>5</sub>
S <sub>0</sub> S <sub>2</sub> S <sub>5</sub>	\$	Reduce 1) S → XX, goto S <sub>1</sub>
S <sub>0</sub> S <sub>1</sub>	\$	Accept

Now, let's consider what the states mean. S<sub>4</sub> is where X → b is completed; S<sub>2</sub> and S<sub>6</sub> is where we are in the middle of processing the 2 a's; S<sub>7</sub> is where we process the final b; S<sub>9</sub> is where we complete the X → aX production; S<sub>5</sub> is where we complete S → XX; and S<sub>1</sub> is where we accept.

### LR(1) grammars

Every SLR(1) grammar is a canonical LR(1) grammar, but the canonical LR(1) parser may have more states than the SLR(1) parser for same grammar. An LR(1) grammar may not be SLR(1), the grammar given earlier is such an example. Those grammars require splitting states to resolve shift-reduce conflicts that would otherwise result if lookaheads were not carried with the states.

A grammar is LR(1) if the following two conditions are satisfied for each configuring set:

3. For any item in the set  $[A \rightarrow \underline{u} \bullet x \underline{v}, a]$  with  $x$  a terminal, there is no item in the set of the form  $[B \rightarrow \underline{u} \bullet, x]$  In the action table, this translates no shift-reduce conflict for any state. The successor function for  $x$  either shifts to a new state or reduces, but not both.
4. The lookaheads for all complete items within the set must be disjoint, e.g. set cannot have both  $[A \rightarrow \underline{u} \bullet, a]$  and  $[B \rightarrow \underline{v} \bullet, a]$  This translates to no reduce-reduce conflict on any state. If more than one non-terminal could be reduced from this set, it must be possible to uniquely determine which is appropriate from the next input token.

As long as there is a unique shift or reduce action on each input symbol from each state, we can parse using an LR(1) algorithm. The above state conditions are similar to what is required for SLR(1), but rather than the looser constraint about disjoint follow sets and so on, canonical LR(1) computes a more precise notion of the appropriate lookahead within a particular context and thus is able to resolve conflicts that SLR(1) would encounter.

### **Bibliography**

- A. Aho, R. Sethi, J. Ullman, Compilers: Principles, Techniques, and Tools. Reading, MA: Addison-Wesley, 1986.
- J.P. Bennett, Introduction to Compiling Techniques. Berkshire, England: McGraw-Hill, 1990.
- K. Loudon, Compiler Construction. Boston, MA: PWS, 1997
- A. Pyster, Compiler Design and Construction. New York, NY: Van Nostrand Reinhold, 1988.