

Compilers Course

Lecture 6: Bottom-Up (Shift-Reduce) Parsing

A shift-reduce parser works as follows:

- Read tokens one by one and push them on the stack (shift) .
- If the top N symbols on the stack equal β in a production $X \rightarrow \beta$, and the prefix of the stack "expects" an X, replace β with X on the stack (reduce) .
- Accept if the stack is just the start symbol S and the input is empty .

Note: the difference with LL(1), here choices are based on having seen complete right-hand sides. LL(1) has to choose having only seen a single token that may be at the start of a right-hand side.

Example: Consider this grammar G1:

$S \rightarrow \text{if } E \text{ then } S$
 $\quad \quad \quad | \text{id} \quad \quad \quad (\text{proc. call})$
 $E \rightarrow \text{id} \mid \text{num}$

Stack	Input	Action
\$	if id then id\$	shift
\$if	id then id\$	shift
\$if id	then id\$	
Problem: reduce $S \rightarrow \text{id}$ or $E \rightarrow \text{id}$? Answer: reduce $E \rightarrow \text{id}$		
\$if E	then id\$	shift
\$if E then	id\$	shift
\$if E then id	\$	
Problem: reduce $S \rightarrow \text{id}$ or $E \rightarrow \text{id}$? Answer: reduce $S \rightarrow \text{id}$		
\$if E then S	\$	reduce $S \rightarrow \text{if } E \text{ then } S$
\$\$	\$	accept

This corresponds to the right-most derivation:

$S \Rightarrow_{\text{rm}} \text{if } E \text{ then } S \Rightarrow_{\text{rm}} \text{if } E \text{ then id} \Rightarrow_{\text{rm}} \text{if id then id}$

which the parser traced in reverse order.

Choosing between reductions

The central question for a shift-reduce parser is:

“If two productions match the top of the stack, how does the parser choose which one to reduce?”

(If no reduction is applicable, the machine shifts.)

Tracking Candidates

In steps 3 and 6 above, the top of the stack matched both $S \rightarrow id$ and $E \rightarrow id$. Why did the parser choose the second in step 3 and the first in step 6? The answer is that by scanning the stack the parser can identify the context to the left of the production it might reduce. That tells it if the production is in fact valid for that context.

Consider a transition system for a shift-reduce parser where states are pairs $\langle \text{stack of symbols, a candidate production with a "progress" point} \rangle$.

In the initial state the stack is empty, and the candidate is the start symbol. We augment the grammar (add $S' \rightarrow S$), and describe this as:

$\langle \epsilon, S' \rightarrow .S \rangle$

If the candidate in a state has a progress point before a symbol X , then the parser can shift that symbol onto the stack and move the progress point past that symbol:

$\langle \gamma \alpha, A \rightarrow \alpha . X \beta \rangle \xrightarrow{X} \langle \gamma \alpha X, A \rightarrow \alpha X . \beta \rangle$

If the candidate in a state has a progress point before a non-terminal B , then the productions of B are new candidates:

$\langle \gamma \alpha, A \rightarrow \alpha . B \beta \rangle \xrightarrow{\epsilon} \langle \gamma \alpha, B \rightarrow . \delta \rangle$
for every production $B \rightarrow \delta$

Finally, if at some point the machine arrives at a state $\langle \gamma \alpha, A \rightarrow \alpha . \rangle$, then $A \rightarrow \alpha$ is both "complete" and known to be valid. So the machine could reduce the stack to γA at that point.

The stacks don't actually affect the transitions but instead describe the path taken from the start state to the current state. Erasing the stacks turns this machine into an NFA_ϵ , with "candidate productions" as states.

Formal definitions

Consider a step in a right-most derivation:

$$S \Rightarrow_{*_{rm}} \delta A w \Rightarrow_{rm} \delta \beta w$$

w: remaining input

β : the handle (top of the stack)

$\delta \beta$: a viable prefix (the entire stack)

Every prefix of $\delta \beta$ is a viable prefix.

An LR(0) item is a production $A \rightarrow X_1 \dots X_n$ with a dot before any X_i or after X_n .

Example: $S \rightarrow abC$ has items $S \rightarrow . a b C$, $S \rightarrow a . b C$, $S \rightarrow a b . C$, and $S \rightarrow a b C .$

An item with the dot before X_1 is an initial item.

An item with the dot after X_n is a complete item.

An item $A \rightarrow \alpha . \beta$ is valid for a viable prefix $\delta \alpha$ if

$$S \Rightarrow_{*_{rm}} \delta A w \Rightarrow_{rm} \delta \alpha \beta w$$

If $A \rightarrow \delta .$ is a complete item valid for $\delta \alpha$ (the stack), then

$$\delta A w \Rightarrow_{rm} \delta \alpha w$$

could have been the last step in a right-most derivation of $\delta \alpha w$. It could be that $(A \rightarrow \alpha .)$ is valid for some other remaining input w' , or there could be other complete items valid for $\delta \alpha$.

Intuitively, a grammar is defined to be LR(0) if $(A \rightarrow \alpha .)$ is the only complete item valid for $\delta \alpha$.

Parsing a string w then becomes a process of finding handles and doing backwards derivations until w has been reduced to S .

Recognising valid items

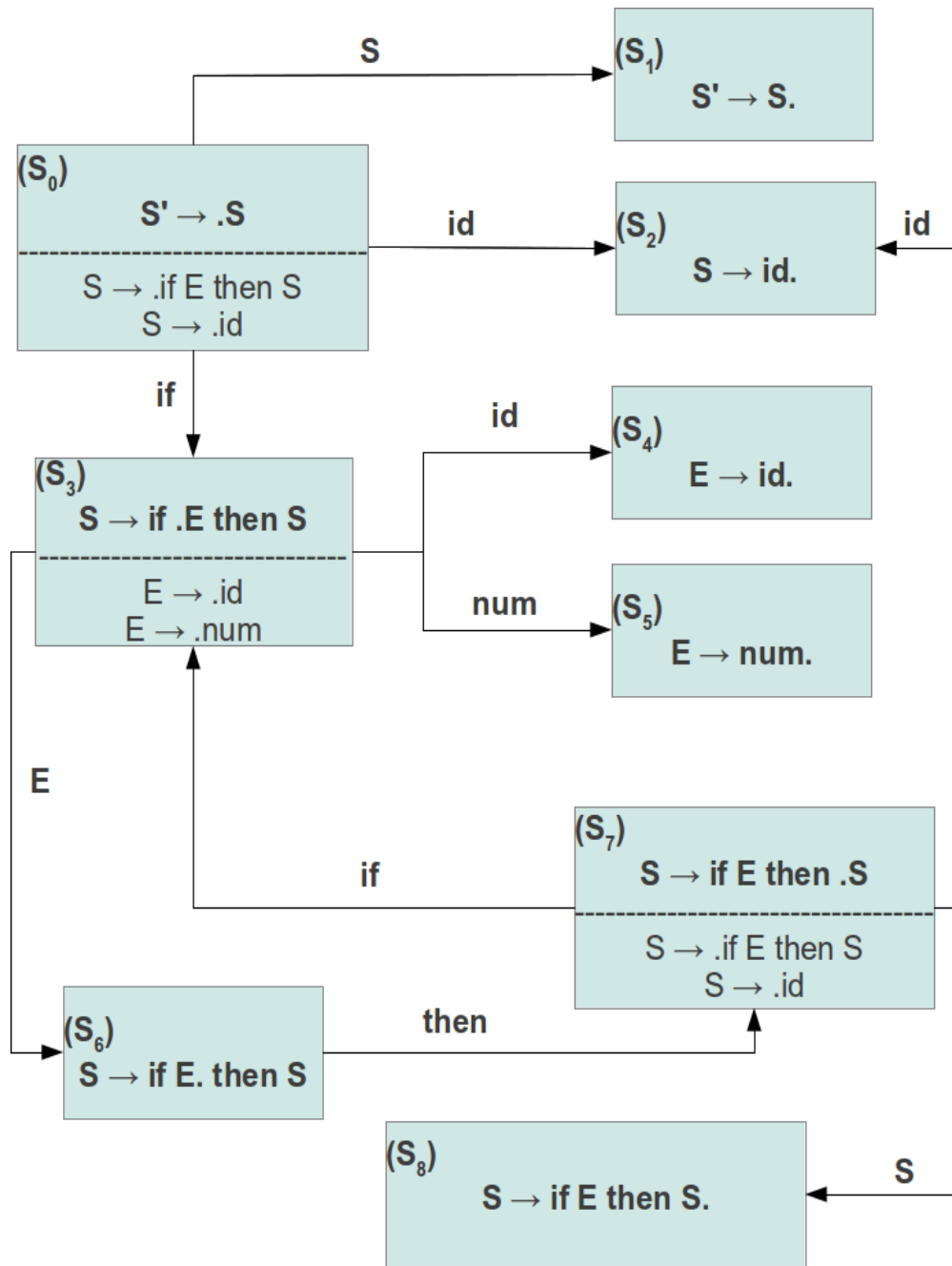
For each viable prefix $\delta \alpha$ (stack of symbols), we must now be able to compute its set of valid items.

Theorem: for every CFG G , its set of viable prefixes is a regular language recognized by an NFA_ϵ M whose states are LR(0) items for G . M is defined as follows:

1. $\delta[q_0, \epsilon] = \{ S' \rightarrow . S \}$ $S' \rightarrow S$ is added to the grammar and S' is made the new start symbol; this is called augmenting the grammar .
2. $\delta[A \rightarrow \alpha . B \beta, \epsilon] = \{ B \rightarrow . \gamma \mid B \rightarrow \gamma \text{ is a production} \}$ dot before variable implies ϵ transitions to initial items for each production of that variable .
3. $\delta[A \rightarrow \alpha . X \beta, X] = \{ A \rightarrow \alpha X . \beta \}$ this is called "moving the dot" .

Example: (grammar G1)**Augmented grammar:** $(p_0) S' \rightarrow S$ $(p_1) S \rightarrow \text{if } E \text{ then } S$ $(p_2) S \rightarrow \text{id}$ $(p_3) E \rightarrow \text{id}$ $(p_4) E \rightarrow \text{num}$ **Items/States:** $I_0: S' \rightarrow \cdot S$ $I_1: S' \rightarrow S \cdot$ $I_2: S \rightarrow \cdot \text{if } E \text{ then } S$ $I_3: S \rightarrow \text{if } \cdot E \text{ then } S$ $I_4: S \rightarrow \text{if } E \cdot \text{ then } S$ $I_5: S \rightarrow \text{if } E \text{ then } \cdot S$ $I_6: S \rightarrow \text{if } E \text{ then } S \cdot$ $I_7: S \rightarrow \cdot \text{id}$ $I_8: S \rightarrow \text{id} \cdot$ $I_9: E \rightarrow \cdot \text{id}$ $I_{10}: E \rightarrow \text{id} \cdot$ $I_{11}: E \rightarrow \cdot \text{num}$ $I_{12}: E \rightarrow \text{num} \cdot$ **Transitions:** $\delta[I_0, \varepsilon] \rightarrow \{I_2, I_7\}$ $\delta[I_0, S] \rightarrow I_1$ $\delta[I_2, \text{if}] \rightarrow I_3$ $\delta[I_3, \varepsilon] \rightarrow \{I_9, I_{11}\}$ $\delta[I_3, E] \rightarrow I_4$ $\delta[I_4, \text{then}] \rightarrow I_5$ $\delta[I_5, \varepsilon] \rightarrow \{I_2, I_7\}$ $\delta[I_5, S] \rightarrow I_6$ $\delta[I_7, \text{id}] \rightarrow I_8$ $\delta[I_9, \text{id}] \rightarrow I_{10}$ $\delta[I_{11}, \text{num}] \rightarrow I_{12}$

Conversion to DFA:



The items above "-----" are called the "kernel", the items below come from ε-transitions and are called the "closure".

At each step, the shift-reduce parser applies this automaton to its stack. Let's say it ends at some state (set of items) I :

- If I contains a complete item ($A \rightarrow \alpha \cdot$), then reduce
- If I contains a non-complete item, the current input terminal is T , and I has a transition on T , then shift
- If I contains the complete item ($S' \rightarrow S \cdot$) and the input is empty, accept

Example 1: (grammar G1)

\$	if id then id\$	State S ₀ : shift on if
\$if	id then id\$	State S ₃ : shift on id
\$if id	then id\$	State S ₄ : reduce $E \rightarrow id$
\$if E	then id\$	State S ₆ : shift on then
\$if E then	id\$	State S ₇ : shift on id
\$if E then id	\$	State S ₂ : reduce $S \rightarrow id$
\$if E then S	\$	State S ₈ : reduce \rightarrow if E then S
\$\$	\$	State S ₁ : accept

LR(0) action/goto tables

Given a DFA recognizing viable prefixes, the LR(0) decision tables are defined as follows:

- **Action:** matrix indexed by *states* and *terminals*, containing *actions*:

Action[q, X] = shift q'

if X is a terminal and $q \rightarrow q'$ on X in the DFA

Action[q, X] = reduce $A \rightarrow \alpha$

if $A \rightarrow \alpha \cdot$ is a complete item in q

for all terminals X

Action[q, \$] = accept

if $S' \rightarrow S \cdot$ is a complete item in q

- **Goto:** matrix indexed by *states* and *variables*, containing *states*

Goto[q, X] = q'

if X is a variable and $q \rightarrow q'$ on X

Action						GOTO	
	if	then	id	num	\$	S	E
0	s ₃		s ₂			g ₁	
1					A		
2	r ₂	r ₂	r ₂	r ₂	r ₂		
3			s ₄	s ₅			g ₆
4	r ₃	r ₃	r ₃	r ₃	r ₃		
5	r ₄	r ₄	r ₄	r ₄	r ₄		
6		s ₇					
7	s ₃		s ₂			g ₈	
8	r ₁	r ₁	r ₁	r ₁	r ₁		

(S_N means "shift to state N", r_N means "reduce by production N", A means "accept", g_N means "goto state N")

Implementation note 1: instead of scanning the entire stack in each step, push the new state after pushing each symbol. Base decisions only on top-most state and current input symbol.

Stack	Input	Action
\$0	if id then id\$	Shift 3
\$0 if 3	id then id\$	Shift 4
\$0 if 3 id 4	then id\$	Reduce E → id pop 2 * id symbols, exposing state 3 push E, goto[3, E] = 6
\$0 if 3 E 6	then id\$	Shift 7
\$0 if 3 E 6 then 7	id\$	Shift 2
\$0 if 3 E 6 then 7 id 2	\$	Reduce S → id pop 2 * id symbols, exposing state 7 push S, goto[7, S] = 8
\$0 if 3 E 6 then 7 S 8	\$	Reduce S → if E then S pop 2 * if E then S symbols, exposing state 0 push S, goto[0, S]-1
\$0 S 1	\$	accept

Implementation note 2: the symbols on the stack are now redundant, so only push states.

LR(0) Parsing Algorithm

Let S be a stack of states, and SP the current position in S.

Init: SP := 0; S[SP] := q0

Loop:

```

let q = S[SP] and t = current input terminal
if action[q,t] = shift q' then:
    push q' (SP := SP+1; S[SP] := q')
    remove t from input
    continue at Loop
if action[q,t] = reduce A → β then:
    pop β from S (SP := SP - |β|)
    let q' = S[SP] and q'' = goto[q',A]
    push q'' (SP := SP+1; S[SP] := q'')
    continue at Loop
if action[q,t] = accept then
    done

```

Semantic actions, construction of AST

Each position in the state stack corresponds to a terminal, or a variable that is the result of reducing some production.

Both terminals and variables can have attributes.

We will add a "semantic stack" SEM that parallels the state stack S, such that SEM[SP] contains the attribute(s) for the symbol that branched to the state S[SP].

At actions shift q':

push t's attribute on SEM (SEM[SP] := t's attribute)

At actions reduce $A \rightarrow \beta$:

evaluate synthesized attribute before popping S
after pushing the new state, push attribute on SEM

Example:

$T \rightarrow T \times F$ (MKBINARY(MUL,T,F))

T := SEM[SP-2]

F := SEM[SP-0]

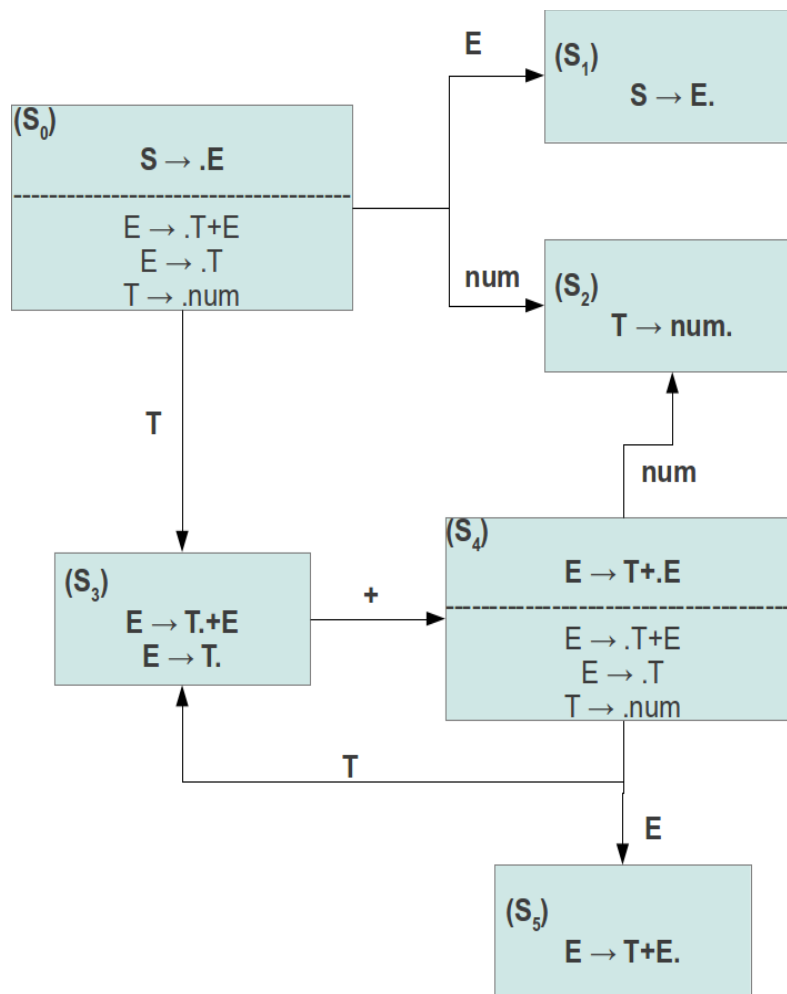
SP := (SP - 3) + 1

SEM[SP] := MKBINARY(MUL,T,F)

Simple LR, SLR

LR(0) is still fairly weak; consider:

$S \rightarrow E$
 $E \rightarrow T+E$
 $E \rightarrow T$
 $T \rightarrow \text{num}$



When defining the Action table,

Action[S₃, +] is shift 4, BUT Action[S₃, +] is reduce $E \rightarrow T$ because $E \rightarrow T$ is a complete item, and LR(0) reduces for all terminals when seeing a complete item.

This is a classic shift/reduce conflict, caused in this case by LR(0)'s failure to use lookahead (the current input terminal) to decide whether to reduce or not.

Simple LR (SLR):

Action[q, X] = reduce $A \rightarrow \alpha$
 if $(A \rightarrow \alpha \cdot)$ is a complete item in q
 $\forall X \in \text{FOLLOW}(A)$

In the example, $\text{FOLLOW}(S) = \{\$\}$, $\text{FOLLOW}(E) = \text{FOLLOW}(S) = \{\$\}$, and $\text{FOLLOW}(T) = \{+, \$\}$.
 So in S_3 we choose shift 4 for + and reduce $E \rightarrow T$ for \$.

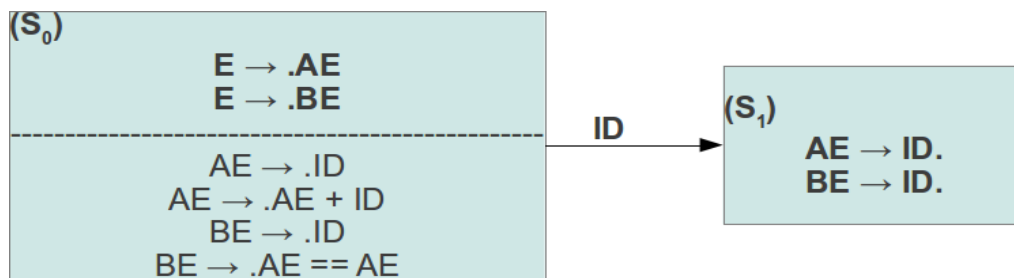
SLR is strong enough to handle many actual programming languages.

Reduce/Reduce Conflicts

This grammar attempts to describe arithmetic and boolean expressions as separate entities:

$E \rightarrow AE \mid BE$
 $AE \rightarrow ID \mid AE + ID$
 $BE \rightarrow ID \mid AE == AE$

Now consider the LR(0) DFA for this grammar:



Now, $\text{FOLLOW}(E) = \{\$\}$, $\text{FOLLOW}(AE) = \{+, ==, \$\}$, and $\text{FOLLOW}(BE) = \text{FOLLOW}(E) = \{\$\}$. So when the current input terminal is \$ (in S_1), we can reduce either $AE \rightarrow ID$ or $BE \rightarrow ID$.

This is called a reduce/reduce conflict, and is usually a sign of an ambiguous grammar.

In this case, the ambiguity comes from the fact that some arithmetic and boolean expressions look the same, so the grammar is unable to decide which interpretation to make.

The typical solution is to delay attempts at "type checking" to a separate type checking pass that runs after the parser has produced a parse tree. A better grammar for this language is:

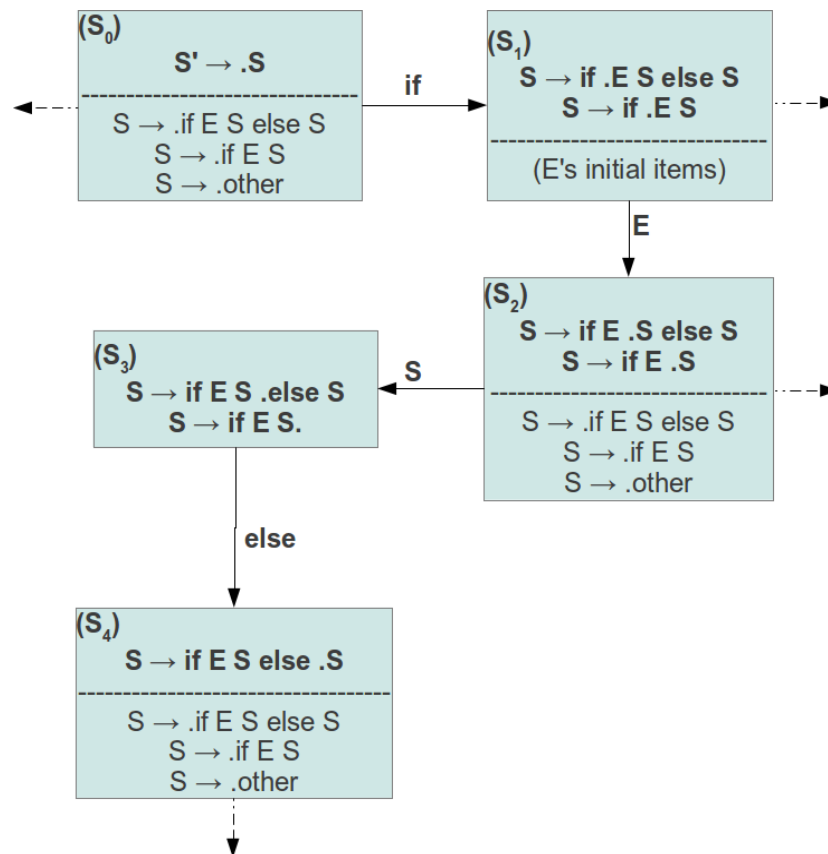
$E \rightarrow BE$
 $BE \rightarrow AE == AE \mid AE$
 $AE \rightarrow AE + ID \mid ID$

Dangling-Else

$S' \rightarrow S$
 $S \rightarrow \text{if } E \text{ } S \text{ else } S$
 $S \rightarrow \text{if } E \text{ } S$
 $S \rightarrow \text{other}$

This is an ambiguous grammar for the dangling-else construct.

Constructing the LR(0) DFA for this grammar results in:



Now, $\text{FOLLOW}(S) = \{\text{else}\} \cup \text{FOLLOW}(S') = \{\text{else}, \$\}$, so in S_3 we should shift to S_4 on else, but we should also reduce $S \rightarrow \text{if } E \text{ } S$ on else.

So the well-known dangling-else ambiguity results in a shift/reduce conflict in the LR parser.

One solution is to make the grammar unambiguous, as shown earlier.

However, this case is very common so LR parser generators have a heuristic which is to *prefer shifting over reducing* whenever there is a shift/reduce conflict. In general that's not correct, but for the

dangling-else it *is* the right solution since it leads to the else part being connected to the nearest preceding if that doesn't already have an else part.

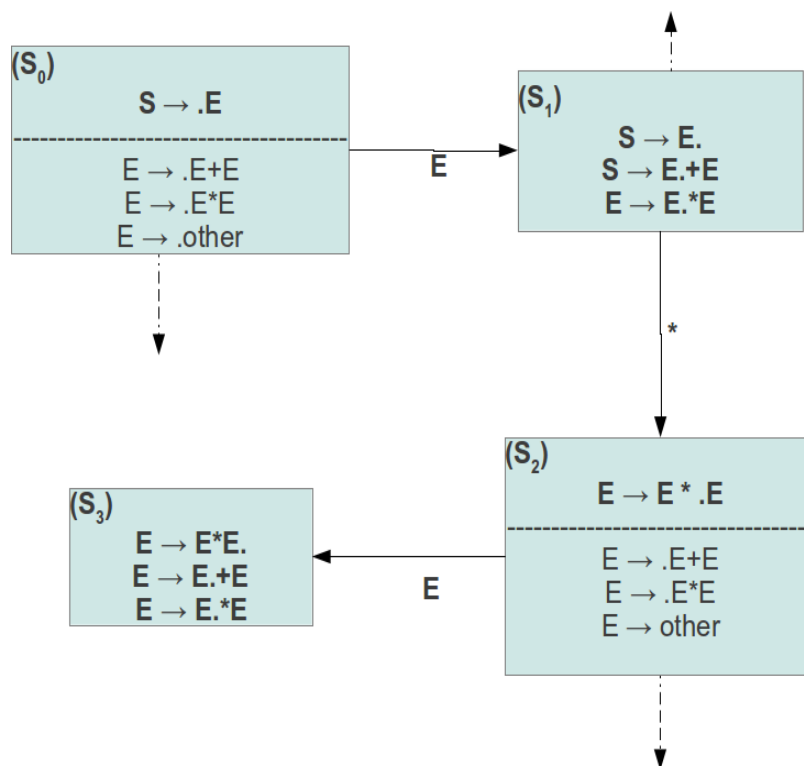
Precedence and Associativity

Consider the ambiguous grammar:

$$S \rightarrow E$$

$$E \rightarrow E + E \mid E * E \mid \text{other}$$

Let's see what happens in the LR(0) DFA:



Now, since $\text{FOLLOW}(E)$ includes $+$ and $*$, we will have two shift/reduce conflicts in S_3 :

1. On $*$ shift or reduce $E \rightarrow E^*E$ this conflict comes from the grammar's failure to express associativity of $*$
2. On $+$ shift or reduce $E \rightarrow E^*E$ this conflict comes from the grammar's failure to express precedence between $+$ and $*$

Yacc and similar LR parser generators accept declarations to resolve these kinds of conflicts:

```
%left <token>
```

When there is a choice between reducing a rule containing the $\langle \text{token} \rangle$ or shifting that $\langle \text{token} \rangle$, the

parser will choose the reduction, which effectively makes the $\langle \text{token} \rangle$ left-associative. So "%left *" would solve ambiguity number 1 above. Similarly, %right $\langle \text{token} \rangle$ causes the parser choose to shift instead of reduce, making the $\langle \text{token} \rangle$ right-associative. The order of the %left and %right declarations cause the tokens to be given different precedences:

```
%left +
%left *
```

means that * is given higher precedence than +: when there is a choice between reducing a rule involving one token or shifting another token, the choice involving the higher-precedence token is the one taken.

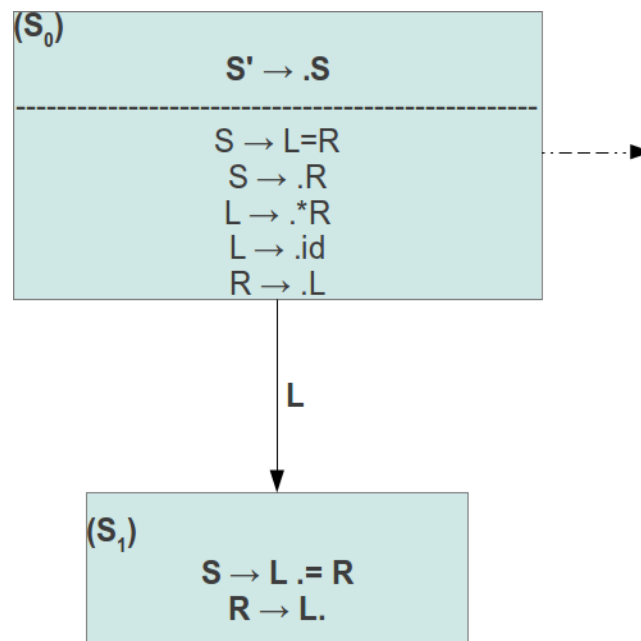
LR(1)

SLR is stronger than LR(0) since it uses lookahead to decide whether to reduce or not, which eliminates some shift/reduce and reduce/reduce conflicts.

Example: The grammar in Example 4.48 of the Purple Dragon book:

```
S' → S
S → L=R | R
L → *R | id
R → L
```

LR(0) DFA:



Now, $FOLLOW(R) = \{ =, \$ \}$ (because $S \rightarrow L=R$ and $L \rightarrow *R$), so SLR wants to shift on = and also

reduce $R \rightarrow L$ on $=$.

The problem with SLR is that the lookahead information used is FOLLOW, which is a global property for *all* uses of the non-terminal. It is in general an over-approximation for a specific DFA state with a complete item.

LR(1) changes the NFA construction to explicitly include and propagate accurate lookahead information:

1. Items are pairs $\langle \text{LR}(0) \text{ item, lookahead terminal symbol} \rangle$
2. $\delta[q_0, \epsilon] = \{ \langle S' \rightarrow .S, \$ \rangle \}$
3. $\delta[\langle A \rightarrow \alpha . B \beta, Z \rangle, \epsilon] = \{ \langle B \rightarrow . \gamma, y \rangle \mid B \rightarrow \gamma \text{ is a production, and } y \in \text{FIRST}(\beta Z) \}$
 B is followed by β in this context, so $\text{FIRST}(\beta)$ provides the possible lookahead symbols for B 's initial items. If β derives the empty string, the context's lookahead Z is propagated instead. Hence $\text{FIRST}(\beta Z)$.
4. $\delta[\langle A \rightarrow \alpha . X \beta, Z \rangle, X] = \{ \langle A \rightarrow \alpha X . \beta, Z \rangle \}$ moving the dot doesn't change the lookahead symbol

In the action table:

Action $[q, X] = \text{reduce } A \rightarrow \alpha$ if $\langle A \rightarrow \alpha ., X \rangle$ is an item in q . The item provides its own lookahead symbol, and the reduction is only done if the input matches the lookahead.

LR(1) is much stronger than SLR or LR(0), in fact, LR(1) corresponds exactly to the class of "deterministic push-down automata".

Unfortunately, the number of LR(1) states can be orders of magnitude larger than the number of SLR/LR(0) states, so LR(1) parsers are rarely used in practice.

LALR(1)

LALR(1) stands for "lookahead LR(1)". In principle:

1. Build the LR(1) DFA
2. Merge states with identical LR(0) items, and use the union of their lookaheads

The result is a DFA with the same number of states as SLR and LR(0) would have, but with items $\langle \text{LR}(0) \text{ item, set of lookahead symbols} \rangle$.

In practice, the LALR(1) DFA is constructed from the LR(0) DFA by a process of lookahead propagation.

LALR(1) is implemented by tools like Yacc and its clones (Bison, ML-Yacc), and is very popular among language designers and compiler implementors.