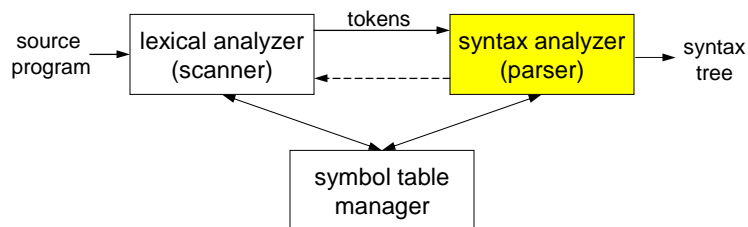# CSc 453
# Syntax Analysis (Parsing)

Saumya Debray

*The University of Arizona*

*Tucson*

---

## Overview



*Main Task*: Take a token sequence from the scanner and verify that it is a syntactically correct program.

*Secondary Tasks*:

- Process declarations and set up symbol table information accordingly, in preparation for semantic analysis.
- Construct a syntax tree in preparation for intermediate code generation.

## Context-free Grammars

- A *context-free grammar* for a language specifies the syntactic structure of programs in that language.
- Components of a grammar:
  - a finite set of tokens (obtained from the scanner);
  - a set of variables representing "related" sets of strings, e.g., *declarations*, *statements*, expressions.
  - a set of rules that show the structure of these strings.
  - an indication of the "top-level" set of strings we care about.

## Context-free Grammars: Definition

Formally, a context-free grammar $G$ is a 4-tuple $G = (V, T, P, S)$, where:

- V is a finite set of <u>*variables*</u> (or <u>*nonterminals*</u>). These describe sets of "related" strings.
- T is a finite set of <u>*terminals*</u> (i.e., tokens).
- P is a finite set of <u>*productions*</u>, each of the form
    $$A \rightarrow \alpha$$
  where $A \in V$ is a variable, and $\alpha \in (V \cup T)^*$ is a sequence of terminals and nonterminals.
- $S \in V$ is the <u>*start symbol*</u>.

# Context-free Grammars: An Example

A grammar for palindromic bit-strings:

$G = (V, T, P, S)$, where:

- $V = \{ S, B \}$
- $T = \{0, 1\}$
- $P = \{ S \rightarrow B,$
  - $S \rightarrow \varepsilon,$
  - $S \rightarrow 0\ S\ 0,$
  - $S \rightarrow 1\ S\ 1,$
  - $B \rightarrow 0,$
  - $B \rightarrow 1$
  - $\}$

---

# Context-free Grammars: Terminology

- *Derivation*: Suppose that
  - $\alpha$ and $\beta$ are strings of grammar symbols, and
  - $A \rightarrow \gamma$ is a production.

  Then, $\alpha A \beta \Rightarrow \alpha\gamma\beta$ ("$\alpha A \beta$ *derives* $\alpha\gamma\beta$").

- $\Rightarrow$ : "derives in one step"

  $\Rightarrow^*$ : "derives in 0 or more steps"

  $\alpha \Rightarrow^* \alpha$                                  (0 steps)

  $\alpha \Rightarrow^* \beta$ if $\alpha \Rightarrow \gamma$ and $\gamma \Rightarrow^* \beta$     ($\geq 1$ steps)

# Derivations: Example

- Grammar for palindromes: G = (V, T, P, S),
  - V = {S},
  - T = {0, 1},
  - P = { S → 0 S 0 | 1 S 1 | 0 | 1 | ε }.
- A derivation of the string 10101:

  S
  ⇒ 1 S 1      (using S → 1S1)
  ⇒ 1 0S0 1    (using S → 0S0)
  ⇒ 10101     (using S → 1)

# Leftmost and Rightmost Derivations

- A *leftmost derivation* is one where, at each step, the leftmost nonterminal is replaced.
  (analogous for *rightmost derivation*)
- *Example*: a grammar for arithmetic expressions:
  E → E + E | E * E | **id**
  - *Leftmost derivation*:
    E ⇒ E * E ⇒ E + E * E ⇒ **id** + E * E ⇒ **id** + **id** * E ⇒ **id** + **id** * **id**
  - *Rightmost derivation*:
    E ⇒ E + E ⇒ E + E * E ⇒ E + E * **id** ⇒ E + **id** * **id** ⇒ **id** + **id** * **id**
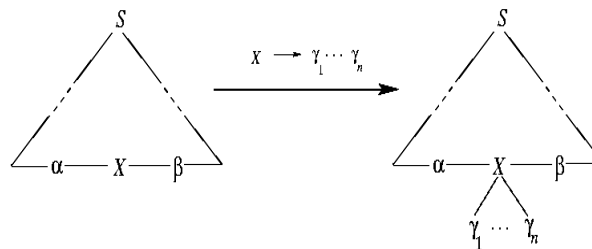
# Context-free Grammars: Terminology

- The *language* of a grammar $G = (V,T,P,S)$ is

    $L(G) = \{ w \mid w \in T^* \text{ and } S \Rightarrow^* w \}$.

    The language of a grammar contains only strings of terminal symbols.

- Two grammars $G_1$ and $G_2$ are *equivalent* if

    $L(G_1) = L(G_2)$.

---

# Parse Trees

- A *parse tree* is a tree representation of a derivation.
- Constructing a parse tree:
  - The root is the start symbol S of the grammar.
  - Given a parse tree for $\alpha\, X\, \beta$, if the next derivation step is
    $\alpha\, X\, \beta \Rightarrow \alpha\, \gamma_1 \ldots \gamma_n\, \beta$ then the parse tree is obtained as:

## Approaches to Parsing

- *Top-down parsing*:
  - attempts to figure out the derivation for the input string, starting from the start symbol.

- *Bottom-up parsing*:
  - starting with the input string, attempts to "derive in reverse" and end up with the start symbol;
  - forms the basis for parsers obtained from parser-generator tools such as yacc, bison.
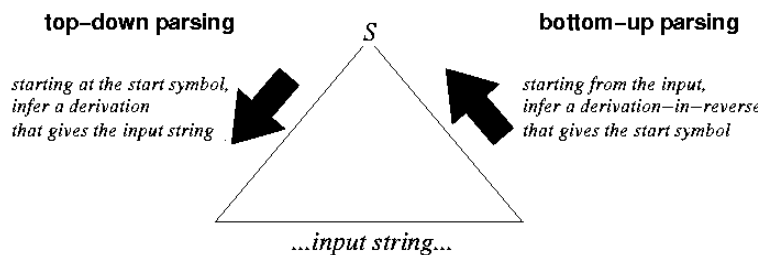
## Top-down Parsing

- "top-down:" starting with the start symbol of the grammar, try to derive the input string.
- *Parsing process*: use the current state of the parser, and the next input token, to guide the derivation process.
- *Implementation*: use a finite state automaton augmented with a runtime stack ("*pushdown automaton*").

## Bottom-up Parsing

- "bottom-up:" work backwards from the input string to obtain a derivation for it.
- *Parsing process*: use the parser state to keep track of:
  - what has been seen so far, and
  - given this, what the rest of the input might look like.
- *Implementation*: use a finite state automaton augmented with a runtime stack ("*pushdown automaton*").

## Parsing: Top-down vs. Bottom-up

top−down parsing                    $S$                    bottom−up parsing

starting at the start symbol,          starting from the input,
infer a derivation                     infer a derivation−in−reverse
that gives the input string            that gives the start symbol

...input string...

## Parsing Problems: Ambiguity

- A grammar G is *ambiguous* if some string in L(G) has more than one parse tree.
- Equivalently: if some string in L(G) has more than one leftmost (rightmost) derivation.
- *Example*: The grammar

  $E \rightarrow E + E \mid E * E \mid$ **id**

  is ambiguous, since "id+id*id" has multiple parses:

## Dealing with Ambiguity

1. Transform the grammar to an equivalent unambiguous grammar.
2. Use *disambiguating rules* along with the ambiguous grammar to specify which parse to use.

*Comment*: It is not possible to determine algorithmically whether:

- Two given CFGs are equivalent;
- A given CFG is ambiguous.

## Removing Ambiguity: Operators

- *Basic idea*: use additional nonterminals to enforce associativity and precedence:
    - Use one nonterminal for each precedence level:
        - $E \rightarrow E * E \mid E + E \mid$ id

        needs 2 nonterminals (2 levels of precedence).
    - Modify productions so that "lower precedence" nonterminal is in direction of precedence:

        $E \rightarrow E + E \quad \Rightarrow \quad E \rightarrow E + T$ (+ is left-associative)

## Example

- *Original grammar*:

    $E \rightarrow E * E \mid E / E \mid E + E \mid E - E \mid ( E ) \mid$ **id**

    precedence levels: $\{ *, / \} > \{ +, - \}$

    associativity: *, /, +, – are all left-associative.

- *Transformed grammar*:

    $E \rightarrow E + T \mid E - T \mid T$    (precedence level for: +, -)
    $T \rightarrow T * F \mid T / F \mid F$    (precedence level for: *, /)
    $F \rightarrow ( E ) \mid$ **id**

# Bottom-up parsing: Approach

1. Preprocess the grammar to compute some info about it.
   (FIRST and FOLLOW sets)
2. Use this info to construct a pushdown automaton for the grammar:
   - the automaton uses a table ("parsing table") to guide its actions;
   - constructing a parser amounts to constructing this table.

---

# FIRST Sets

_Defn_: For any string of grammar symbols $\alpha$,
  - FIRST($\alpha$) = { **a** | **a** is a terminal and $\alpha \Rightarrow^* $ **a**$\beta$}.
  - if $\alpha \Rightarrow^* \varepsilon$ then $\varepsilon$ is also in FIRST($\alpha$).

- _Example_:  E → T E′
    E′ → + T E′ | ε
    T → F T′
    T′ → * F T′ | ε
    F → ( E ) | **id**

  FIRST(E) = FIRST(T) = FIRST(F) = { (, **id** }
  FIRST(E′) = { +, ε }
  FIRST(T′) = { *, ε }

# Computing FIRST Sets

Given a sequence of grammar symbols $A$:
- **if** $A$ is a terminal or $A = \varepsilon$ **then** FIRST($A$) = {$A$}.
- **if** $A$ is a nonterminal with productions $A \rightarrow \alpha_1 \mid \ldots \mid \alpha_n$ **then**:
  - FIRST($A$) = FIRST($\alpha_1$) $\cup \ldots \cup$ FIRST($\alpha_n$).
- **if** $A$ is a sequence of symbols $Y_1 \ldots Y_k$ **then**:
  - **for** $i$ = 1 to $k$ **do**:
    - add each $a \in$ FIRST($Y_i$), such that $a \neq \varepsilon$, to FIRST($A$).
    - **if** $\varepsilon \notin$ FIRST($Y_i$) **then** break;
  - **if** $\varepsilon$ is in each of FIRST($Y_1$), ..., FIRST($Y_k$) **then** add $\varepsilon$ to FIRST($A$).

# Computing FIRST sets: cont'd

- For each nonterminal $A$ in the grammar, initialize FIRST($A$) = $\varnothing$.
- **repeat** {

  for each nonterminal $A$ in the grammar {

    compute FIRST($A$);   /* as described previously */

  }

  } **until** there is no change to any FIRST set.

# Example (FIRST Sets)

$$X \rightarrow YZ \ | \ \mathbf{a}$$
$$Y \rightarrow \mathbf{b} \ | \ \varepsilon$$
$$Z \rightarrow \mathbf{c} \ | \ \varepsilon$$

- $X \rightarrow \mathbf{a}$, so add **a** to FIRST(X).
- $X \rightarrow YZ$, $\mathbf{b} \in$ FIRST(Y), so add **b** to FIRST(X).
- $Y \rightarrow \varepsilon$, i.e. $\varepsilon \in$ FIRST(Y), so add non-$\varepsilon$ symbols from FIRST(Z) to FIRST(X).
  - ► add **c** to FIRST(X).
- $\varepsilon \in$ FIRST(Y) and $\varepsilon \in$ FIRST(Z), so add $\varepsilon$ to FIRST(X).

*Final*: FIRST(X) = { **a**, **b**, **c**, $\varepsilon$ }.

---

# FOLLOW Sets

*Definition*: Given a grammar G = (V, T, P, S), for any nonterminal $A \in$ V:

- FOLLOW($A$) = { $\mathbf{a} \in$ T | S $\Rightarrow^*$ $\alpha A \mathbf{a} \beta$ for some $\alpha, \beta$}.
  i.e., FOLLOW(A) contains those terminals that can appear after A in something derivable from the start symbol S.
- if S $\Rightarrow^*$ $\alpha A$ then $ is also in FOLLOW($A$).
  ($ $\equiv$ EOF, "end of input.")

*Example*:

$$E \rightarrow E + E \ | \ \mathbf{id}$$
$$FOLLOW(E) \ = \ \{ +, \$ \}.$$

# Computing FOLLOW Sets

Given a grammar G = (V, T, P, S):

1. add \$ to FOLLOW(S);

2. **repeat** {

   - for each production $A \rightarrow \alpha B\beta$ in P, add every non-$\varepsilon$ symbol in FIRST($\beta$) to FOLLOW($B$).

   - for each production $A \rightarrow \alpha B\beta$ in P, where $\varepsilon \in$ FIRST($\beta$), add everything in FOLLOW($A$) to FOLLOW($B$).

   - for each production $A \rightarrow \alpha B$ in P, add everything in FOLLOW($A$) to FOLLOW($B$).

   } **until** no change to any FOLLOW set.

# Example (FOLLOW Sets)

$$X \rightarrow YZ \ | \ \textbf{a}$$
$$Y \rightarrow \textbf{b} \ | \ \varepsilon$$
$$Z \rightarrow \textbf{c} \ | \ \varepsilon$$

- X is start symbol: add \$ to FOLLOW(X);
- $X \rightarrow YZ$, so add everything in FOLLOW(X) to FOLLOW(Z).
  - ► add \$ to FOLLOW(Z).
- $X \rightarrow YZ$, so add every non-$\varepsilon$ symbol in FIRST(Z) to FOLLOW(Y).
  - ► add **c** to FOLLOW(Y).
- $X \rightarrow YZ$ and $\varepsilon \in$ FIRST(Z), so add everything in FOLLOW(X) to FOLLOW(Y).
  - ► add \$ to FOLLOW(Y).

# Shift-reduce Parsing

- An instance of bottom-up parsing
- *Basic idea*: repeat
    1. in the string being processed, find a substring α such that $A \rightarrow α$ is a production;
    2. replace the substring α by *A* (i.e., reverse a derivation step).

    until we get the start symbol.
- *Technical issues*: Figuring out
    1. which substring to replace; and
    2. which production to reduce with.

# Shift-reduce Parsing: Example

*Grammar*:   $S \rightarrow$ **a**$AB$**e**

$A \rightarrow A$**bc**  |  **b**

$B \rightarrow$ **d**

*Input*:   **abbcde**        (using $A \rightarrow$ **b**)

$\Rightarrow$   **a**$A$**bcde**       (using $A \rightarrow A$**bc**)

$\Rightarrow$   **a**$A$**de**        (using $B \rightarrow$ **d**)

$\Rightarrow$   **a**$AB$**e**        (using $S \rightarrow$ **a**$AB$**e**)

$\Rightarrow$   $S$

## Shift-Reduce Parsing: cont'd

- Need to choose reductions carefully:

  **a**b**bcde** $\Rightarrow$ **a***A***bcde** $\Rightarrow$ **a***A***bc***B***e** $\Rightarrow$ ...

  doesn't work.

- A *handle* of a string *s* is a substring $\alpha$ s.t.:
  - $\alpha$ matches the RHS of a rule $A \rightarrow \alpha$; and
  - replacing $\alpha$ by $A$ (the LHS of the rule) represents a step in the <u>reverse</u> of a <u>rightmost derivation</u> of *s*.
- For shift-reduce parsing, reduce only handles.

## Shift-reduce Parsing: Implementation

- *Data Structures*:
  - a stack, its bottom marked by '**$**'. Initially empty.
  - the input string, its right end marked by '**$**'. Initially $w$.
- *Actions*:

  **repeat**
  1. *Shift* some ($\geq 0$) symbols from the input string onto the stack, until a handle $\beta$ appears on top of the stack.
  2. *Reduce* $\beta$ to the LHS of the appropriate production.

  **until** ready to accept.
  - *Acceptance*: when input is empty and stack contains only the start symbol.

# Example

| _Stack_ (→) | _Input_ | _Action_ |
|:---:|:---:|:---|
| **$** | **abbcde$** | shift |
| **$a** | **bbcde$** | shift |
| **$ab** | **bcde$** | reduce: $A \rightarrow$ **b** |
| **$a**_A_ | **bcde$** | shift |
| **$a**_A_**b** | **cde$** | shift |
| **$a**_A_**bc** | **de$** | reduce: $A \rightarrow A$**bc** |
| **$a**_A_ | **de$** | shift |
| **$a**_A_**d** | **e$** | reduce: $B \rightarrow$ **d** |
| **$a**_A__B_ | **e$** | shift |
| **$a**_A__B_**e** | **$** | reduce: $S \rightarrow$ **a**_AB_**e** |
| **$**_S_ | **$** | accept |

_Grammar_ :

$S \rightarrow$ **a**_AB_**e**

$A \rightarrow A$**bc** | **b**

$B \rightarrow$ **d**

---

# Conflicts

- Can't decide whether to shift or to reduce — both seem OK (*"shift-reduce conflict"*).

  _Example:_ $S \rightarrow$ **if** $E$ **then** $S$ | **if** $E$ **then** $S$ **else** $S$ | ...

- Can't decide which production to reduce with — several may fit (*"reduce-reduce conflict"*).

  _Example: Stmt_ $\rightarrow$ **id (** _args_ **)** | _Expr_
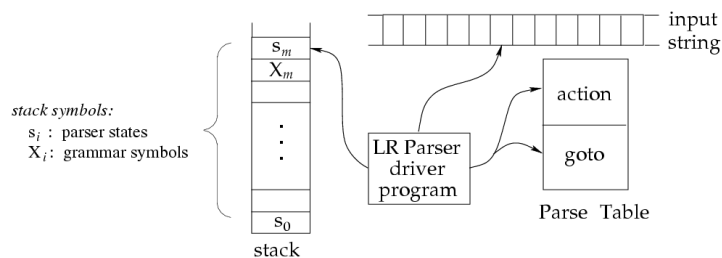
          _Expr_ $\rightarrow$ **id (** _args_ **)**

# LR Parsing

- A kind of shift-reduce parsing. An LR($k$) parser:
  - scans the input L-to-R;
  - produces a Rightmost derivation (in reverse); and
  - uses $k$ tokens of lookahead.
- Advantages:
  - very general and flexible, and handles a wide class of grammars;
  - efficiently implementable.
- Disadvantages:
  - difficult to implement by hand (use tools such as *yacc* or *bison*).

# LR Parsing: Schematic



*stack symbols:*
$s_i$ : parser states
$X_i$ : grammar symbols

- The driver program is the same for all LR parsers (SLR(1), LALR(1), LR(1), …). Only the parse table changes.
- Different LR parsing algorithms involve different tradeoffs between parsing power, parse table size.

# LR Parsing: the parser stack

- The parser stack holds strings of the form

  $s_0 X_1 s_1 X_2 s_2 \ ... \ X_m s_m$   ($s_m$ is on top)

  where $s_i$ are parser states and $X_i$ are grammar symbols.

  (**Note**: the $X_i$ and $s_i$ always come in pairs, with the state component $s_i$ on top.)

- A parser configuration is a pair

  ⟨stack contents, unexpended input⟩

# LR Parsing: Roadmap

- LR parsing algorithm:
  - parse table structure
  - parsing actions

- Parse table construction:
  - viable prefix automaton
  - parse table construction from this automaton
  - improving parsing power: different LR parsing algorithms

## LR Parse Tables

- The parse table has two parts: the **<u>action</u>** function and the **<u>goto</u>** function.

- At each point, the parser's next move is given by **<u>action</u>**$[s_m, \mathbf{a}_i]$, where:
  - $s_m$ is the state on top of the parser stack, and
  - $\mathbf{a}_i$ the next input token.

- The **<u>goto</u>** function is used only during *reduce* moves.

## LR Parser Actions: shift

- Suppose:
  - the parser configuration is $\langle s_0\, X_1 s_1 \ldots X_m s_m, \;\; \mathbf{a}_i \ldots \mathbf{a}_n \rangle$, and
  - **<u>action</u>**$[s_m, \mathbf{a}_i] =$ '*<u>shift</u>* $s_n$'.
- Effects of shift move:
  1. push the next input symbol $\mathbf{a}_i$; and
  2. push the state $s_n$

- New configuration: $\langle s_0\, X_1 s_1 \ldots X_m s_m\, \mathbf{a}_i\, s_n, \;\; \mathbf{a}_{i+1} \ldots \mathbf{a}_n \rangle$

# LR Parser Actions: reduce

- Suppose:
  - the parser configuration is $\langle s_0 X_1 s_1 ... X_m s_m, \ \mathbf{a}_i ... \mathbf{a}_n \rangle$, and
  - **action**$[s_m, \mathbf{a}_i]$ = '*reduce* $A \to \beta$'.
- Effects of reduce move:
  1. pop $n$ states and $n$ grammar symbols off the stack ($2n$ symbols total), where $n = |\beta|$.
  2. suppose the (newly uncovered) state on top of the stack is $t$, and **goto**$[t, A] = u$.
  3. push $A$, then $u$.

- New configuration: $\langle s_0 X_1 s_1 ... X_{m-n} s_{m-n} \ A \ u, \ \mathbf{a}_i ... \mathbf{a}_n \rangle$

---

# LR Parsing Algorithm

1. set `ip` to the start of the input string $w$**$**.
2. **while** TRUE **do**:
   1. let $s$ = state on top of parser stack, **a** = input symbol pointed at by `ip`.
   2. if **action**$[s,\mathbf{a}]$ == '*shift* $t$' then: (*i*) push the input symbol **a** on the stack, then the state $t$; (*ii*) advance `ip`.
   3. if **action**$[s,\mathbf{a}]$ == '*reduce* $A \to \beta$' then: (*i*) pop $2*|\beta|$ symbols off the stack; (*ii*) suppose $t$ is the state that now gets uncovered on the stack; (*iii*) push the LHS grammar symbol $A$ and the state $u$ = **goto**$[A, t]$.
   4. if **action**$[s,\mathbf{a}]$ == '*accept*' then accept;
   5. else signal a syntax error.

# LR parsing: Viable Prefixes

- **_Goal_**: to be able to identify handles, and so produce a rightmost derivation in reverse.
- Given a configuration $\langle s_0\, X_1 s_1 \ldots X_m s_m, \mathbf{a}_i \ldots \mathbf{a}_n \rangle$:
  - $X_1\, X_2 \ldots X_m\, \mathbf{a}_i \ldots \mathbf{a}_n$ is obtainable on a rightmost derivation.
  - $X_1\, X_2 \ldots X_m$ is called a _viable prefix_.
- The set of viable prefixes of a grammar are recognizable using a finite automaton.
  This automaton is used to recognize handles.

# Viable Prefix Automata

- An LR(0) item of a grammar G is a production of G with a dot "•" somewhere in the RHS.
  - _Example_: The rule $A \rightarrow \mathbf{a}\, A\, \mathbf{b}$ gives these LR(0) items:
    - $A \rightarrow \bullet\, \mathbf{a}\, A\, \mathbf{b}$
    - $A \rightarrow \mathbf{a}\, \bullet\, A\, \mathbf{b}$
    - $A \rightarrow \mathbf{a}\, A\, \bullet\, \mathbf{b}$
    - $A \rightarrow \mathbf{a}\, A\, \mathbf{b}\, \bullet$
- Intuition: '$A \rightarrow \alpha \bullet \beta$' denotes that:
  - we've seen something derivable from $\alpha$; and
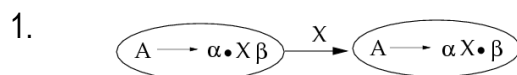  - it would be legal to see something derivable from $\beta$ at this point.

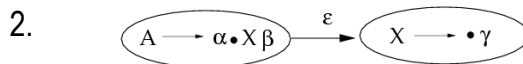# Overall Approach

Given a grammar G with start symbol $S$:

- Construct the augmented grammar by adding a new start symbol $S'$ and a new production $S' \rightarrow S$.

- Construct a finite state automaton whose start state is labeled by the LR(0) item $S' \rightarrow \bullet S$.

- Use this automaton to construct the parsing table.

# Viable Prefix NFA for LR(0) items

- Each state is labeled by an LR(0) item. The initial state is labeled $S' \rightarrow \bullet S$.

- Transitions:

    1. 
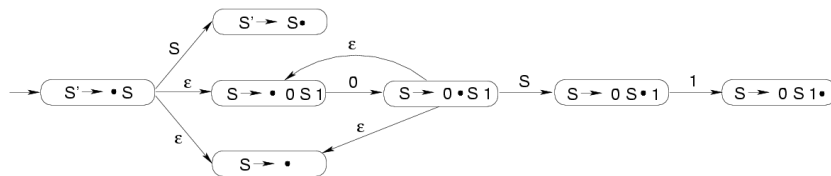    
    
    
    where *X* is a terminal or nonterminal.

    2. 
    
    
    
    where *X* is a nonterminal, and $X \rightarrow \gamma$ is a production.

# Viable Prefix NFA: Example

*Grammar* :

$S \rightarrow 0\,S\,1$

$S \rightarrow \varepsilon$

---

# Viable Prefix NFA $\Rightarrow$ DFA

- Given a set of LR(0) items *I*, the set *closure*(*I*) is constructed as follows:

  **repeat**
  1. add every item in *I* to *closure*(*I*);
  2. if $A \rightarrow \alpha \bullet B\beta \in$ *closure*(I) and *B* is a nonterminal, then for each production $B \rightarrow \gamma$, add the item $B \rightarrow \bullet\, \gamma$ to *closure*(*I*).

  **until** no new items can be added to *closure*(*I*).

- *Intuition*:

  $A \rightarrow \alpha \bullet B\beta \in$ *closure*(I) means something derivable from B$\beta$ is legal at this point. This means that something derivable from B (and thus $\gamma$) is also legal.
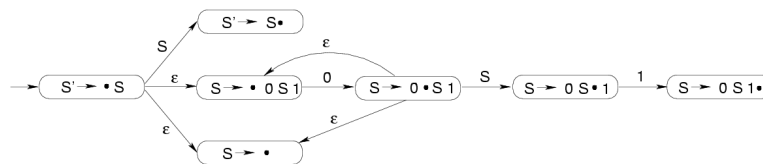
# Viable Prefix NFA $\Rightarrow$ DFA (cont'd)

- Given a set of LR(0) items *I*, the set *goto*(*I,X*) is defined as

  *goto*(*I*, *X*) = *closure*({ $A \to \alpha X \bullet \beta \mid A \to \alpha \bullet X \beta \in I$ })

- ***Intuition***:
  - if $A \to \alpha \bullet X \beta \in I$ then (*a*) we've seen something derivable from $\alpha$; and (*b*) something derivable from $X\beta$ would be legal at this point.
  - Suppose we now see something derivable from *X*.
    The parser should "go to" a state where (*a*) we've seen something derivable from $\alpha X$; and (*b*) something derivable from $\beta$ would be legal.

---

# Example



- Let $I_0$ = {S′ → •S}.
- $I_1$ = closure($I_0$) = { S′ → •S,                    /* from $I_0$ */

  S → • **0** S **1**, S → • }
- goto ($I_1$, **0**) = closure( { S → **0** • S **1** } )

  = {S → **0** • S **1**, S → • **0** S **1**, S → • }

# Viable Prefix DFA for LR(0) Items

1. Given a grammar G with start symbol $S$, construct the augmented grammar with new start symbol $S'$ and new production $S' \rightarrow S$.

2. $C = \{ \underline{closure}(\{ S' \rightarrow \bullet S \}) \}$;    // $C$ = a <u>set of sets</u> of items = set of parser states

3. **repeat** {
     **for** each set of items $I \in C$ {
       **for** each grammar symbol $X$ {
         **if** ( $\underline{goto}(I,X) \neq \varnothing$ && $\underline{goto}(I,X) \notin C$ ) {      // new state
           add $\underline{goto}(I,X)$ to $C$;
         }
       }
     }
   } **until** no change to $C$;

4. **return** $C$.

# SLR(1) Parse Table Construction I

Given a grammar $G$ with start symbol $S$:

- Construct the augmented grammar $G'$ with start symbol $S'$.
- Construct the set of states $\{I_0, I_1, \ldots, I_n\}$ for the Viable Prefix DFA for the augmented grammar $G'$.
- Each DFA state $I_i$ corresponds to a parser state $s_i$.
- The initial parser state $s_0$ coresponds to the DFA state $I_0$ obtained from the item $S' \rightarrow \bullet S$.
- The parser actions in state $s_i$ are defined by the items in the DFA state $I_i$.

# SLR(1) Parse Table Construction II

Parsing action for parser state $s_i$:

- action table entries:
  - if DFA state $I_i$ contains an item $A \rightarrow \alpha \bullet \mathbf{a} \beta$ where $\mathbf{a}$ is a terminal, and $\underline{goto}(I_i, \mathbf{a}) = I_j$ :  set **action**[$i$, $\mathbf{a}$] = *shift j*.
  - if DFA state $I_i$ contains an item $A \rightarrow \alpha \bullet$,  where $A \neq S'$: for each $\mathbf{b} \in$ FOLLOW($A$), set **action**[$i$, $\mathbf{b}$] = *reduce $A \rightarrow \alpha$.*
  - if state $I_i$ contains the item $S' \rightarrow S \bullet$: set **action**[$i$, **$**] = *accept*.
- goto table entries:
  - for each nonterminal $A$, if $\underline{goto}(I_i, A) = I_j$, then **goto**[$i$, $A$] = $j$.
- any entry not defined by these steps is an *error state*.
- *if any state has multiple entries, the grammar is not SLR(1).*

---

# SLR(1) Shortcomings

- SLR(1) parsing uses *reduce* actions too liberally.  Because of this it fails on many reasonable grammars.

- <u>Example</u> *(simple pointer assignments):*

  $S \rightarrow R$  |  $L = R$

  $L \rightarrow$ *R  |  **id**

  $R \rightarrow L$

  The SLR parse table has a state { $S \rightarrow L \bullet = R$, $R \rightarrow L \bullet$ }, and FOLLOW(L) = { =, $ }.

  $\Rightarrow$ *shift-reduce conflict.*

## Improving LR Parsing

- SLR(1) parsing weaknesses can be addressed by incorporating *lookahead* into the LR items in parser states.

  The lookahead makes it possible to remove some "spurious" reduce actions in the parse table.

  *The LALR(1) parsers produced by **bison** and **yacc** incorporate such lookahead items.*

- This improves parsing power, but at the cost of larger parse tables.

## Error Handling

Possible reactions to lexical and syntax errors:

- *ignore the error.*  Unacceptable!
- *crash, or quit, on first error.*  Unacceptable!
- *continue to process the input.*  No code generation.
- *attempt to repair the error:* transform an erroneous program into a similar but legal input.
- *attempt to correct the error:* try to guess what the programmer meant.  Not worthwhile.

# Error Reporting

- Error messages should refer to the source program.

  prefer "line 11: X redefined" to "conflict in hash bucket 53"

- Error messages should, as far as possible, indicate the location and nature of the error.

  avoid "syntax error" or "illegal character"

- Error messages should be specific.

  prefer "x not declared in function foo" to "missing declaration"

- They should not be redundant.

---

# Error Recovery

- *Lexical errors*: pass the illegal character to the parser and let it deal with the error.

- *Syntax errors*: "*panic mode* error recovery"
  - ***Essential idea***: *skip part of the input and pretend as though we saw something legal, then hope to be able to continue.*
  - Pop the stack until we find a state $s$ such that **goto**[$s,A$] is defined for some nonterminal $A$.
  - discard input tokens until we find some token **a** that can legitimately follow $A$ (i.e., **a** $\in$ FOLLOW($A$)).
  - push the state **goto**[$s,A$] and continue parsing.