

CS415 Compilers

Syntax Analysis Part 5

These slides are based on slides copyrighted by
Keith Cooper, Ken Kennedy & Linda Torczon at Rice
University

Roadmap for the remainder of the course

- Fourth homework:
Due Monday, March 28
- Project #2 - Bottom-up parser and compiler
Will be posted next week, due April 13 (tentative)
- Project #3 - Peephole optimizer for ILOC
Will be posted April 13, due May 2 (tentative)
- Second midterm on Wednesday, April 6 (60 minutes in class)
- Final exam on May 10 (60 minutes at assigned location)
- At least 3 more homeworks

Bottom-up Parsing (Syntax Analysis)

EAC Chapters 3.4

```
stack.push(INVALID); stack.push( $s_0$ );
not_found = true;
token = scanner.next_token();
do while (not_found) {
    s = stack.top();
    if ( ACTION[s,token] == "reduce  $A \rightarrow \beta$ " ) then {
        stack.popnum( $2 * |\beta|$ ); // pop  $2 * |\beta|$  symbols
        s = stack.top();
        stack.push(A);
        stack.push(GOTO[s,A]);
    }
    else if ( ACTION[s,token] == "shift s" ) then {
        stack.push(token); stack.push(s);
        token ← scanner.next_token();
    }
    else if ( ACTION[s,token] == "accept"
              & token == EOF )
        then not_found = false;
    else report a syntax error and recover;
}
report success;
```

The skeleton parser

- uses ACTION & GOTO tables
- does $|words|$ shifts
- does $|derivation|$ reductions
- does 1 accept

How do we generate the ACTION and GOTO tables?

- Use the grammar to build a model of the DFA
- Use the model to build ACTION & GOTO tables
- If construction succeeds, the grammar is LR(1)

*Terminal or
non-terminal*

The Big Picture

- Model the state of the parser
- Use two functions *goto(s, X)* and *closure(s)*
 - *goto()* is analogous to *move()* in the subset construction
 - *closure()* adds information to round out a state
- Build up the states and transition functions of the DFA
- Use this information to fill in the ACTION and GOTO tables

The LR(1) table construction algorithm uses LR(1) items to represent valid configurations of an LR(1) parser

An LR(k) item is a pair $[P, \delta]$, where

P is a production $A \rightarrow \beta$ with a \cdot at some position in the *rhs*

δ is a lookahead string of length $\leq k$ (words or EOF)

The \cdot in an item indicates the position of the top of the stack

LR(1):

$[A \rightarrow \cdot \beta \gamma, \underline{a}]$ means that the input seen so far is consistent with the use of $A \rightarrow \beta \gamma$ immediately after the symbol on top of the stack

$[A \rightarrow \beta \cdot \gamma, \underline{a}]$ means that the input seen so far is consistent with the use of $A \rightarrow \beta \gamma$ at this point in the parse, and that the parser has already recognized β .

$[A \rightarrow \beta \gamma \cdot, \underline{a}]$ means that the parser has seen $\beta \gamma$, and that a lookahead symbol of \underline{a} is consistent with reducing to A .

The production $A \rightarrow \beta$, where $\beta = B_1 B_2 B_3$ with lookahead \underline{a} , can give rise to 4 items

$[A \rightarrow \cdot B_1 B_2 B_3, \underline{a}]$, $[A \rightarrow B_1 \cdot B_2 B_3, \underline{a}]$, $[A \rightarrow B_1 B_2 \cdot B_3, \underline{a}]$, & $[A \rightarrow B_1 B_2 B_3 \cdot, \underline{a}]$

The set of LR(1) items for a grammar is **finite**

What's the point of all these lookahead symbols?

- Carry them along to choose the correct reduction, *if there is a choice*
 - Lookaheads are bookkeeping, unless item has \cdot at the right end
 - Has no direct use in $[A \rightarrow \beta \cdot \gamma, \underline{a}]$
 - In $[A \rightarrow \beta \cdot, \underline{a}]$, a lookahead of \underline{a} implies a reduction by $A \rightarrow \beta$
 - For $\{ [A \rightarrow \beta \cdot, \underline{a}], [B \rightarrow \gamma \cdot c, \underline{b}] \}$, $\underline{a} \Rightarrow \text{reduce to } A$; $\underline{c} \Rightarrow \text{shift}$
- ⇒ **Limited right context** is enough to pick the actions (unique, i.e., deterministic choice)

High-level overview

1 Build the **canonical collection of sets of LR(1) Items, I** a Begin in an appropriate state, s_0

- ◆ Assume: $S' \rightarrow S$, and S' is unique start symbol that does not occur on any RHS of a production (extended CFG - ECFG)
- ◆ $[S' \rightarrow \cdot S, \underline{\text{EOF}}]$, along with any equivalent items
- ◆ Derive equivalent items as $\text{closure}(s_0)$

b Repeatedly compute, for each s_k , and each X , $\text{goto}(s_k, X)$

- ◆ If the set is not already in the collection, add it
- ◆ Record all the transitions created by $\text{goto}()$

This eventually reaches a fixed point

2 Fill in the table from the collection of sets of LR(1) items

The canonical collection completely encodes the transition diagram for the handle-finding DFA

$Closure(s)$ adds all the items implied by items already in s

- Any item $[A \rightarrow \beta \bullet B \delta, \underline{a}]$ implies $[B \rightarrow \bullet \tau, \underline{x}]$ for each production with B on the *lhs*, and each $x \in FIRST(\delta \underline{a})$

The algorithm

```

Closure( s )
  while ( s is still changing )
     $\forall$  items  $[A \rightarrow \beta \bullet B \delta, \underline{a}] \in s$ 
       $\forall$  productions  $B \rightarrow \tau \in P$ 
         $\forall \underline{b} \in FIRST(\delta \underline{a})$  //  $\delta$  might be  $\epsilon$ 
          if  $[B \rightarrow \bullet \tau, \underline{b}] \notin s$ 
            then add  $[B \rightarrow \bullet \tau, \underline{b}]$  to  $s$ 
  
```

- Classic fixed-point method
 - Halts because $s \subset ITEMS$
- Closure "fills out" a state*

$Goto(s, x)$ computes the state that the parser would reach if it recognized an X while in state s

- $Goto(\{ [A \rightarrow \beta \bullet X \delta, \underline{a}] \}, X)$ produces $[A \rightarrow \beta X \bullet \delta, \underline{a}]$ (easy part)
- Should also includes $closure([A \rightarrow \beta X \bullet \delta, \underline{a}])$ (fill out the state)

The algorithm

```
 $Goto(s, X)$   
   $new \leftarrow \emptyset$   
   $\forall \text{ items } [A \rightarrow \beta \bullet X \delta, \underline{a}] \in s$   
     $new \leftarrow new \cup [A \rightarrow \beta X \bullet \delta, \underline{a}]$   
  return  $closure(new)$ 
```

- Not a fixed-point method!
- Straightforward computation
- Uses $closure()$

$Goto()$ moves forward

Start from $s_0 = \text{closure}([S' \rightarrow S, \underline{\text{EOF}}])$

Repeatedly construct new states, until all are found

The algorithm

```
 $cc_0 \leftarrow \text{closure}([S' \rightarrow \bullet S, \underline{\text{EOF}}])$   
 $CC \leftarrow \{ cc_0 \}$   
while ( new sets are still being added to CC )  
  for each unmarked set  $cc_j \in CC$   
    mark  $cc_j$  as processed  
    for each  $x$  following a  $\bullet$  in an item in  $cc_j$   
       $temp \leftarrow \text{goto}(cc_j, x)$   
      if  $temp \notin CC$   
        then  $CC \leftarrow CC \cup \{ temp \}$   
      record transitions from  $cc_j$  to temp on  $x$ 
```

- Fixed-point computation
(worklist version)
- Loop adds to CC
- $CC \subseteq 2^{\text{ITEMS}}$,
so CC is finite

Simplified, right recursive expression grammar

1: $Goal \rightarrow Expr$
2: $Expr \rightarrow Term - Expr$
3: $Expr \rightarrow Term$
4: $Term \rightarrow Factor * Term$
5: $Term \rightarrow Factor$
6: $Factor \rightarrow \underline{ident}$

Symbol	FIRST
Goal	{ <u>ident</u> }
Expr	{ <u>ident</u> }
Term	{ <u>ident</u> }
Factor	{ <u>ident</u> }
-	{ - }
*	{ * }
<u>ident</u>	{ <u>ident</u> }

1: $Goal \rightarrow Expr$
 2: $Expr \rightarrow Term - Expr$
 3: $Expr \rightarrow Term$
 4: $Term \rightarrow Factor * Term$
 5: $Term \rightarrow Factor$
 6: $Factor \rightarrow \underline{ident}$

Symbol	FIRST
Goal	{ <u>ident</u> }
Expr	{ <u>ident</u> }
Term	{ <u>ident</u> }
Factor	{ <u>ident</u> }
-	{ - }
*	{ * }
<u>ident</u>	{ <u>ident</u> }

$s_0 \leftarrow \text{closure}([Goal \rightarrow \cdot Expr, EOF]) =$
 $\{ [Goal \rightarrow \cdot Expr, EOF], [Expr \rightarrow \cdot Term - Expr, EOF],$
 $[Expr \rightarrow \cdot Term, EOF], [Term \rightarrow \cdot Factor * Term, EOF],$
 $[Term \rightarrow \cdot Factor * Term, -], [Term \rightarrow \cdot Factor, EOF],$
 $[Term \rightarrow \cdot Factor, -], [Factor \rightarrow \cdot \underline{ident}, EOF],$
 $[Factor \rightarrow \cdot \underline{ident}, -], [Factor \rightarrow \cdot \underline{ident}, *] \}$

Iteration 1

 $s_1 \leftarrow \text{goto}(s_0, \text{Expr})$ $s_2 \leftarrow \text{goto}(s_0, \text{Term})$ $s_3 \leftarrow \text{goto}(s_0, \text{Factor})$ $s_4 \leftarrow \text{goto}(s_0, \underline{\text{ident}})$

Iteration 2

 $s_5 \leftarrow \text{goto}(s_2, -)$ $s_6 \leftarrow \text{goto}(s_3, *)$

Iteration 3

 $s_7 \leftarrow \text{goto}(s_5, \text{Expr})$ $s_8 \leftarrow \text{goto}(s_6, \text{Term})$

$$S_0 : \{ [Goal \rightarrow \cdot Expr, EOF], [Expr \rightarrow \cdot Term - Expr, EOF], \\ [Expr \rightarrow \cdot Term, EOF], [Term \rightarrow \cdot Factor * Term, EOF], \\ [Term \rightarrow \cdot Factor * Term, -], [Term \rightarrow \cdot Factor, EOF], \\ [Term \rightarrow \cdot Factor, -], [Factor \rightarrow \cdot \underline{ident}, EOF], \\ [Factor \rightarrow \cdot \underline{ident}, -], [Factor \rightarrow \cdot \underline{ident}, *] \}$$

$$S_1 : \{ [Goal \rightarrow Expr \cdot, EOF] \}$$

$$S_2 : \{ [Expr \rightarrow Term \cdot - Expr, EOF], [Expr \rightarrow Term \cdot, EOF] \}$$

$$S_3 : \{ [Term \rightarrow Factor \cdot * Term, EOF], [Term \rightarrow Factor \cdot * Term, -], \\ [Term \rightarrow Factor \cdot, EOF], [Term \rightarrow Factor \cdot, -] \}$$

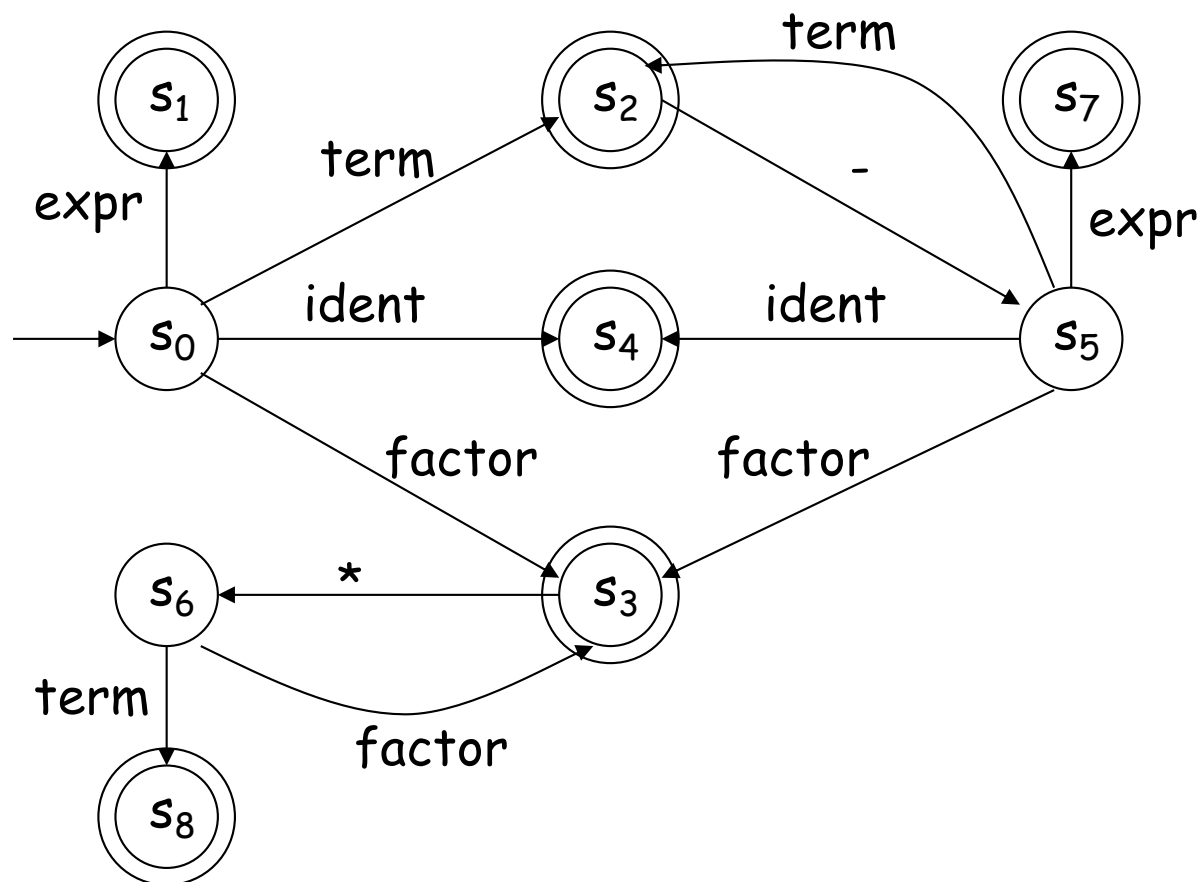
$$S_4 : \{ [Factor \rightarrow \underline{ident} \cdot, EOF], [Factor \rightarrow \underline{ident} \cdot, -], [Factor \rightarrow \underline{ident} \cdot, *] \}$$

$$S_5 : \{ [Expr \rightarrow Term - \cdot Expr, EOF], [Expr \rightarrow \cdot Term - Expr, EOF], \\ [Expr \rightarrow \cdot Term, EOF], [Term \rightarrow \cdot Factor * Term, -], \\ [Term \rightarrow \cdot Factor, -], [Term \rightarrow \cdot Factor * Term, EOF], \\ [Term \rightarrow \cdot Factor, EOF], [Factor \rightarrow \cdot \underline{ident}, *], \\ [Factor \rightarrow \cdot \underline{ident}, -], [Factor \rightarrow \cdot \underline{ident}, EOF] \}$$

$$S_6 : \{ [Term \rightarrow Factor * \cdot Term, EOF], [Term \rightarrow Factor * \cdot Term, -], \\ [Term \rightarrow \cdot Factor * Term, EOF], [Term \rightarrow \cdot Factor * Term, -], \\ [Term \rightarrow \cdot Factor, EOF], [Term \rightarrow \cdot Factor, -], \\ [Factor \rightarrow \cdot \underline{ident}, EOF], [Factor \rightarrow \cdot \underline{ident}, -], [Factor \rightarrow \cdot \underline{ident}, *] \}$$
$$S_7 : \{ [Expr \rightarrow Term - Expr \cdot, EOF] \}$$
$$S_8 : \{ [Term \rightarrow Factor * Term \cdot, EOF], [Term \rightarrow Factor * Term \cdot, -] \}$$

The Goto Relationship (*from the construction*)

State	Expr	Term	Factor	-	*	<u>ident</u>
0	1	2	3			4
1						
2				5		
3					6	
4						
5	7	2	3			4
6		8	3			4
7						
8						



x is the number of the state for s_x

The algorithm

```

 $\forall$  set  $s_x \in S$ 
   $\forall$  item  $i \in s_x$ 
    if  $i$  is  $[A \rightarrow \beta \cdot \underline{a}d, \underline{b}]$  and  $\text{goto}(s_x, \underline{a}) = s_k, \underline{a} \in T$ 
      then  $\text{ACTION}[x, \underline{a}] \leftarrow \text{"shift } k\text{"}$ 
    else if  $i$  is  $[S' \rightarrow S \cdot, \text{EOF}]$ 
      then  $\text{ACTION}[x, \text{EOF}] \leftarrow \text{"accept"}$ 
    else if  $i$  is  $[A \rightarrow \beta \cdot, \underline{a}]$ 
      then  $\text{ACTION}[x, \underline{a}] \leftarrow \text{"reduce } A \rightarrow \beta\text{"}$ 
   $\forall n \in NT$ 
    if  $\text{goto}(s_x, n) = s_k$ 
      then  $\text{GOTO}[x, n] \leftarrow k$ 

```

Many items
generate no
table entry

e.g., $[A \rightarrow \beta \cdot B \alpha, \underline{a}]$
does not, but
closure ensures
that all the rhs'
for B are in s_x

The algorithm produces the following table

	ACTION				GOTO		
	<u>ident</u>	-	*	EOF	<i>Expr</i>	<i>Term</i>	<i>Factor</i>
0	s 4				1	2	3
1				acc			
2		s 5		r 3			
3		r 5	s 6	r 5			
4		r 6	r 6	r 6			
5	s 4				7	2	3
6	s 4					8	3
7				r 2			
8		r 4		r 4			

Plugs into the skeleton LR(1) parser

What if set s contains $[A \rightarrow \beta \cdot \underline{a} \gamma, \underline{b}]$ and $[B \rightarrow \beta \cdot, \underline{a}]$?

- First item generates "shift", second generates "reduce"
- Both define $\text{ACTION}[s, \underline{a}]$ — cannot do both actions
- This is a fundamental ambiguity, called a *shift/reduce error*
- Modify the grammar to eliminate it *(if-then-else)*
- Shifting will often resolve it correctly

EaC includes a
worked example

What if set s contains $[A \rightarrow \gamma' \cdot, \underline{a}]$ and $[B \rightarrow \gamma \cdot, \underline{a}]$?

- Each generates "reduce", but with a different production
- Both define $\text{ACTION}[s, \underline{a}]$ — cannot do both reductions
- This fundamental ambiguity is called a *reduce/reduce error*
- Modify the grammar to eliminate it

In either case, the grammar is not LR(1)

Closure(s) adds all the items implied by items already in s

- Any item $[A \rightarrow \beta \bullet B \delta]$ implies $[B \rightarrow \bullet \tau]$ for each production with B on the *lhs*

The algorithm

```
Closure(s)
  while (s is still changing)
     $\forall$  items  $[A \rightarrow \beta \bullet B \delta] \in s$ 
       $\forall$  productions  $B \rightarrow \tau \in P$ 

        if  $[B \rightarrow \bullet \tau] \notin s$ 
          then add  $[B \rightarrow \bullet \tau]$  to s
```

- Classic fixed-point method
 - Halts because $s \subset \text{ITEMS}$
- Closure "fills out" a state*

Context-Sensitive Analysis