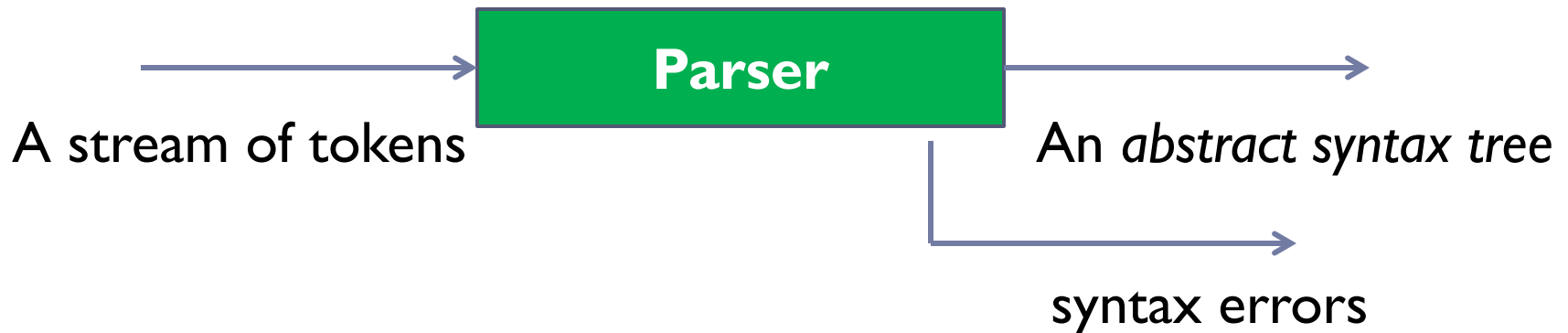# Compiler Techniques

## 3.    Syntax Analysis

Huang Shell Ying

# Overview

▸ The syntax analyzer (parser) checks that the program is **syntactically well-formed** and transforms it from a sequence of tokens into an *abstract syntax tree (AST)*.
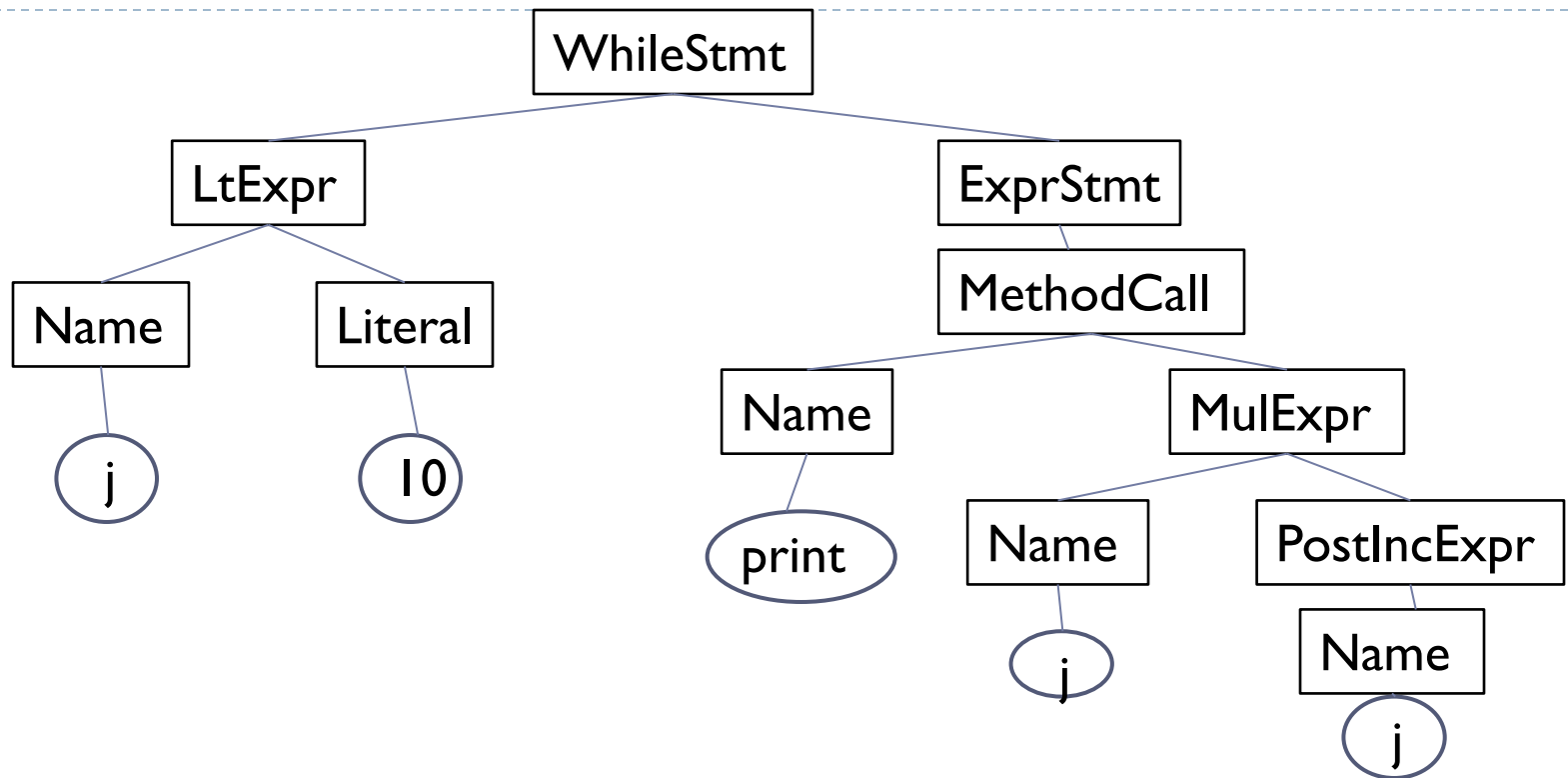


A stream of tokens → **Parser** → An *abstract syntax tree*

syntax errors

For example,

**while ( j < 10 ) print ( j * ( j ++ ) ) ;**

| WHILE | LPAR | ID | LT | LIT | RPAR | ID | LPAR | ID | MUL | LPAR | … | back |
|-------|------|----|----|-----|------|----|------|----|----|------|----|------|

# AST for while ( j < 10 ) print ( j * ( j ++ ) ) ;



- The abstract syntax tree (AST) captures the **structure of the program**.

- Later stages of the compiler use the AST for **semantic analysis** and **code generation**.

Syntax Analysis   CZ3007

# Overview

▸ To check whether a program is well-formed, we need a tool – The grammar of a language.

▸ A language's grammar serves as a concise definition of how meaningful sentences in a language can be constructed.

▸ In this chapter, we will learn

   ▸ Context-free grammar

   ▸ The parsing problem

      ▪ Top-down parsing

      ▪ Bottom-up parsing

# Context-free Grammars

▸ Programming language syntax can be described by a **context-free grammar** (can we use regular expressions?).

▸ A context-free grammar $G = (N, T, P, S)$ consists of four components:

1. A finite set $T$ of **terminals** (token types)
2. A finite set $N$ of **nonterminals** such that $T \cap N = \varnothing$
3. A **start symbol** $S \in N$
4. A finite set $P$ of **rules** of the form $A \rightarrow s_1 \ldots s_n$ where $A \in N$, $n \geq 0$, and $\forall i \in \{1, \ldots, n\}, s_i \in T \cup N$.
   If $n = 0$, we write the rule as $A \rightarrow \lambda$.

# Example of a context free grammar

$Expr \rightarrow Expr$ plus $Term$
$\quad | \quad Expr$ minus $Term$
$\quad | \quad Term$

$Term \rightarrow Term$ mul $Factor$
$\quad | \quad Term$ div $Factor$
$\quad | \quad Factor$

$Factor \rightarrow$ number
$\quad | \quad$ id
$\quad | \quad$ lparen $Expr$ rparen

"|" is used to group multiple rules for the same nonterminal.

$Factor \rightarrow$ number
$Factor \rightarrow$ id
$Factor \rightarrow$ lparen $Expr$ rparen

back

# Notation Adopted

| Names Beginning With | Examples | Represent Symbols in |
|---|---|---|
| Upper case | A, B, C, Expr, Stmt | $N$ |
| Lower case and punctuation | a, b, c, if, then, plus | $T$ |
| $X, Y$ | $X_1, X_2$ | $N \cup T$ |
| Other Greek letters | $\alpha, \beta, \gamma$ | $(N \cup T)^*$ |

back

# Deriving Sentences

▸ The language defined by a grammar is <u>the set of all sentences</u> that can be derived from its start symbol.

▸ Example of one sentence derived using the rules on slide 6:

*Expr* ⇒ *Expr* plus *Term*
⇒ *Term* plus *Term*
⇒ *Factor* plus *Term*
⇒ *id* plus *Term*
⇒ *id* plus *Term mul Factor*

⇒ *id* plus *Factor mul Factor*
⇒ *id* plus *id mul Factor*
⇒ *id* plus *id mul id*

▸ Since there may be multiple rules for a nonterminal, derivation is non-deterministic: we may derive many different sentences from the same initial phrase. (try to derive 3 more sentences for the Expr language)

# Deriving Sentences

- If $A \rightarrow \beta$ is a rule of G and $\alpha, \gamma \in (N \cup T)^*$ are phrases of G, then $\alpha A \gamma \Rightarrow \alpha \beta \gamma$ is a **one-step derivation**

  E.g. *Expr* plus *Term* $\Rightarrow$ *Expr* minus *Term* plus *Term*

- $\alpha \Rightarrow +\beta$ means that $\beta$ is derived in *one or more* steps from $\alpha$,

  $\alpha \Rightarrow *\beta$ in *zero or more* steps

- The set of terminal strings derivable from *S* are the **set of all sentences** of the language, denoted *L(G)*

# Derivation Sequences

▸ If there are multiple nonterminals in a phrase, there is a choice as to which nonterminal should be expanded next.

**Term** $\Rightarrow$ **Term** mul **Term**

▸ In a *leftmost* derivation, we always expand the first nonterminal, in a *rightmost* one the last nonterminal.

▸ Ultimately, the derivation order does not matter: we can derive any sentence in $L(G)$ using any strategy.

▸ What is important, however, is which rule is applied at each nonterminal occurrence.

$$
\begin{aligned}
Term \quad &\rightarrow \quad Term \text{ mul } Factor \\
&| \quad\quad Term \text{ div } Factor \\
&| \quad\quad Factor
\end{aligned}
$$

# Parse Trees (this is not the AST)

- A *parse tree* **represents a derivation**, but abstracts away from the order of derivations.

- Every node in a parse tree is labelled with a symbol:
  - the **root** node is labelled with the **start** symbol;
  - **leaf** nodes are labelled with **terminal** symbols or $\lambda$;
  - **inner nodes** are labelled with **nonterminal** symbols.

- The labelling obeys the following requirement:

  A node labelled with $A$ which has children labelled $s_1 \ldots s_n$, if and only if there is a rule $A \rightarrow s_1 \ldots s_n$

# Example Parse Tree

Expr → Expr plus Term
    | Expr minus Term
    | Term
Term → Term mul Factor
    | Term div Factor
    | Factor
Factor → number
    | lparen Expr rparen



Sentence derived:

number plus lparen number plus number rparen mul number

Syntax Analysis    CZ3007

# Ambiguous Grammars

▸ Consider the grammar *G* defined by the following two rules:

Term  →  Term div Term                    // rule 1

        |    number                       // rule 2

**A parse tree for number div number div number**

**The string: 100/10/2**

# Ambiguous Grammars

**Another parse tree for**
**number div number div number**

**The string:  100/10/2**

# Ambiguous Grammars

- **A grammar that produces more than one leftmost derivation or more than one rightmost derivation for the same sentence are called ambiguous.**

- Ambiguous grammars are not useful for compilers.

- There is no algorithm that can check an arbitrary context free grammar for ambiguity.

# The Parsing Problem

| Formal Statement |
| --- |
| Given a context-free grammar $G = (T, N, S, P)$ and a sentence $w \in T^*$, decide whether or not $w \in L(G)$. |

▸ ***Top-down parsing***: generates a parse tree by starting at the **root** of the tree (the **start** symbol $S$), expanding the tree by applying rules in a depth-first manner.

▸ ***Bottom-up parsing***: generates a parse tree by starting at the tree's **leaves** (and $w$) and working towards its root. A node is inserted into the tree only after its children have been inserted.

# Top-down Parsing

▸ This parsing technique is known by a few names:

1. **Top-down**, because it begins with the grammar's start symbol and grows a parse tree from its root to its leaves.

2. **Predictive**, because it predicts at each step which grammar rule is to be used.

3. **LL(k)**, because it scans the input from left to right, producing a leftmost derivation, using k symbols of lookahead. We will consider only **LL(1)**.

4. **Recursive descent**, because it can be implemented by a collection of mutually recursive procedures.

# Recursive Descent Parsing

▸ In a recursive descent parser, for every nonterminal A there is a corresponding method parseA that can parse sentences derived from A.

| The grammar: | The corresponding methods: |
|---|---|
| S $\rightarrow$ A C | parseS {…} |
| C $\rightarrow$ c | parseC {…} |
| $\quad$ \| $\quad \lambda$ | |
| A $\rightarrow$ a B C d | parseA {…} |
| $\quad$ \| $\quad$ B Q | |
| B $\rightarrow$ b B | parseB {…} |
| $\quad$ \| $\quad \lambda$ | |
| Q $\rightarrow$ q | parseQ {…} |
| $\quad$ \| $\quad \lambda$ | |

Syntax Analysis　CZ3007

# Recursive Descent Parsing

▸ If there is more than one rule for A, parseA inspects the next input token(s) and choose a production rule among the rules for A to apply.

Slide 12

▸ The code for applying a production rule is obtained by processing the RHS of the rule, symbol by symbol

$$A \rightarrow X_1 X_2 \ldots X_m$$

Slide 21

▸ If the next symbol $X_i$ is a terminal $t$, confirms the next input token is $t$.

▸ If it is a nonterminal B, call the parsing function parseB.

▸ The code for applying a production rule $A \rightarrow \lambda$ will do nothing and simply return.

back

# Recursive Descent Parsing

▸ The parsing of the whole program starts from the parse method for the start symbol.

▸ Recursive descent parsers use one token of lookahead to determine which rule to use.

▸ Lookahead has to be unambiguous: there cannot be more than one rule (for the same nonterminal) whose RHS starts with the same token.

▸ A grammar that fulfills this condition is an *LL*(1) *grammar*.

▸ A language for which there exists an *LL*(1) grammar is an *LL*(1) *language*.

# Example

**parseS(ts)**
{  // ts is the input token stream
    if (ts.peek() $\in$ **predict**($p_1$))
            **parseA(ts);  parseC(ts);**
    else /* syntax error */        }

**parseA(ts)**
{  if (ts.peek() $\in$ **predict**($p_4$))
            **match(a);  parseB(ts);**
            **parseC(ts);  match(d);**
    else if (ts.peek() $\in$ **predict**($p_5$))
            **parseB(ts); parseQ(ts);**
    else /* syntax error */        }

/* **peek()** examines the next input token
    without advancing the input  */

$S \rightarrow A\ C$        $p_1$

$C \rightarrow c$            $p_2$

$\quad |\ \ \lambda$         $p_3$

$A \rightarrow \textbf{a B C d}$     $p_4$

$\quad |\ \ \textbf{B Q}$       $p_5$

$B \rightarrow b\ B$          $p_6$

$\quad |\ \ \lambda$          $p_7$

$Q \rightarrow q$           $p_8$

$\quad |\ \ \lambda$          $p_9$

/* **match()**
confirms a token */

# Computing predict(p)

Consider a production rule p: $X \rightarrow X_1 X_2 \ldots X_m$, $m \geq 0$.

- The set of tokens that **predicts** rule p includes
  - The set of possible tokens that are first produced in some derivation from $X_1 X_2 \ldots X_m$
    - The set of first tokens in $X_1$
    - If $X_1$ may be empty, the set of first tokens in $X_2$, and so on.
  - Those tokens that can follow X in some derivation from $X_1 X_2 \ldots X_m$
    - If $X_1 X_2 \ldots X_m$ may be empty, the first tokens that may follow X

| | |
|---|---|
| $S \rightarrow A\ C$ | $P_1$ |
| $C \rightarrow c$ | $P_2$ |
| $\mid \quad \lambda$ | $P_3$ |
| $A \rightarrow$ **a B C d** | $P_4$ |
| $\mid \quad$ **B Q** | $P_5$ |
| $B \rightarrow b\ B$ | $P_6$ |
| $\mid \quad \lambda$ | $P_7$ |
| $Q \rightarrow q$ | $P_8$ |
| $\mid \quad \lambda$ | $P_9$ |

# Computing predict(p)

Example:  p: $A \rightarrow B\ Q$

▸ The set of possible tokens that are first produced in some derivation from B Q for rule $A \rightarrow B\ Q$ is { b, q}

▸ Those tokens that can follow A in some derivation from B Q  is {c}

The set of tokens that predicts rule p:
  {b,  q, c}

| | | | |
|---|---|---|---|
| S | $\rightarrow$ | A C | $P_1$ |
| C | $\rightarrow$ | c | $P_2$ |
| | | $\lambda$ | $P_3$ |
| A | $\rightarrow$ | **a B C d** | $P_4$ |
| | | **B Q** | $P_5$ |
| B | $\rightarrow$ | b B | $P_6$ |
| | | $\lambda$ | $P_7$ |
| Q | $\rightarrow$ | q | $P_8$ |
| | | $\lambda$ | $P_9$ |

# Computing predict(p)

- To compute the set of tokens that predict rule p, we need to know whether or not

  - a nonterminal can derive empty
  - The RHS of a rule can derive empty

- Two boolean arrays are used:

  - symbolDerivesEmpty[$X$] for $X \in N$
  - ruleDerivesEmpty[p] for $p \in P$

| | |
|---|---|
| $S \rightarrow A\ C$ | $P_1$ |
| $C \rightarrow c$ | $P_2$ |
| $\quad\| \quad \lambda$ | $P_3$ |
| $A \rightarrow \textbf{a B C d}$ | $P_4$ |
| $\quad\| \quad \textbf{B Q}$ | $P_5$ |
| $B \rightarrow b\ B$ | $P_6$ |
| $\quad\| \quad \lambda$ | $P_7$ |
| $Q \rightarrow q$ | $P_8$ |
| $\quad\| \quad \lambda$ | $P_9$ |

# Computing predict(p)

▶ For the grammar on the right

    ▶ symbolDerivesEmpty[$X$] for $X \in N$

| S | C | A | B | Q |
|---|---|---|---|---|
| T | T | T | T | T |

▶ ruleDerivesEmpty[p] for p $\in P$:

| $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_5$ | $P_6$ | $P_7$ | $P_8$ | $P_9$ |
|---|---|---|---|---|---|---|---|---|
| T | F | T | F | T | F | T | F | T |

| | | |
|---|---|---|
| S → A C | | $P_1$ |
| C → c | | $P_2$ |
| &#124; λ | | $P_3$ |
| A → **a B C d** | | $P_4$ |
| &#124; **B Q** | | $P_5$ |
| B → b B | | $P_6$ |
| &#124; λ | | $P_7$ |
| Q → q | | $P_8$ |
| &#124; λ | | $P_9$ |

back

# Computing predict(p)

**predict**($p$: $X \rightarrow X_1 X_2 \dots X_m$) // returns a set of tokens    Slide 21

{    *ans* = **first**($X_1 X_2 \dots X_m$);

      if ruleDerivesEmpty[$p$] then  // when $X_1 X_2 \dots X_m$ may be empty

         *ans* = *ans* $\cup$ **follow**($X$);    slide 22

      return *ans*;

}

back

▸   **first**($X_1 X_2 \dots X_m$) returns a set of tokens each of which is the **first token in a sentence** derived from $X_1 X_2 \dots X_m$. Formally:

   first($X_1 X_2 \dots X_m$) = {t $\in$ $T$ | $\exists w \in T^*$, [$X_1 X_2 \dots X_m \Rightarrow^* tw$] }

# Computing first $(X_1X_2\ldots X_m)$

**first** $(X_1X_2\ldots X_m)$ // returns a set of tokens

{   for each nonterminal X in the language

   visitedFirst[X] = false;

*ans* = **internalFirst**$(X_1X_2\ldots X_m)$;

return *ans*;

}

# Computing first $(X_1 X_2 \ldots X_m)$

The main ideas for computing **internalFirst**$(X_1 X_2 \ldots X_{m:})$:

1. If $X_1 X_2 \ldots X_m = \lambda$, there is no first token. Return empty set.

   internalFirst$(\lambda)$ returns $\varnothing$

2. If $X_1$ is a terminal symbol, the first token is this symbol. Return $\{X_1\}$.

   internalFirst$(b\ B)$ returns b

3. If $X_1$ is a nonterminal

   i. Look at every rule for $X_1$ and find the first tokens of $X_1$.

   What does internalFirst$(A\ C)$ do?    Slide 25

   ii. If $X_1$ may derive empty, find the first tokens for $X_2 \ldots X_m$.

**internalFirst**$(X_1 X_2 \ldots X_m)$  //returns a set of tokens

```
{    if (m == 0)  return ∅;                    /* 1 */
     if (X₁ is a terminal symbol)  return {X₁} /* 2 */
     /* X₁ is a nonterminal */
     ans = ∅;
     if not visitedFirst[X₁]
          visitedFirst[X₁] = true;
          for the RHS of each rule for X₁        /* 3.i */
               ans = ans ∪ internalFirst(RHS);
     if symbolDerivesEmpty[X₁]                  /* 3.ii */
          ans = ans ∪ internalFirst(X₂…Xₘ);
```

$S \rightarrow A\ C$
$C \rightarrow c$
$\quad | \quad \lambda$
$A \rightarrow a\ B\ C\ d$
$\quad | \quad B\ Q$
$B \rightarrow b\ B$
$\quad | \quad \lambda$
$Q \rightarrow q$
$\quad | \quad \lambda$

```
     return ans;          }
```

internalFirst(B Q) returns {b, q}

# Example first(*B*)

| | | |
|---|---|---|
| A | → | B |
| | \| | a |
| B | → | A |
| | \| | b |

**first** ($X_1 X_2 \ldots X_m$) // returns a set of tokens
{     for each nonterminal X in the language
        visitedFirst[X] = false;
    *ans* = **internalFirst**($X_1 X_2 \ldots X_m$);
    return *ans*;
}

visitedFirst

| A | B |
|---|---|
| F | F |

internalFirst(*B*) ➡ {*a,b*}

internalFirst(*B*)

/* $X_1$ = B i.e. a nonterminal */
*ans* = $\varnothing$;
if not visitedFirst[$X_1$]
    visitedFirst[$X_1$] = true;
    for the *RHS* of each rule for $X_1$
        *ans* = *ans* $\cup$ **internalFirst**(*RHS*);

visitedFirst

| A | B |
|---|---|
| F | T |

internalFirst(*A*) $\cup$
internalFirst(*b*)

$$A \quad \rightarrow B$$
$$| \quad a$$
$$B \quad \rightarrow A$$
$$| \quad b$$

internalFirst(*A*)

```
/* X₁ = A i.e. a nonterminal */
ans = ∅;
if not visitedFirst[X₁]
    visitedFirst[X₁] = true;
    for the RHS of each rule for X₁
        ans = ans ∪ internalFirst(RHS);
if symbolDerivesEmpty[X₁]
    ans = ans ∪ internalFirst(X₂…Xₘ);

return ans;
```

visitedFirst:

| A | B |
|---|---|
| T | T |

internalFirst(*B*) ∪
internalFirst(*a*)

symbolsDerivesEmpty:

| A | B |
|---|---|
| F | F |

internalFirst(*A*) ➡ **{a}**

$$A \quad \rightarrow B$$
$$| \quad a$$
$$B \quad \rightarrow A$$
$$| \quad b$$

internalFirst($B$)

/* $X_I$ = B i.e. a nonterminal */
ans = $\varnothing$;
if not visitedFirst[$X_I$]
    visitedFirst[$X_I$] = true;
    for the *RHS* of each rule for $X_I$
        ans = ans $\cup$ **internalFirst**(*RHS*);
if symbolDerivesEmpty[$X_I$]
    ans = ans $\cup$ **internalFirst**($X_2 \ldots X_m$);

return *ans*;  ⟹  $\varnothing$

visitedFirst:

| A | B |
|---|---|
| T | T |

symbolsDerivesEmpty:

| A | B |
|---|---|
| F | F |

**internalFirst(*B*)** ➡ $\varnothing$

# Computing follow(X)

▸ **follow**(X) returns a set of tokens that can appear right behind the nonterminal X in a phrase derived from the start symbol S. Formally,

follow(X) = {t $\in$ T | $\exists \alpha, \beta \in$ (N $\cup$ T)*, [S $\Rightarrow$* $\alpha$Xt$\beta$ ]}

**follow** (*X*) // returns a set of tokens that may follow *X*

{    for each nonterminal *Y* in the language

      visitedFollow[*Y*] = false;

   *ans* = **internalFollow**(*X*);

   return *ans*;

}

back

# Main ideas of internalFollow(*X*)

How do we find what tokens may follow *X*?

1. Find each occurrence of *X* in all RHS, E.g. what may follow B:

   $$A \to a\ B\ C\ d$$
   $$A \to B\ Q$$
   $$B \to b\ B$$

2. For each such occurrence, find the **first tokens** of the string after *X*. If this string derives empty, call internalFollow(LHS) to find what tokens follow the LHS nonterminal. E.g. in A $\to$ B Q.

$$S \to A\ C$$
$$C \to c$$
$$\quad | \quad \lambda$$
$$A \to a\ B\ C\ d$$
$$\quad | \quad B\ Q$$
$$B \to b\ B$$
$$\quad | \quad \lambda$$
$$Q \to q$$
$$\quad | \quad \lambda$$

back

# Computing follow(X)

**internalFollow**(*Y*) // *Y* is a nonterminal

{    *ans* = $\varnothing$;

    if not visitedFollow[*Y*]

        visitedFollow[*Y*] = true;

        for each occurrence of *Y* in the RHS of all rules

            tail = stream of symbols that appear after *Y* ;

            ans = ans $\cup$ **first**(tail);

            if **allDeriveEmpty**(tail)

                target = LHS of the rule;

                ans = ans $\cup$ **internalFollow**(target);

    return ans;

}

Occurrence of B in
'A $\rightarrow$ a B C d'

tail = 'C d'
first(tail) returns **{c, d}**
**ans =** $\varnothing \cup$ **{c, d}**

Occurrence of B in
'A $\rightarrow$ B Q'

tail = 'Q'
first(Q) returns **{q}**
**ans = {c, d}** $\cup$ **{q}**
target = 'A'
internalFollow(A)
          returns **{c}**
**ans = {c, d, q}** $\cup$ **{c}**

Occurrence of B in
'B $\rightarrow$ b B'

…
**ans = {c, d, q}** $\cup \varnothing$

# Computing follow(A)

**allDeriveEmpty**($\beta$) // $\beta$ is a stream of symbols

{

    for each symbol $X$ in $\beta$

        if $X$ is a terminal or not symbolDerivesEmpty[$X$]

           return false;

    return true;

}

E.g. allDerivesEmpty(tail) is called when tail = "C d".  It will return false.

| | S | C | A | B | Q |
|---|---|---|---|---|---|
| symbolDerivesEmpty | T | T | T | T | T |

# Obtaining LL(1) Grammars

▸ LL(1) requires a unique combination of a nonterminal and a lookahead symbol to decide which rule to use. ┌─────────┐ slide21 └─────────┘

▸ Two common categories of production rules make a grammar not LL(1): *common prefixes* and *left recursion*.

## Common Prefixes

▸ If the RHSs of two rules for the same nonterminal start with the same lookahead symbol, the grammar is not LL(1).

| Example |
|---|
| *Expr* → number plus *Expr* <br>     \|    number |

# Obtaining LL(1) Grammars

> Example

*Expr* $\rightarrow$ number plus *Expr* | *Factor*
*Factor* $\rightarrow$ number

▸ One way is to eliminate common prefixes is by introducing new nonterminals (left factoring a grammar):

> Example

*Expr* $\rightarrow$ number *Expr'*
*Expr'* $\rightarrow$ plus Expr | $\lambda$

This grammar accepts the same language as the one on the previous slide, and this language is LL(1).

# Obtaining LL(1) Grammars

## Left Recursion

▸ If the RHS of a rule starts with the LHS nonterminal, the grammar is not LL(1):

| Example |
| --- |
| $StmtList \rightarrow StmtList$ semicolon $Stmt$ <br> $\qquad\qquad$ \| $Stmt$ <br><br> … |

The method that parses *StmtList*, parseStmtList() will call itself repeatedly 'forever'.

Slide 21

Syntax Analysis    CZ3007

# Obtaining LL(1) Grammars

1. Change left recursion to right recursion:

   $StmtList \rightarrow Stmt$ semicolon $StmtList$
   $\qquad\qquad |\quad Stmt$

2. Remove the common prefix:

   $StmtList \rightarrow Stmt \quad StmtListTail$
   $StmtListTail \rightarrow$ semicolon $StmtList$
   $\qquad\qquad\qquad |\quad \lambda$

   Slide 17

   Slide 21

3. May want to remove mutual recursion:

   $StmtList \rightarrow Stmt \quad StmtListTail$
   $StmtListTail \rightarrow$ semicolon $Stmt \quad StmtListTail$
   $\qquad\qquad\qquad |\quad \lambda$

Syntax Analysis    CZ3007

# Syntactic Error Recovery

▸ A compiler should produce a useful set of diagnostic messages when presented with a faulty program.

▸ Thus after an error is detected it is desirable to recover from it and continue the syntax analysis.

▸ Semantic analysis and code generation will be disabled.

▸ In a simple form of error recovery, the parser skips input tokens until it finds a delimiter (e.g. a semicolon) to end the parsing of the current nonterminal.

▸ The method for parsing a nonterminal is augmented with an extra parameter that is a set of delimiters.

# Example

parseA(ts, termset)

{ // ts is the input token stream

  if (ts.peek() $\in$ **{a}**)

      **match(a);  parseB(ts, {d} $\cup$ termset);**

      **match(d);  match(e);**

  else if (ts.peek() $\in$ **{b}**)

      **parseB(ts, {q} $\cup$ termset);  …**

  else

      error("expected an a or b");

      skip input till a symbol in termset is found

}

q should follow B

$A \rightarrow$ **a B d e**

    | **B Q e**

$B \rightarrow b$

$Q \rightarrow q$

**End-of-file** symbol is in the *termset* of every parsing method

# Bottom-up Parsing

▸ Recall:  Slide 16

▸ Bottom-up parsers are commonly used in the syntax analysis phase of a compiler because of its power, efficiency and ease of construction.

▸ Grammar features like common prefixes and left recursion need to be addressed before top-down parsing can be used.  But they can be accommodated without issue in bottom-up parsing.

# Bottom-up Parsing

▸ This parsing technique is known by a few names:

1. **Bottom-up**, because it works its way from the terminal symbols to the grammar's start symbol.

2. **Shift-reduce**, because the two prevalent actions taken by the parser are to *shift* symbols onto the parse stack and to *reduce* a string of such symbols at the top-of-stack to one of the grammar's nonterminals.

3. **LR(k)**, because it scans the input from left to right, producing a rightmost derivation in reverse, using k symbols of lookahead.

# Rightmost Derivation in Reverse

| Rule | Derivation |
|------|------------|
| 1 | Start $\Rightarrow$ S  $ |
| 2 | $\Rightarrow$ A C $ |
| 3 | $\Rightarrow$ A c $ |
| 5 | $\Rightarrow$ a B C d c $ |
| 4 | $\Rightarrow$ a B d c $ |
| 7 | $\Rightarrow$ a b B d c $ |
| 7 | $\Rightarrow$ a b b B d c $ |
| 8 | $\Rightarrow$ a b b d c $ |

1. Start $\rightarrow$ S $
2. S $\rightarrow$ A C
3. C $\rightarrow$ c
4.     | $\lambda$
5. A $\rightarrow$ a B C d
6.     | B Q
7. B $\rightarrow$ b B
8.     | $\lambda$
9. Q $\rightarrow$ q
10.     | $\lambda$

back

# LR Parsing Engine

▸ The parsing engine is driven by a table.

▸ The table is indexed using the parser's *current state* and the *next input symbol*.

▸ At each step, the engine looks up the table based on the current state and the next input symbol for an action.

▸ The table entry indicates the action to perform (either a shift or a reduce, till the final action which is accept).

```
call Stack.PUSH(StartState)
accepted ← false
while not accepted do
    action ← Table[Stack.TOS( )][InputStream.PEEK( )]          ①
    if action = shift s
    then
        call Stack.PUSH(s)                                      ②
        if s ∈ AcceptStates                                     ③
        then   accepted ← true
        else   call InputStream.ADVANCE( )
    else
        if action = reduce A→γ
        then                          // apply rule A → γ
            call Stack.POP(|γ|)                                 ④
            call InputStream.PREPEND(A)                         ⑤
        else
            call ERROR( )                                       ⑥
```

Figure 6.3: Driver for a bottom-up parser.

input: a b b d c $

0

input: b b d c $

| 0 | 3 |
|---|---|
|   | a |

input: b d c $

| 0 | 3 | 2 |
|---|---|---|
|   | a | b |

input: d c $

| 0 | 3 | 2 | 2 |
|---|---|---|---|
|   | a | b | b |

shift state 3

Reduce by rule 8

| State | a | b | c | d | q | $ | Start | S | A | B | C | Q |
|-------|---|---|---|---|---|---|-------|---|---|---|---|---|
| 0 | 3 | 2 | 8 |   | 8 | 8 | accept | 4 | 1 | 5 |   |   |
| 1 |   |   | 11 |   |   | 4 |   |   |   |   | 14 |   |
| 2 |   | 2 | 8 | 8 | 8 | 8 |   |   |   | 13 |   |   |
| 3 |   | 2 | 8 | 8 |   |   |   |   |   | 9 |   |   |
| 4 |   |   |   |   |   | 8 |   |   |   |   |   |   |
| 5 |   |   | 10 |   | 7 | 10 |   |   |   |   |   | 6 |
| 6 |   |   | 6 |   |   | 6 |   |   |   |   |   |   |
| 7 |   |   | 9 |   |   | 9 |   |   |   |   |   |   |
| 8 |   |   |   |   |   | 1 |   |   |   |   |   |   |
| 9 |   |   | 11 | 4 |   |   |   |   |   |   | 10 |   |
| 10 |   |   |   | 12 |   |   |   |   |   |   |   |   |
| 11 |   |   |   | 3 |   | 3 |   |   |   |   |   |   |
| 12 |   |   | 5 |   |   | 5 |   |   |   |   |   |   |
| 13 |   |   | 7 | 7 | 7 | 7 |   |   |   |   |   |   |
| 14 |   |   |   |   |   | 2 |   |   |   |   |   |   |

back

Figure 6.5: Parse table for the grammar shown in **Slide 45**.

8.    B → λ

input: B d c $

| 0 | 3 | 2 | 2 |
|---|---|---|---|
|   | a | b | b |

input: d c $

| 0 | 3 | 2 | 2 | 13 |
|---|---|---|---|---|
|   | a | b | b | B |

7.    B → b B

input: B d c $

| 0 | 3 | 2 |
|---|---|---|
|   | a | b |

input: d c $

| 0 | 3 | 2 | 13 |
|---|---|---|---|
|   | a | b | B |

**shift state 3**

**Reduce by rule 8**

| State | a | b | c | d | q | $ | Start | S | A | B | C | Q |
|-------|---|---|---|---|---|---|-------|---|---|---|---|---|
| 0 | 3 | 2 | 8 |   | 8 | 8 | accept | 4 | 1 | 5 |   |   |
| 1 |   |   | 11 |   |   | 4 |   |   |   |   | 14 |   |
| 2 |   | 2 | 8 | 8 | 8 | 8 |   |   |   | 13 |   |   |
| 3 |   | 2 | 8 | 8 |   |   |   |   |   | 9 |   |   |
| 4 |   |   |   |   |   | 8 |   |   |   |   |   |   |
| 5 |   |   | 10 |   | 7 | 10 |   |   |   |   |   | 6 |
| 6 |   |   | 6 |   |   | 6 |   |   |   |   |   |   |
| 7 |   |   | 9 |   |   | 9 |   |   |   |   |   |   |
| 8 |   |   |   |   |   | 1 |   |   |   |   |   |   |
| 9 |   |   | 11 | 4 |   |   |   |   |   |   | 10 |   |
| 10 |   |   |   | 12 |   |   |   |   |   |   |   |   |
| 11 |   |   |   | 3 |   | 3 |   |   |   |   |   |   |
| 12 |   |   | 5 |   |   | 5 |   |   |   |   |   |   |
| 13 |   |   | 7 | 7 | 7 | 7 |   |   |   |   |   |   |
| 14 |   |   |   |   |   | 2 |   |   |   |   |   |   |

Figure 6.5: Parse table for the grammar shown in **Slide 45**

back

Syntax Analysis    CZ3007

7.  B → b B

input: B d c $

| 0 | 3 |
|---|---|
|   | a |

input: d c $

| 0 | 3 | 9 |
|---|---|---|
|   | a | B |

4.  C → λ

input: C d c $

| 0 | 3 | 9 |
|---|---|---|
|   | a | B |

input: d c $

| 0 | 3 | 9 | 10 |
|---|---|---|----|
|   | a | B | C  |

50



*shift state 3* — *Reduce by rule 8*

| State | a | b | c | d | q | $ | Start | S | A | B | C | Q |
|-------|---|---|---|---|---|---|-------|---|---|---|---|---|
| 0 | 3 | 2 | 8 |  | 8 | 8 | accept | 4 | 1 | 5 |  |  |
| 1 |  |  | 11 |  |  | 4 |  |  |  |  | 14 |  |
| 2 |  | 2 | 8 | 8 | 8 | 8 |  |  |  | 13 |  |  |
| 3 |  | 2 | 8 | 8 |  |  |  |  |  | 9 |  |  |
| 4 |  |  |  |  |  | 8 |  |  |  |  |  |  |
| 5 |  |  | 10 |  | 7 | 10 |  |  |  |  |  | 6 |
| 6 |  |  | 6 |  |  | 6 |  |  |  |  |  |  |
| 7 |  |  | 9 |  |  | 9 |  |  |  |  |  |  |
| 8 |  |  |  |  |  | 1 |  |  |  |  |  |  |
| 9 |  |  | 11 | 4 |  |  |  |  |  |  | 10 |  |
| 10 |  |  |  | 12 |  |  |  |  |  |  |  |  |
| 11 |  |  |  | 3 |  | 3 |  |  |  |  |  |  |
| 12 |  |  | 5 |  |  | 5 |  |  |  |  |  |  |
| 13 |  |  | 7 | 7 | 7 | 7 |  |  |  |  |  |  |
| 14 |  |  |  |  |  | 2 |  |  |  |  |  |  |

Figure 6.5: Parse table for the grammar shown in **Slide 45**

**Slide 47**

input: c $

| 0 | 3 | 9 | 10 | 12 |
|---|---|---|----|----|
|   | a | B | C  | d  |

5.  A → a B C d

input: A c $

| 0 |
|---|

input: c $

| 0 | 1 |
|---|---|
|   | A |

input: $

| 0 | 1 | 11 |
|---|---|----|
|   | A | c  |

back

shift state 3

Reduce by rule 8

| State | a | b | c | d | q | $ | Start | S | A | B | C | Q |
|-------|---|---|---|---|---|---|-------|---|---|---|---|---|
| 0  | 3 | 2 | 8 |    | 8 | 8 | accept | 4 | 1 | 5 |    |   |
| 1  |   |   | 11 |   |   | 4 |        |   |   |   | 14 |  |
| 2  |   | 2 | 8 | 8 | 8 | 8 |        |   |   | 13 |   |  |
| 3  |   | 2 | 8 | 8 |   |   |        |   |   | 9 |    |   |
| 4  |   |   |   |   |   | 8 |        |   |   |   |    |   |
| 5  |   |   | 10 |  | 7 | 10 |       |   |   |   |    | 6 |
| 6  |   |   | 6 |   |   | 6 |        |   |   |   |    |   |
| 7  |   |   | 9 |   |   | 9 |        |   |   |   |    |   |
| 8  |   |   |   |   |   | 1 |        |   |   |   |    |   |
| 9  |   |   | 11 | 4 |  |   |        |   |   |   | 10 |  |
| 10 |   |   |   | 12 |  |   |        |   |   |   |    |   |
| 11 |   |   |   | 3 |  | 3 |        |   |   |   |    |   |
| 12 |   |   | 5 |   |   | 5 |        |   |   |   |    |   |
| 13 |   |   | 7 | 7 | 7 | 7 |        |   |   |   |    |   |
| 14 |   |   |   |   |   | 2 |        |   |   |   |    |   |

Figure 6.5: Parse table for the grammar shown in **Slide 45**

Syntax Analysis    CZ3007

3.     C → c

input: C $

| 0 | I |
|---|---|
|   | A |

input: $

| 0 | I | 14 |
|---|---|----|
|   | A | C  |

2.     S → A C

input: S $

| 0 |
|---|

input: $

| 0 | 4 |
|---|---|
|   | S |

shift state 3

Reduce by rule 8

| State | a | b | c | d | q | $ | Start | S | A | B | C | Q |
|-------|---|---|---|---|---|---|-------|---|---|---|---|---|
| 0 | 3 | 2 | 8 |   | 8 | 8 | accept | 4 | 1 | 5 |   |   |
| 1 |   |   | 11 |   |   | 4 |   |   |   |   | 14 |   |
| 2 |   | 2 | 8 | 8 | 8 | 8 |   |   |   | 13 |   |   |
| 3 |   | 2 | 8 | 8 |   |   |   |   |   | 9 |   |   |
| 4 |   |   |   |   |   | 8 |   |   |   |   |   |   |
| 5 |   |   | 10 |   | 7 | 10 |   |   |   |   |   | 6 |
| 6 |   |   | 6 |   |   | 6 |   |   |   |   |   |   |
| 7 |   |   | 9 |   |   | 9 |   |   |   |   |   |   |
| 8 |   |   |   |   |   | 1 |   |   |   |   |   |   |
| 9 |   |   | 11 | 4 |   |   |   |   |   |   | 10 |   |
| 10 |   |   |   | 12 |   |   |   |   |   |   |   |   |
| 11 |   |   |   | 3 |   | 3 |   |   |   |   |   |   |
| 12 |   |   | 5 |   |   | 5 |   |   |   |   |   |   |
| 13 |   |   | 7 | 7 | 7 | 7 |   |   |   |   |   |   |
| 14 |   |   |   |   |   | 2 |   |   |   |   |   |   |

back

Figure 6.5: Parse table for the grammar shown in **Slide 45**

**Slide 47**

Syntax Analysis    CZ3007

input: $

| 0 | 4 | 8 |
|---|---|---|
|   | S | $ |

I. Start → S $

input: Start $

| 0 |
|---|
|   |

accept

| State | a | b | c | d | q | $ | Start | S | A | B | C | Q |
|-------|---|---|---|---|---|---|-------|---|---|---|---|---|
| 0 | 3 | 2 | 8 |   | 8 | 8 | accept | 4 | 1 | 5 |    |    |
| 1 |   |   | 11 |   |   | 4 |        |   |   |   | 14 |   |
| 2 |   | 2 | 8 | 8 | 8 | 8 |        |   |   | 13 |   |    |
| 3 |   | 2 | 8 | 8 |   |   |        |   |   | 9 |    |    |
| 4 |   |   |   |   |   | 8 |        |   |   |   |    |    |
| 5 |   |   | 10 |   | 7 | 10 |       |   |   |   |    | 6 |
| 6 |   |   | 6 |   |   | 6 |        |   |   |   |    |    |
| 7 |   |   | 9 |   |   | 9 |        |   |   |   |    |    |
| 8 |   |   |   |   |   | 1 |        |   |   |   |    |    |
| 9 |   |   | 11 | 4 |   |   |        |   |   |   | 10 |   |
| 10 |  |   |   | 12 |   |   |        |   |   |   |    |    |
| 11 |  |   |   | 3 |   | 3 |        |   |   |   |    |    |
| 12 |  |   | 5 |   |   | 5 |        |   |   |   |    |    |
| 13 |  |   | 7 | 7 | 7 | 7 |        |   |   |   |    |    |
| 14 |  |   |   |   |   | 2 |        |   |   |   |    |    |

Figure 6.5: Parse table for the grammar shown in **Slide 45**

Slide 47

# Classes of Bottom-up Parsers

▸ In practice, bottom-up parsers only depend on a single token of lookahead.

▸ They can also be described as finite automata (though of a more complicated kind than DFAs)

▸ There are several methods for constructing parse tables:

  ▸ LR(0): simplest, fails for many practical grammars;

  ▸ LR(1): quite general, can handle almost all practically interesting grammars;

  ▸ LALR(1): faster, slightly weaker variant of LR(1), used by most bottom-up parser generators.

Syntax Analysis    CZ3007

# LR(0) Automata and Table Construction

▸ The table construction process analyzes the grammar.

▸ Each state corresponds to a row of the parser table.

▸ Each symbol in the terminal and nonterminal sets corresponds to a column of the table.

▸ During parsing, we want to keep track of where we are in the grammar.

▸ To do this, we use **LR(0) items**:  an LR(0) item is **a grammar rule with a marker "•"** showing the current progress of the parser in recognizing the RHS of the rule.

# LR(0) Automata and Table Construction

▸ Symbols before the • have already been seen. The first symbol after • is what we expect next. For examples,
E → • plus E E,    E → plus • E E,    E → plus E • E

▸ For an item A→α • β, the item is called *initial* if α →λ, and *final* if β →λ. A final item for the start symbol is called accepting.

E. g.

Start → • S $        // an initial item

S → A C •            // a final item

# LR(0) Automata and Table Construction

▸ An LR(0) state is **a set of LR(0) items**, which is *closed* in the sense that if the state contains an item with a nonterminal A immediately following the marker, we add the initial items for A, i.e., items for all rules of A with the marker at the beginning of the RHS.  This is called taking *the closure of the item*.

**E.g.,  For an item S → • A,**

    **we add**

    **A → • a B d, A→ • B Q,**

    **B → • b**

**i.e. closure = {S → • A, A → • a B d,**

        **A→ • B Q, B → • b}**

| |
|---|
| S → A |
| A → a B d |
|     &#124; B Q |
| B → b |
| Q → q |

# LR(0) Automata and Table Construction

▸ We describe the parsing process as a finite automaton.

▸ The start state is the closure of the initial items of the start symbol:

For the grammar on the right, there is only one initial item for the Start symbol:

$$Start \rightarrow \bullet \ E \ \$$$

```
Start → E  $
E → plus  E  E
    |    num
```

Taking the closure, we get the set of items in state 0:

0
```
Start → • E  $
E → • plus  E  E
E → • num
```

In slides 48, 51, 52, state 0 is on TOS

# LR(0) Automata and Table Construction

▸ For each symbol γ that appears to the right of the marker of an item (or items), we can **shift** over it and transition to a new state:

   ▸ replace an item of the form A→α • γ β with A→α γ • β

   ▸ throw away all other items;

   ▸ take closure

▸ Transitions are labelled by γ which is either a terminal or a nonterminal.

▸ If there is a final item of the form A →α •, we can **reduce**.

▸ An accepting state is one that contains a final item for start symbol.

Syntax Analysis    CZ3007

# LR(0) Automata and Table Construction

▶ Continue from State 0 of our finite automata in slide 58, each of the three items will lead to a new state. Each new state leads to more new states.



Slide 48

**0**
Start → • E  $
E → • plus  E  E
E → • num

**3**
Start → E • $

**4**
Start → E $ •

**2**
E → num •

**1**
E → plus • E  E
E → • plus  E  E
E → • num

**5**
E → plus E • E
E → • plus  E  E
E → • num

**6**
E → plus E E  •

Start → E  $
E → plus  E  E
        |    num

# LR(0) Automata and Table Construction

▸ The LR(0) automaton needs to be used together with a stack of states; this kind of automaton is known as a *pushdown automaton.*

▸ The state of the parser is the state on top of the parser stack.

Slide 51

▸ Actions taken during bottom-up parsing are of four types:

  ▸ **shift**(i): consume next input symbol, push state i onto state stack;

Slide 47

  ▸ **reduce**(A→γ): reduce by rule A→γ, i.e., pop |γ| states from the stack and consider A as next input symbol;

  ▸ **accept**: report that input was parsed successfully;

  ▸ **error**: report a parse error

# LR(0) Automata and Table Construction

▸ An LR(0) parse table is a compact representation of an LR(0) automaton

▸ Rows are indexed by states, columns by symbols; cell in row r , column c contain a single parsing action to take when encountering input symbol c in state r

▸ Constructing a parse table from an LR(0) automaton is easy:

  ▸ For every transition from state s to state s' labelled with symbol x, enter **shift**(s') into the cell in row s, column x

  ▸ If state s contains a final item $A \rightarrow \beta \bullet$, enter **reduce**$(A \rightarrow \beta)$ into all cells of row s

  ▸ For cell (0, StartSymbol), enter **accept**

  ▸ Enter **error** into any remaining empty cells

back

# Example



0
Start → • E $
E → • plus E E
E → • num

3
Start → E • $

4
Start → E $ •

2
E → num •

1
E → plus • E E
E → • plus E E
E → • num

5
E → plus E • E
E → • plus E E
E → • num

6
E → plus E E •

back

Parser Table

| state | num | plus | $ | Start | E |
|-------|-----|------|---|-------|---|
| 0 | Shift 2 | Shift 1 | error | accept | Shift 3 |
| 1 | Shift 2 | Shift 1 | error | error | Shift 5 |
| 2 | Reduce (E → num ) | | | | |
| 3 | error | error | Shift 4 | error | error |
| 4 | Reduce (Start → E $ ) | | | | |
| 5 | Shift 2 | Shift 1 | error | error | Shift 6 |
| 6 | Reduce(E → plus E E ) | | | | |

# Conflicts

▸ Sometimes when trying to construct an LR(0) parse table we end up with two different actions in the same cell; this is known as a conflict.

▸ There are two kinds of conflicts:

   ▸ **shift-reduce conflict**: the same cell contains both a **shift()** action and a **reduce()** action;

   ▸ **reduce-reduce conflict**: the same cell contains two different **reduce()** actions.

▸ **Question: Can there be a shift-shift conflict?**

Syntax Analysis    CZ3007

# Shift-reduce Conflict

▸ **shift-reduce conflict**: if a state contains both a non-final item A → β • γ (a **shift()** action) and a final item A → β • (a **reduce()** action)

> E.g. A state has these two items:
>
> > IfStatement → IF Cond THEN StatList • ELSE StatList
> >
> > IfStatement → IF Cond THEN StatList •

▸ Example:

Start → E $

E → num plus E

   |   num

<div style="text-align:right;">

back

</div>

# Example: shift-reduce conflict

**0**

Start → • E  $

E → • num plus  E

E → • num

**2**

Start → E • $

**4**

Start → E $ •

E → E

$ →

**num**

**1**

E → num • plus  E

E → num •

**num**

**3**

E → num plus • E

E → • num plus  E

E → • num

**E**

**5**

E → num plus E •

**plus**

| state | num | plus | $ | Start | E |
|-------|-----|------|------|-------|------|
| **0** | Shift 1 | error | error | accept | Shift 2 |
| **1** | | Shift/reduce? | | | |
| **2** | … | | | | |

back

# Shift-reduce Conflict

▸ **A shift-reduce conflict** may be eliminated by rewriting the grammar.  For example,

Rewrite                                         to

Start → E $                                Start → E $
E → num plus  E                          E → E  plus num
   |   num                                     |   num

# An Ambiguous Grammar

▸ Ambiguous grammars always lead to conflicts. For example,

$$Start \rightarrow E \ \$$$

$$E \rightarrow E \ plus \ E$$

$$| \ num$$

**Rewriting the grammar will solve the problem in this case.**

LR(0) automaton:



Syntax Analysis    CZ3007

# LALR(1) – Look Ahead LR with one token lookahead

- Due to its balance of power and efficiency, LALR(1) is the most popular LR table-building method.
  - For every transition from state s to state s' labelled with symbol x, enter **shift**(s') into the cell in row s, column x
  - **If state s contains a final item A→ β •, enter reduce(A→β) into the cells of row s for each token T ∈ itemFollow(*s, item*)**

  - For cell (0, StartSymbol), enter **accept**
  - Enter **error** into any remaining empty cells

- **itemFollow()** keeps track of the symbols that can follow the *item* after reduction *from this state*.

Slide 66

# LALR Propagation Graph

- Consider the LR(0) table, the pair $(s, A \rightarrow \alpha \bullet \beta)$ suffices to identify an item $A \rightarrow \alpha \bullet \beta$ that occurs in state $s$.

Slide 63

- Each item in an LR(0) table is represented by a vertex in the **LALR propagation graph**.

- We will compute itemFollow() for each item.

- The LALR table is constructed with reference to the itemFollow() computed for the items.

- The propagation graph will not be retained after constructing the LALR table.

# Generating the Propagation Graph

A. Setup:

1. Create the LR(0) finite automaton.
2. For each (state, item), create a vertex *v* in the graph.
3. Initialize all itemFollow[*v*] to $\varnothing$.
4. Initialize itemFollow[(0, StartSymbol Productions)] = {$}

B. Build the propagation graph

C. Propagate itemFollow[ ]

back

# Generating the Propagation Graph (Example)

$P_1$     Start $\rightarrow$ S $

$P_2$     S     $\rightarrow$ A  B

$P_3$             |  a  c

$P_4$             |  x  A c

$P_5$     A     $\rightarrow$ a

$P_6$     B     $\rightarrow$ b

$P_7$             | $\lambda$

back

# A: Setup (create LR(0) finite automaton)

**0**

Start → • S $
S → • A  B
S → • x  A  c
S → • a  c
A → • a

**2**

S → A • B
B → • b
B → • λ

**8**

S → A  B •

**7**

B → b •

**11**

S → x  A  c •

**1**

S → x • A  c
A → • a

**9**

S → x  A • c

**10**

A → a •

**3**

S → a • c
A →  a •

**4**

Start → S • $

**6**

S → a c •

**5**

Start → S $ •

Arrows labeled: A (0→2), x (0→1), a (0→3), S (0→4), B (2→8), b (2→7), A (1→9), a (1→10), c (9→11), c (3→6), $ (4→5)

# A: Setup (create vertices, initialize itemFollow)

| 1 | 0, Start → • S $ | {$} |

| 8 | 2, S → A • B | ∅ |

| 17 | 8, S → A B • | ∅ |

| 2 | 0, S → • A  B | ∅ |

| 9 | 2, B → • b | ∅ |

| 16 | 7, B → b • | ∅ |

| 10 | 2, B → • λ | ∅ |

| 3 | 0, S → • x  A  c | ∅ |

| 20 | 11, S → x  A  c • | ∅ |

| 4 | 0, S → • a  c | ∅ |

| 6 | 1, S → x • A  c | ∅ |

| 18 | 9, S → x  A • c | ∅ |

| 5 | 0, A → • a | ∅ |

| 7 | 1, A → • a | ∅ |

| 19 | 10, A → a • | ∅ |

| 11 | 3, S → a • c | ∅ |

| 13 | 4, Start → S • $ | ∅ |

| 14 | 5, Start → S $ • | ∅ |

| 12 | 3, A → a • | ∅ |

| 15 | 6, S → a c • | ∅ |

# B:Building the Propagation Graph (ideas)

- For item $(s, A \rightarrow \alpha \bullet B\gamma)$, any symbol in First$(\gamma)$ can follow each closure item $(s, B \rightarrow \bullet\delta)$ $(2)$ <span style="border:1px solid">Slide 66</span>
  - Put First$(\gamma)$ into the itemFollow set of each initial item of B
- A propagation edge is placed from an item $(s, A \rightarrow \alpha \bullet B\gamma)$ to an item $(t, A \rightarrow \alpha B \bullet \gamma)$ $(1)$ <span style="border:1px solid">back</span>
- For item $(s, A \rightarrow \alpha \bullet B\gamma)$, when $\gamma \Rightarrow^* \lambda$, any symbol that can follow A can also follow B $(3)$
  - Place a propagation edge from $(s, A \rightarrow \alpha \bullet B\gamma)$ to $(s, B \rightarrow \bullet\delta)$

> **itemFollow()** keeps track of the symbols that can follow the **item** after reduction *from this state*

**itemFollow()** keeps track of the symbols that can follow the *item* after reduction *from this state*

**1**

S → x • A  c

A → • a

**0**

Start → • S $

S → • A  B

S → • x  A  c

S → • a  c

A → • a

x

A

S

a

**3**

S → a • c

A →  a •

c

**6**

S → a c •

**2**

S → A • B

B → • b

B → • λ

**4**

Start → S • $

Input scenarios:
1)  a $
2)  a b $
3)  a c $

Part of the LR(0) parse table

Shift 4

| state | a | b | c | x | $ | A | S | ... |
|-------|----|-----|-----|-----|----|----|----|-----|
| 0 | S3 | err | err | S1 | err | S2 | S4 | ... |
| 1 | | | | | | | | |
| 2 | | | | | | | | |
| 3 | err | R5 | S6 | err | R5 | err | err | |

Syntax Analysis    CZ3007

# B: Building the Propagation Graph

| 1 | 0, Start → • S $ | {$} |

| 2 | 0, S → • A  B | {$} |

| 3 | 0, S → • x  A  c | {$} |

| 4 | 0, S → • a  c | {$} |

| 5 | 0, A → • a | {b} |

| 8 | 2, S → A • B | ∅ |

| 9 | 2, B → • b | ∅ |

| 10 | 2, B → • λ | ∅ |

| 6 | 1, S → x • A  c | ∅ |

| 7 | 1, A → • a | {c} |

| 17 | 8, S → A  B • | ∅ |

| 16 | 7, B → b • | ∅ |

| 20 | 11, S → x  A  c • | ∅ |

| 18 | 9, S → x  A • c | ∅ |

| 19 | 10, A → a • | ∅ |

| 11 | 3, S → a • c | ∅ |

| 13 | 4, Start → S • $ | ∅ |

| 14 | 5, Start → S $ • | ∅ |

| 12 | 3, A →  a • | ∅ |

| 15 | 6, S → a c • | ∅ |

# B: Building the Propagation Graph

For each state $s$

For each item of form $A \rightarrow \alpha \bullet B \gamma$ in $s$

vertex $u = (s, A \rightarrow \alpha \bullet B \gamma)$

$s' = \text{Table}[s][B]$

vertex $v = (s', A \rightarrow \alpha B \bullet \gamma)$, successor of $u$ on shift(B)

Add edge $(u, v)$       1

For each vertex $w$ of form $(s, B \rightarrow \bullet \delta)$

Add **first**$(\gamma)$ to itemFollow$[w]$      2

Slide 36

if **allDeriveEmpty**$(\gamma)$

add edge $(u, w)$      3

# C: Propagating itemFollw [ ]

While (making progress)

    For each edge (*u,v*)

        add ItemFollow(*u*) to ItemFollow(*v*)  ← – – – – – – – – –  4

▸ In general, multiple passes can be required

▸ In practice, LALR(1) lookahead computations converge quickly, usually in one or two passes

▸ LALR(1) is a powerful parsing method

▸ LALR(1) grammars are available for all popular programming languages

# C: Propagating itemFollw [ ]

1 | 0, Start → • S $ | {$}

8 | 2, S → A • B | ∅

17 | 8, S → A B • | ∅

2 | 0, S → • A B | {$}

9 | 2, B → • b | ∅

16 | 7, B → b • | ∅

3 | 0, S → • x A c | {$}

10 | 2, B → • λ | ∅

20 | 11, S → x A c • | ∅

4 | 0, S → • a c | {$}

6 | 1, S → x • A c | ∅

18 | 9, S → x A • c | ∅

5 | 0, A → • a | {b}

7 | 1, A → • a | {c}

19 | 10, A → a • | ∅

11 | 3, S → a • c | ∅

13 | 4, Start → S • $ | {$}

14 | 5, Start → S $ • | ∅

12 | 3, A → a • | ∅

15 | 6, S → a c • | ∅

# C: Propagating itemFollw [ ]

1 | 0, Start → • S $ | {$}

8 | 2, S → A • B | {$}

17 | 8, S → A B • | ∅

2 | 0, S → • A B | **{$}**

9 | 2, B → • b | ∅

16 | 7, B → b • | ∅

3 | 0, S → • x A c | **{$}**

10 | 2, B → • λ | ∅

20 | 11, S → x A c • | ∅

4 | 0, S → • a c | **{$}**

6 | 1, S → x • A c | ∅

18 | 9, S → x A • c | ∅

5 | 0, A → • a | **{b, $}**

7 | 1, A → • a | **{c}**

19 | 10, A → a • | ∅

11 | 3, S → a • c | ∅

13 | 4, Start → S • $ | {$}

14 | 5, Start → S $ • | ∅

12 | 3, A → a • | ∅

15 | 6, S → a c • | ∅

# C: Propagating itemFollw [ ]

1 | 0, Start → • S $ | {$}

8 | 2, S → A • B | {$}

17 | 8, S → A B • | ∅

2 | 0, S → • A B | {$}

9 | 2, B → • b | ∅

16 | 7, B → b • | ∅

3 | 0, S → • x A c | {$}

10 | 2, B → • λ | ∅

20 | 11, S → x A c • | ∅

4 | 0, S → • a c | {$}

6 | 1, S → x • A c | {$}

18 | 9, S → x A • c | ∅

5 | 0, A → • a | {b, $}

7 | 1, A → • a | {c}

19 | 10, A → a • | ∅

11 | 3, S → a • c | ∅

13 | 4, Start → S • $ | {$}

14 | 5, Start → S $ • | ∅

12 | 3, A → a • | ∅

15 | 6, S → a c • | ∅

# C: Propagating itemFollw [ ]

1 | 0, Start → • S $ | {$}

2 | 0, S → • A  B | **{$}**

3 | 0, S → • x  A  c | **{$}**

4 | 0, S → • a  c | **{$}**

5 | 0, A → • a | **{b, $}**

8 | 2, S → A • B | {$}

9 | 2, B → • b | ∅

10 | 2, B → • λ | ∅

6 | 1, S → x • A  c | {$}

7 | 1, A → • a | **{c}**

11 | 3, S → a • c | {$}

13 | 4, Start → S • $ | {$}

15 | 6, S → a c • | ∅

12 | 3, A →  a • | {b,$}

17 | 8, S → A  B • | ∅

16 | 7, B → b • | ∅

20 | 11, S → x  A  c • | ∅

18 | 9, S → x  A • c | ∅

19 | 10, A → a • | ∅

14 | 5, Start → S $ • | ∅

# C: Propagating itemFollw [ ]

1  0, Start → • S $  {$}

2  0, S → • A  B  **{$}**

3  0, S → • x  A  c  **{$}**

4  0, S → • a  c  **{$}**

5  0, A → • a  **{b, $}**

11  3, S → a • c  {$}

12  3, A → a •  {b,$}

8  2, S → A • B  {$}

9  2, B → • b  {$}

10  2, B → • λ  {$}

6  1, S → x • A  c  {$}

7  1, A → • a  **{c}**

13  4, Start → S • $  {$}

15  6, S → a c •  {$}

17  8, S → A  B •  {$}

16  7, B → b •  ∅

20  11, S → x  A  c •  ∅

18  9, S → x  A • c  ∅

19  10, A → a •  ∅

14  5, Start → S $ •  ∅

# C: Propagating itemFollw [ ]

1 | 0, Start → • S $ | {$}

2 | 0, S → • A  B | **{$}**

3 | 0, S → • x  A  c | **{$}**

4 | 0, S → • a  c | **{$}**

5 | 0, A → • a | **{b, $}**

11 | 3, S → a • c | {$}

12 | 3, A → a • | {b,$}

8 | 2, S → A • B | {$}

9 | 2, B → • b | {$}

10 | 2, B → • λ | {$}

6 | 1, S → x • A  c | {$}

7 | 1, A → • a | **{c}**

13 | 4, Start → S • $ | {$}

15 | 6, S → a c • | {$}

17 | 8, S → A  B • | {$}

16 | 7, B → b • | {$}

20 | 11, S → x  A  c • | {$}

18 | 9, S → x  A • c | {$}

19 | 10, A → a • | {c}

14 | 5, Start → S $ • | {$}

# The Result:

The LALR(1) parse table:

Start → S $
S     → A  B
      |  a  c
      |  x  A  c
A     → a
B     → b
      |  λ

**0**
Start → • S  $
S → • A  B
S → • a  c
S → • x  A  c
A → • a

**a**

**3**
S → a • c
A → a •

**c**

**6**
S → a c •

Slide 69

**itemFollow(3, A → a •) = {b, \$}**

| state | … | b | c | \$ | … |
|-------|---|---|---|-----|---|
| … | … | | | | |
| **3** | | Reduce(A→ a) | Shift 6 | Reduce(A→ a) | |
| … | … | | | | |

back

# Parser Generator Beaver

▸ Beaver is a LALR(1) parser generator that generates parsers written in Java.

▸ Beaver's syntax is very similar to the notation we have been using for context free grammars, except that Beaver uses = where we have used $\rightarrow$. $\boxed{\text{Slide 6}}$

▸ The rules for a nonterminal must be terminated by a semicolon.

▸ The directive **%terminals** on the first line declares the set of terminals used by the grammar.

▸ Beaver implicitly assumes that every name that isn't declared to be a terminal is a nonterminal.

# Example of a Beaver specification: grammar for arithmetic expressions

**%terminals** PLUS, MINUS, MUL, DIV, NUMBER, LPAREN, RPAREN;

**%goal** Expr;

Expr = Expr PLUS Term

    | Expr MINUS Term

    | Term

    ;

Term = Term MUL Factor

    | Term DIV Factor

    | Factor

    ;

Factor = NUMBER

    | LPAREN Expr RPAREN

    ;

# Some review questions/tasks

1. What does a syntax analyzer do?

2. What are the input and output of a syntax analyzer?

3. What is a context free grammar used for in a compiler?

4. Write the pseudocode for a recursive descent parser for the context free grammar of Expr on slide 6 (reference slide 21). Assume the methods peek(), match() and predict() are provided.

5. Follow the bottom-up parsing engine on slide 47 and use the parse table on slide 48 for the grammar on slide 45 to trace the parsing process for the input "adc$".