

# Parsing

- **TD parsing - LL(1)**
  - First and Follow sets
  - Parse table construction
- **BU Parsing**
  - Handle, viable prefix, items, closures, goto's
  - LR(k): SLR(1), LR(1), LALR(1)
    - Problems with SLR
- *Aho, Sethi, Ullman, Compilers : Principles, Techniques and Tools*
- *Aho + Ullman, Theory of Parsing and Compiling, vol II.*

## TD Parsing

**Elimination of left recursion.**

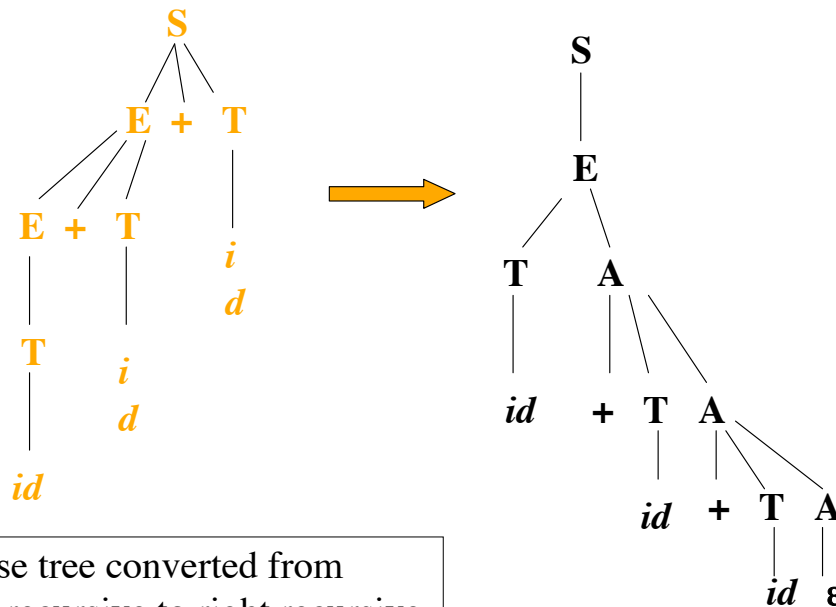
$$\mathbf{E \rightarrow E \alpha \mid \beta \text{ becomes } E \rightarrow \beta A}$$
$$\mathbf{A \rightarrow \alpha A \mid \epsilon}$$

**Example:**

$$\begin{array}{l} \mathbf{S \rightarrow E} \\ \mathbf{E \rightarrow E + T} \\ \mathbf{E \rightarrow T} \\ \mathbf{T \rightarrow id} \end{array} \quad \Longrightarrow \quad \begin{array}{l} \mathbf{S \rightarrow E} \\ \mathbf{E \rightarrow T A} \\ \mathbf{A \rightarrow + T A \mid \epsilon} \\ \mathbf{T \rightarrow id} \end{array}$$

**Can also left factor the grammar removing shared prefixes of right-hand-sides.**

# Parse Tree



# TD Parsing

- **Problem:** predicting which nonterminal to expand next, from a leading string of symbols
- **Idea:** generate parse tree top down so its frontier is always a sentential form
  - Use **First** and **Follow** sets to understand the shape of sentential forms possibly generated by the grammar

# TD Stack Parser, EG

| <u>Stack</u> | <u>Input</u>       | <u>Production</u>        |
|--------------|--------------------|--------------------------|
| \$S          | <i>id+id+id</i> \$ |                          |
| \$E          | <i>id+id+id</i> \$ | $S \rightarrow E$        |
| \$A T        | <i>id+id+id</i> \$ | $E \rightarrow T A$      |
| \$A          | <i>+id+id</i> \$   | $T \rightarrow id$       |
| \$A T        | <i>id+id</i> \$    | $A \rightarrow + T A$    |
| \$A          | <i>+id</i> \$      | $T \rightarrow id$       |
| \$A T        | <i>id</i> \$       | $A \rightarrow + T A$    |
| \$A          | \$                 | $T \rightarrow id$       |
| \$           | \$                 | $A \rightarrow \epsilon$ |

$$\begin{aligned}
 S &\rightarrow E \\
 E &\rightarrow T A \\
 A &\rightarrow + T A \mid \epsilon \\
 T &\rightarrow id
 \end{aligned}$$

See algm in ASU Fig 4.14, p 187

## How to mechanize?

- Define  $\alpha$  to be string of nonterminals and terminals
- **First( $\alpha$ )** is the set of terminals that begin strings derivable from  $\alpha$ .  
If  $\alpha \xRightarrow{*} \epsilon$ , then  $\epsilon$  is in **First( $\alpha$ )**.
- **Follow(A)** is the set of terminals that can appear directly to the right of A in a sentential form  
 $S \xRightarrow{*} \alpha A \beta$  means a is in **Follow(A)**.  
If A can be rightmost symbol in a sentential form, that is,  $X \xRightarrow{*} \alpha A \delta$  where  $\delta \xRightarrow{*} \epsilon$ , then **Follow(A)**  $\supseteq$  **Follow(X)**.

## Example

- $\text{First}(S) = \text{First}(E) = \text{First}(T) = \{id\}$
- $\text{First}(A) = \{+, \epsilon\}$
- $\text{Follow}(S) = \text{Follow}(E) = \text{Follow}(A) = \{\$\}$
- $\text{Follow}(T) = \{+, \$\}$

$$\begin{array}{l} S \rightarrow E \\ E \rightarrow T A \\ A \rightarrow + T A \mid \epsilon \\ T \rightarrow id \end{array}$$

## LL(k) Grammars

- Can choose next production to expand by during TD phase, by looking **k** symbols ahead into input
- Use **First sets** to choose production
- Use **Follow sets** to handle  $\epsilon$  cases

# Example: LL(1)

| <u>Nonterms\Inputs:</u> | <u>id</u> | <u>+</u>  | <u>\$</u> |
|-------------------------|-----------|-----------|-----------|
| S                       | S → E     |           |           |
| E                       | E → T A   |           |           |
| T                       | T → id    |           |           |
| A                       |           | A → + T A | A → ε     |

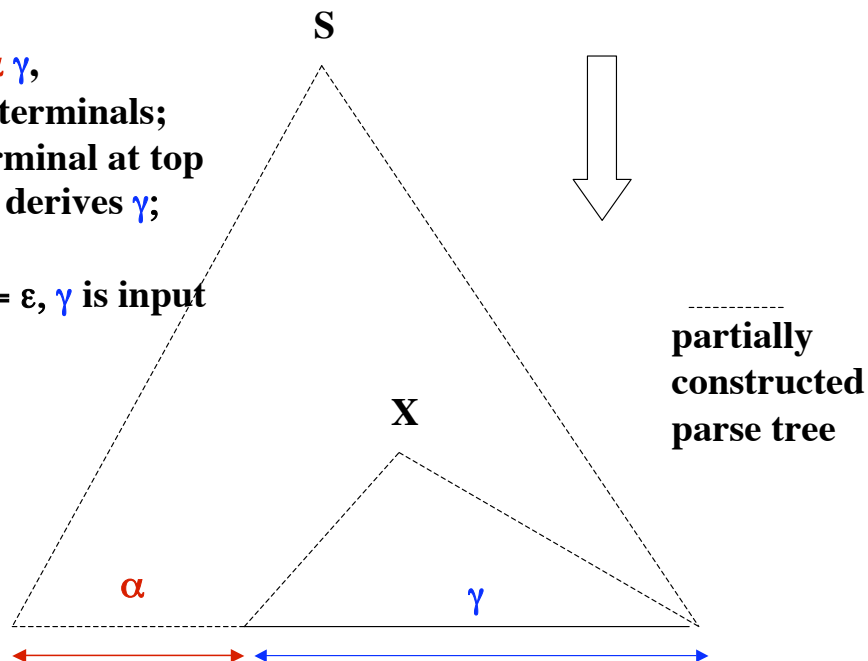
**Ambiguous or left recursive grammars result in multiply defined entries in table**

$\text{First}(S) = \text{First}(E) = \text{First}(T) = \{id\}$   
 $\text{First}(A) = \{+, \epsilon\}$   
 $\text{Follow}(S) = \text{Follow}(E) = \text{Follow}(A) = \{\$\}$   
 $\text{Follow}(T) = \{+, \$\}$

$S \rightarrow E$   
 $E \rightarrow T A$   
 $A \rightarrow + T A \mid \epsilon$   
 $T \rightarrow id$

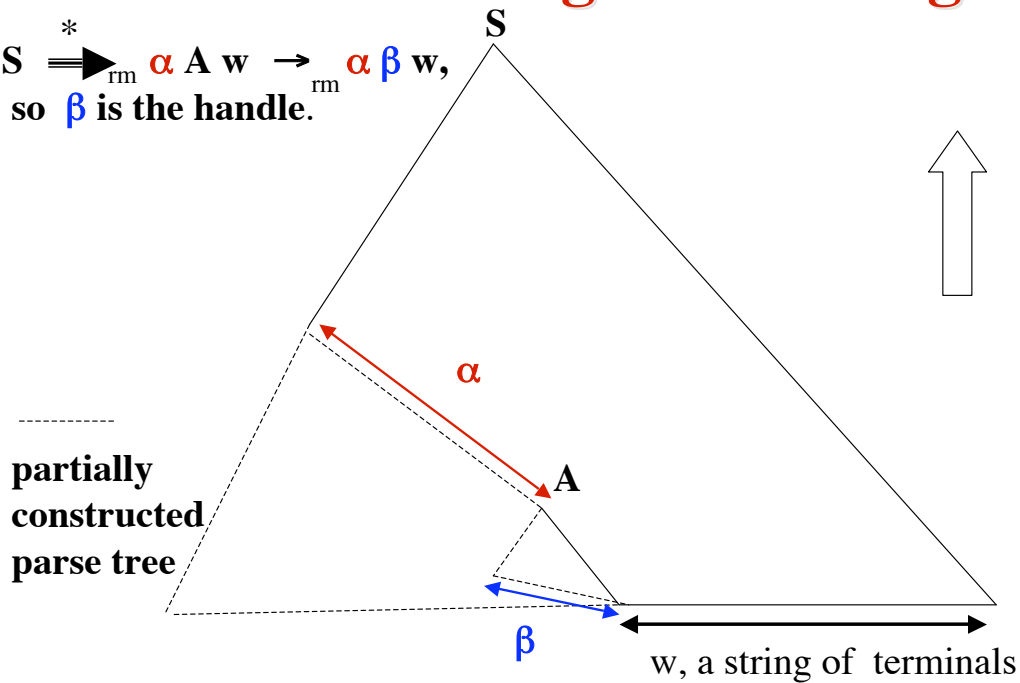
## A View During TD Parsing

**S derives  $\alpha \gamma$ ,**  
**a string of terminals;**  
**X is nonterminal at top**  
**of stack, X derives  $\gamma$ ;**  
**Initially**  
**X==S,  $\alpha == \epsilon$ ,  $\gamma$  is input**



# A View During BU Parsing

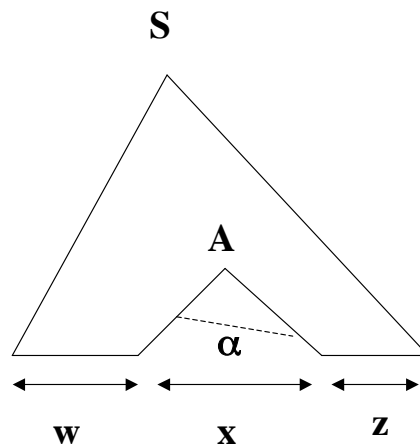
$S \xRightarrow{*}_{rm} \alpha A w \rightarrow_{rm} \alpha \beta w$ ,  
so  $\beta$  is the handle.



Parsing © BGR, Fall2005

11

# Intuitive Comparison



**LR(k)** can recognize  $A \rightarrow \alpha$  knowing  $w$ ,  $x$ , and  $\text{First}_k(z)$ .

**LL(k)** can recognize  $A \rightarrow \alpha$  knowing only  $w$  and  $\text{First}_k(x)$ .

Therefore, the set of languages recognizable by LR(k) contain those recognizable by LL(k).

Parsing © BGR, Fall2005

12

# BU Parsing (Shift-Reduce)

**Handle** - part of sentential form last added in a rightmost derivation.

BU parsing as “*handle hunting*”

- |                           |
|---------------------------|
| (1) $S \rightarrow E$     |
| (2) $E \rightarrow E + T$ |
| (3) $E \rightarrow T$     |
| (4) $T \rightarrow id$    |

Rightmost derivation of  $a+b+c$ , handles in red

- |                            |
|----------------------------|
| $S \rightarrow E$          |
| $\rightarrow E + T$        |
| $\rightarrow E + id$       |
| $\rightarrow E + T + id$   |
| $\rightarrow E + id + id$  |
| $\rightarrow T + id + id$  |
| $\rightarrow id + id + id$ |

## Shift-Reduce Parser, Example

Actions: **shift, reduce, accept, error**

| <u>Stack</u> | <u>Input</u>       | <u>Action</u> |
|--------------|--------------------|---------------|
| \$           | id1 + id2 + id3 \$ | shift         |
| \$ id1       | + id2 + id3 \$     | reduce (4)    |
| \$ T         | + id2 + id3 \$     | reduce (3)    |
| \$ E         | + id2 + id3 \$     | shift         |
| \$ E +       | id2 + id3 \$       | shift         |
| \$ E + id2   | + id3 \$           | reduce(4)     |
| \$ E + T     | + id3 \$           | reduce (2)    |
| \$ E         | + id3 \$           | shift         |
| \$ E +       | id3 \$             | shift         |
| \$ E + id3   | \$                 | reduce (4)    |
| \$ E + T     | \$                 | reduce(2)     |
| \$ E         | \$                 | reduce (1)    |
| \$ S         | \$                 | <b>accept</b> |

- |                           |
|---------------------------|
| (1) $S \rightarrow E$     |
| (2) $E \rightarrow E + T$ |
| (3) $E \rightarrow T$     |
| (4) $T \rightarrow id$    |

# Problems

## Shift-reduce conflicts

$S \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid \text{other}$

On stack: if E then S

Input: else

Should shift trying for 2nd alternative or reduce by first rule?

## Reduce-reduce conflicts

if  $A \rightarrow \alpha$  and  $B \rightarrow \alpha$  both in grammar

When  $\alpha$  on stack, how know which production to choose?

# Predictive Parsing

- Top Down: LL(k), Bottom Up: LR(k)
- Avoids backtracking while parsing by using lookahead into input
- NO cases where more than 1 action possible



# LR(k)

- Left to right scan parsing does a rightmost derivation in reverse, using **k** symbols of lookahead into input
- Three flavors
  - Simple LR, **SLR(1)**
    - Cheap
    - Doesn't always work
  - **LR**
    - Most powerful
    - Most expensive

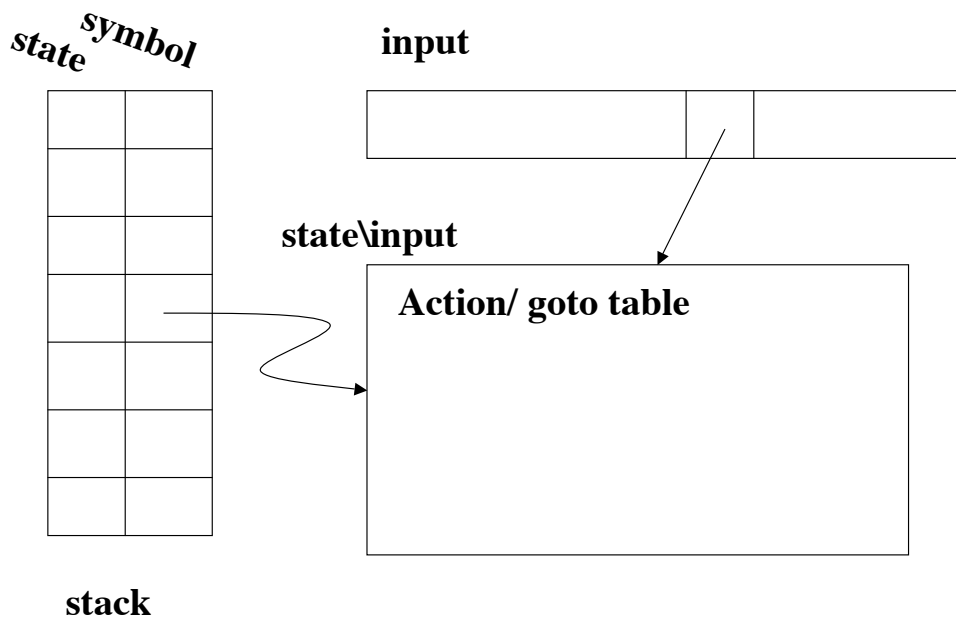
# LR(k)

- **LALR**
  - Intermediate in cost and power
- All SLR(1) languages are also LR(1), but parsers generated by corresponding grammars for the same language will differ in size.
- **LR(k)** catches syntax errors as early as possible in a left-to-right scan of the input
- Covers most programming languages

# LR Parsing

- **DFA is embedded in parser which is a PDA**
- **( $\text{top}_{\text{stack}}$  ,  $\text{input\_symbol}$ ) accesses a particular entry in the parser table**
  - **Shift to state  $s$**
  - **Reduce by  $A \rightarrow \beta$**
  - **Accept**
  - **Error**
- **Goto: (state,  $\text{top}_{\text{stack}}$ )  $\rightarrow$  state**

## LR Parser



# LR Parsing

- **Viable prefix** - set of prefixes of right sentential forms which can appear on a stack of a shift/reduce parser
  - Prefix of right sentential form that doesn't contain symbols beyond the handle
- **Goto** function is transition function of DFA that recognizes viable prefixes of the grammar
- *Idea*: continue to stack inputs until have **handle** on top of stack and then reduce

## Building an SLR Parser

- Need states, goto's, Follow sets
- **Item** - rule with embedded dot  
 $S \rightarrow . E$
- **Closure of item I**  
 $I \cup \{B \rightarrow . \gamma, \text{ if } A \rightarrow \alpha . B \beta \text{ in } I\}$
- States built from items and their closures

## Example - States

|                       |
|-----------------------|
| $S \rightarrow E$     |
| $E \rightarrow E + T$ |
| $E \rightarrow T$     |
| $T \rightarrow id$    |

$I_0: S \rightarrow . E$

|                         |
|-------------------------|
| $E \rightarrow . E + T$ |
| $E \rightarrow . T$     |
| $T \rightarrow . id$    |

Closure of  
 $S \rightarrow . E$

$I_1: S \rightarrow E .$   
 $E \rightarrow E . + T$

$I_2: E \rightarrow T .$

$I_3: T \rightarrow id .$

$I_4: E \rightarrow E + . T$   
 $T \rightarrow . id$

$I_5: E \rightarrow E + T .$

## Example - Goto's + Follow sets

goto (0, E) = 1

goto (0, id) = 3

goto (0, T) = 2

goto (1, +) = 4

goto (4, T) = 5

goto (4, id) = 3

goto ({set of items}, X) =

closure  $\{[A \rightarrow \alpha X . \beta] \mid$

$[A \rightarrow \alpha . X \beta] \text{ in \{set of items\}}\}$

where X is a terminal or nonterminal

---

Follow(S) = {\$}

Follow(E) = Follow(T) = {+, \$}

# Example - Parser Table

si, shift to state I; r(j) reduce by rule j

States\ inputs:

goto's

|   | <i>id</i> | +    | \$            | E | T |
|---|-----------|------|---------------|---|---|
| 0 | s3        |      |               | 1 | 2 |
| 1 |           | s4   | <b>accept</b> |   |   |
| 2 |           | r(3) | r(3)          |   |   |
| 3 |           | r(4) | r(4)          |   |   |
| 4 | s3        |      |               |   | 5 |
| 5 |           | r(2) | r(2)          |   |   |

## Example

| <u>Stack</u>           | <u>input</u>               | <u>action</u>   |
|------------------------|----------------------------|-----------------|
| 0                      | <i>id1</i> + <i>id2</i> \$ | s3              |
| 0 <i>id1</i> 3         | + <i>id2</i> \$            | r(4), goto on T |
| 0 T 2                  | + <i>id2</i> \$            | r(3), goto on E |
| 0 E 1                  | + <i>id2</i> \$            | s4              |
| 0 E 1 + 4              | <i>id2</i> \$              | s3              |
| 0 E 1 + 4 <i>id2</i> 3 | \$                         | r(4), goto on T |
| 0 E 1 + 4 T 5          | \$                         | r(2), goto on E |
| 0 E 1                  | \$                         | <b>accept</b>   |

# SLR(1) Parser Rules

- If  $A \rightarrow \alpha . a \beta$  is in state  $I_j$  and  $\text{goto}(I_j, a)$  is  $I_r$  then  $(I_j, a)$  transitions by shift  $r$  (sr)
- If  $A \rightarrow \alpha .$  is in state  $I_j$ , set action  $[j, a]$  to reduce  $A \rightarrow \alpha$  for all  $a$  in  $\text{Follow}(A)$ 
  - Note:  $A \neq S$
- If  $S \rightarrow E .$  in  $I_j$ , action  $(j, \$)$  is **accept**
- Any table entry not defined is error.

## Problems

- **Shift-reduce** conflicts happen when  $Ab$  can occur in some sentential form and  $b \in \text{Follow}(A)$ .

|  |   |
|--|---|
| $S \rightarrow L = R$<br>$S \rightarrow R$<br>$L \rightarrow * R$<br>$L \rightarrow id$<br>$R \rightarrow L$ | $I_0:$ $S \rightarrow . L = R$<br>$S \rightarrow . R$<br>$R \rightarrow . L$<br>$L \rightarrow . * R$<br>$L \rightarrow . id$ |
|  | $I_1:$ $S \rightarrow L . = R$ (1)<br>$R \rightarrow L .$ (2)   |

In  $I_1$  *shift* when see = in input(item 1); *reduce* on = because = in  $\text{Follow}(R)$  (item 2). Note:  $S \rightarrow L = R \rightarrow * R = R \dots$ , but this is not a rightmost derivation!

## Problems, cont.

Can see that rightmost derivation is:

$$S \rightarrow L = R \rightarrow L = L \rightarrow L = id \rightarrow *R = id \rightarrow \\ *L = id \rightarrow *id = id$$

Therefore, should reduce \*R to L when see =,  
not shift in order to get \*R onto the stack.

*Problem is that we can't distinguish those Follow  
elements corresponding to a rightmost derivation  
in a specific context.*

## Nomenclature in ASU

- An item  $[ A \rightarrow \beta . \gamma ]$  is *valid* for *viable prefix*  $\alpha \beta$  if  $S \xRightarrow[\text{rm}]{*} \alpha A w \xRightarrow[\text{rm}]{} \alpha \beta \gamma w$ .
  - Means can continue towards accumulating an handle on the stack by shifting
  - Previously, shift would have changed viable prefix \*R to nonviable prefix \*R=
- If  $I$  is set of items valid for *viable prefix*  $\beta$  then  $\text{goto}(I, X)$  is set of items valid for *viable prefix*  $\beta X$  where  $X$  is terminal or nonterminal

# LR Parsing

- LR items include a **lookahead symbol**, (into the input) which helps in conflict resolution
- Need new closure rule:
  - For  $[A \rightarrow \alpha . B \gamma, a]$  item add  $[B \rightarrow . \delta, b]$  for every  $b$  in  $\text{First}(\gamma a)$ .

## Example

- $I_0$ :
- $S \rightarrow . E, \$$  - initial item
  - $E \rightarrow . E + T, \$$  - closure initial item
  - $E \rightarrow . T, \$$
  - $E \rightarrow . E + T, +$  - closure 1st red item
  - $E \rightarrow . T, +$
  - $T \rightarrow . id, \$$  - closure 2nd red item
  - $T \rightarrow . id, +$  - closure 2nd blue item

Will write these in more compact form by combining lookaheads.

For  $[A \rightarrow \alpha . B \gamma, a]$  item add  $[B \rightarrow . \delta, b]$   
for every  $b$  in  $\text{First}(\gamma a)$ .



# Example, LR(1) Parser

**I<sub>0</sub>:** S → .E, \$

E → .E + T, \$/+

E → .T, \$/+

T → .id, +/\$

**I<sub>4</sub>:** [goto(I<sub>1</sub>, +)]

E → E + .T, \$/+

T → .id, \$/+

**I<sub>5</sub>:** [goto(I<sub>4</sub>, T)]

E → E + T ., \$/+

**I<sub>1</sub>:** [goto(I<sub>0</sub>, E)]

S → E ., \$

S → E . + T, \$/+

**I<sub>2</sub>:** [goto(I<sub>0</sub>, T)]

E → T ., \$/+

**I<sub>3</sub>:** [goto(I<sub>0</sub>, id)]

T → id ., \$/+

## LR(1) Parser

- Reduce based on lookaheads in item which are a subset of Follow set
- Rules similar to SLR
  - Shift in I<sub>k</sub>, [A → α . a β, b], goto (I<sub>k</sub>, a) = I<sub>j</sub>
  - Reduce [A → α ., b] reduce α to A on b
  - Accept [S → E., \$], accept on \$

# LALR Parsing

- **Idea:** merge all states with common first components in their LR(1) items
- **Implementation problem:** need to reduce number of states to get smaller parser table
- **Reduced size parser will perform**
  - Same as LR on correct inputs (will be parsed by LALR)
  - On incorrect inputs, LR may find error faster; LALR will never do an incorrect shift but may do some wrong reductions

# LALR Parsing

- **Conceptually, build LALR(1) parser from LR(1) parser**
  - Merge all states with same first components
  - Union all goto's of these merged states (goto's are independent of second components)
- **Correctness of conceptual derivation**
  - Can never produce a **shift-reduce conflict** or else  $[A \rightarrow \alpha . , a]$  and  $[B \rightarrow \beta . a \gamma , b]$  existed in some LR(1) state

# LALR, cont.

– But can create **reduce-reduce** conflicts

State 1

[A → c . , d]

[B → c . , e]

After merge:

[A → c . , d/e]

[B → c . , e/d]

State 2

[A → c . , e]

[B → c . , d]

## Handles (Informal)

*Show the handle always remains on the top of the stack in shift-reduce parsing*

Assume  $\gamma B w$  is a right sentential form with  $\gamma B$  on the stack.

1. Assume handle includes  $B$ .

$$S \xRightarrow{*}_{rm} \alpha_1 A \alpha_2 \xRightarrow{rm} \alpha_1 \boxed{\beta_1 B \beta_2} \alpha_2$$

$\gamma$ 
 $w$

where  $\beta_1$  or  $\beta_2$  can be  $\epsilon$ .

# Handles (Informal)

2. Assume handle included in  $w$ .

$$S \xRightarrow[\text{rm}]{*} \alpha_1 B \alpha_2 A \alpha_3 \xRightarrow[\text{rm}]{*} \alpha_1 B \alpha_2 \boxed{\beta} \alpha_3$$

$\gamma$ 
handle  
 $w$

3. Assume handle included in  $\gamma$ .

$$S \xRightarrow[\text{rm}]{*} \alpha_1 A \alpha_2 \xRightarrow[\text{rm}]{*} \alpha_1 \boxed{\beta} B w$$

$\gamma$ 
handle

But if  $\beta$  is handle, then A isn't rightmost nonterminal, B is. This is a contradiction!

# Handles, More Formally

- $A \rightarrow \beta_1 . \beta_2$  valid for viable prefix  $\alpha \beta_1$   
means  $\exists S \xRightarrow[\text{rm}]{*} \alpha A w \xRightarrow[\text{rm}]{*} \alpha \beta_1 \beta_2 w$   
If  $\beta_2 = \epsilon$  then should reduce by  $A \rightarrow \beta_1$   
If  $\beta_2 \neq \epsilon$  then should shift
- Note can have two valid items indicating different actions for same viable prefix  
 $A \rightarrow \beta_1 . \beta_2$  and  $A \rightarrow \beta_1$ . Lookahead chooses which action is taken

## Handles, More Formally

- **Previous argument shows if  $A \rightarrow \beta \cdot$  is valid item for viable prefix  $\alpha \beta$  then  $\alpha A$  is viable prefix (i.e., *we needn't rescan the parse after a reduction*)**

## Ambiguous Grammars

- **Used to build compact parse trees**
  - **Get rid of useless nonterminal to nonterminal productions (e.g.,  $S \rightarrow E \rightarrow T$ )**
- **Conflicts resolvable through desired properties of operators (e.g., precedence)**
- **Generate smaller parsers**
  - **Example of expression grammar**

# Example

$S \rightarrow E$   
 $E \rightarrow E + E$   
 $E \rightarrow id$

$I_0: S \rightarrow \cdot E$   
 $E \rightarrow \cdot E + E$   
 $E \rightarrow \cdot id$

$I_1: \text{goto}(I_0, E)$   
 $S \rightarrow E \cdot$   
 $E \rightarrow E \cdot + E$

$I_2: \text{goto}(I_1, +)$   
 $E \rightarrow E + \cdot E$   
 $E \rightarrow \cdot E + E$   
 $E \rightarrow \cdot id$

$I_3: \text{goto}(I_2, E)$   
 $E \rightarrow E + E \cdot$   
 $E \rightarrow E \cdot + E$

$I_4: \text{goto}(I_0, id)$   
 $E \rightarrow id \cdot$

(reduce on + in Follow(E)  
 shift on +)

Choose reduce action making + left associative; can resolve operator precedences the same way (e.g., + versus \*)

# Grammar Classification

