# LR Parsing
# LALR Parser Generators

# Outline

- Review of bottom-up parsing

- Computing the parsing DFA

- Using parser generators

# Bottom-up Parsing (Review)

- A bottom-up parser rewrites the input string to the start symbol.

- The state of the parser is described as:

$$\alpha \mid \gamma$$

- $\alpha$ is a stack of terminals and non-terminals;
- $\gamma$ is the string of terminals not yet examined.

- Initially: $\mid x_1 x_2 \ldots x_n$

# The Shift and Reduce Actions (Review)

Recall the CFG:  $E \rightarrow E + (E) \mid int$

A bottom-up parser uses two kinds of actions:

- <u>Shift</u> pushes a terminal from input on the stack

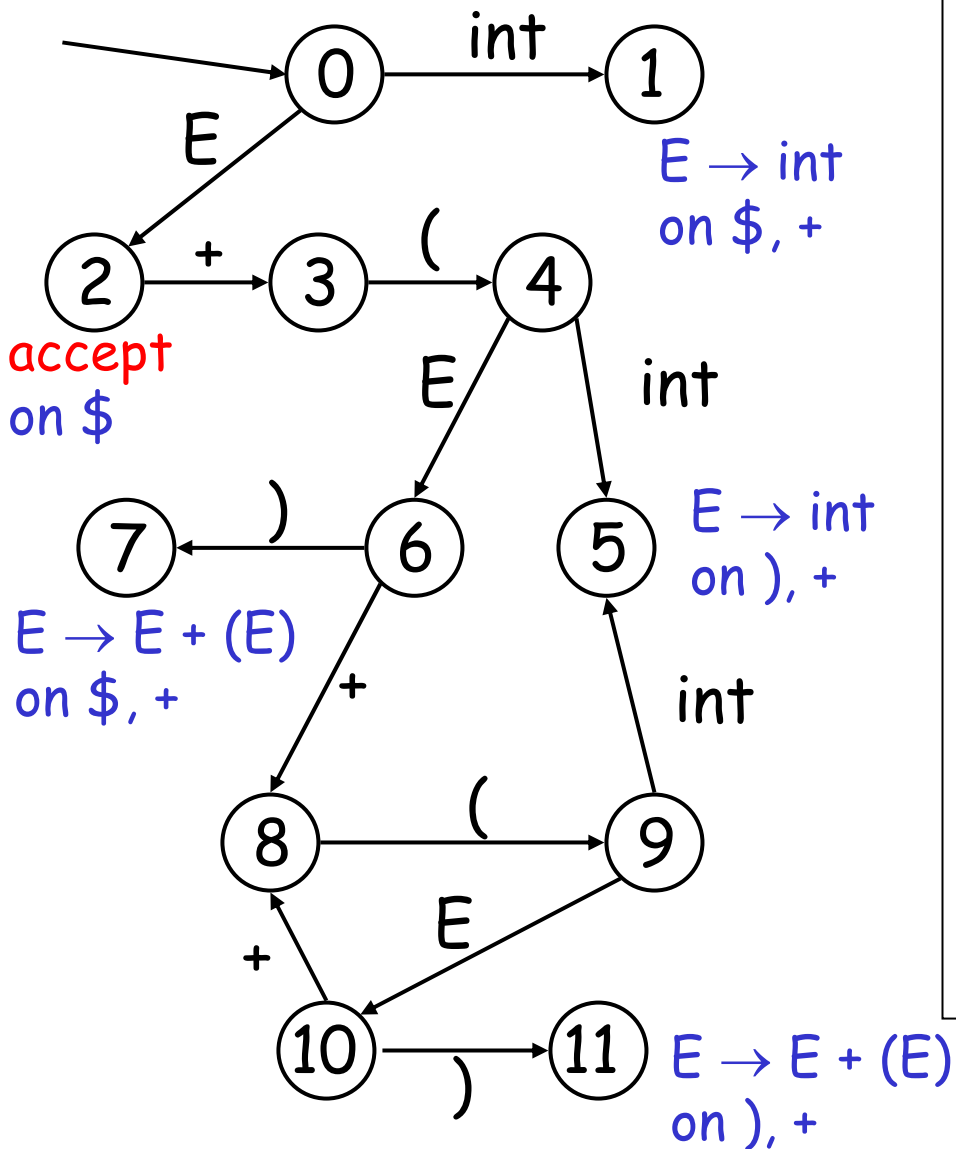$$E + ( \mid int ) \Rightarrow E + ( int \mid )$$

- <u>Reduce</u> pops 0 or more symbols off of the stack (production RHS) and pushes a non-terminal on the stack (production LHS)

$$E + (\underline{E + ( E )} \mid ) \Rightarrow E + ( \underline{E} \mid )$$

# Key Issue: When to Shift or Reduce?

- Idea: use a deterministic finite automaton (DFA) to decide when to shift or reduce
  - The input is the stack
  - The language consists of terminals and non-terminals

- We run the DFA on the stack and we examine the resulting state X and the token tok after I
  - If X has a transition labeled tok then <u>shift</u>
  - If X is labeled with "$A \rightarrow \beta$ on tok" then <u>reduce</u>
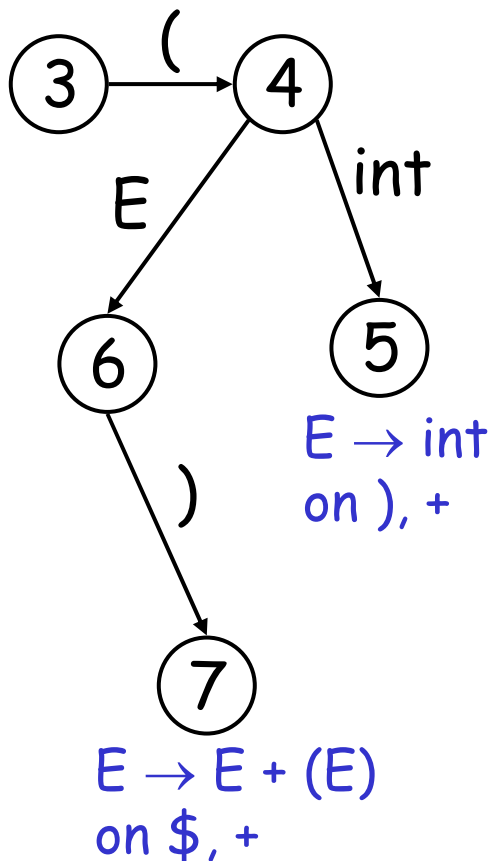
# LR(1) Parsing: An Example



| | |
|---|---|
| I int + (int) + (int)$ | shift |
| int I + (int) + (int)$ | E → int |
| E I + (int) + (int)$ | shift (x3) |
| E + (int I ) + (int)$ | E → int |
| E + (E I ) + (int)$ | shift |
| E + (E) I + (int)$ | E → E+(E) |
| E I + (int)$ | shift (x3) |
| E + (int I )$ | E → int |
| E + (E I )$ | shift |
| E + (E) I $ | E → E+(E) |
| E I $ | accept |

# Representing the DFA

- Parsers represent the DFA as a 2D table.
    (Recall table-driven lexical analysis.)
- Lines correspond to DFA states.
- Columns correspond to terminals and non-terminals.
- Typically columns are split into:
    - Those for terminals: the action table.
    - Those for non-terminals: the goto table.

# Representing the DFA: Example

The table for a fragment of our DFA:



|     | int | + | ( | ) | $ | E |
|-----|-----|---|---|---|---|---|
| ... |     |   |   |   |   |   |
| 3   |     |   | s4 |   |   |   |
| 4   | s5  |   |   |   |   | g6 |
| 5   |     | $r_{E \to int}$ |   | $r_{E \to int}$ |   |   |
| 6   |     | s8 |   | s7 |   |   |
| 7   |     | $r_{E \to E+(E)}$ |   |   | $r_{E \to E+(E)}$ |   |
| ... |     |   |   |   |   |   |

$E \to int$
on ), +

$E \to E + (E)$
on $, +

sk is shift and goto state k
$r_{X \to \alpha}$ is reduce
gk is goto state k

8

# The LR Parsing Algorithm

- After a shift or reduce action we rerun the DFA on the entire stack
  - This is wasteful, since most of the work is repeated

- To avoid this, we remember for each stack element on which state it brings the DFA.

- LR parser maintains a stack

$$\langle\, sym_1, state_1 \,\rangle \ldots \langle\, sym_n, state_n \,\rangle$$

$state_k$ is the final state of the DFA on $sym_1 \ldots sym_k$

# The LR Parsing Algorithm

```
let I = w$ be initial input
let j = 0
let DFA state 0 be the start state
let stack = ⟨ dummy, 0 ⟩
  repeat
    case action[top_state(stack), I[j]] of
      shift k: push ⟨ I[j++], k ⟩
      reduce X → A:
          pop |A| pairs,
          push ⟨ X, goto[top_state(stack), X] ⟩
      accept: halt normally
      error: halt and report error
```

# Key Issue: How is the DFA Constructed?

- The stack describes the context of the parse:
  - What non-terminal we are looking for.
  - What production RHS we are looking for.
  - What we have seen so far from the RHS.

- Each DFA state describes several such contexts.

  E.g., when we are looking for non-terminal E, we might be looking either for an int or an E + (E) RHS.

# LR(0) Items

- An <u>LR(0) item</u> is a production with a "I" somewhere on the RHS.

- The LR(0) items for T → (E) are
  - T → I (E)
  - T → ( I E)
  - T → (E I )
  - T → (E) I

- The only LR(0) item for X → ε is X → I

# LR(0) Items: Intuition

- An item $[X \rightarrow \alpha \mid \beta]$ says that the parser:
  - is looking for an $X$
  - has an $\alpha$ on top of the stack
  - expects to find a string derived from $\beta$ next in the input.

- Notes:
  - $[X \rightarrow \alpha \mid a\beta]$ means that $a$ should follow.
    - Then we can shift it and still have a viable prefix.
  - $[X \rightarrow \alpha \mid]$ means that we could reduce $X$.
    - But this is not always a good idea !

# LR(1) Items

- An <u>LR(1) item</u> is a pair:

$$X \rightarrow \alpha \mid \beta, \ a$$

  - $X \rightarrow \alpha\beta$ is a production.
  - $a$ is a terminal (the lookahead terminal).
  - LR(1) means 1 lookahead terminal.
- $[X \rightarrow \alpha \mid \beta, a]$ describes a context of the parser.
  - We are trying to find an $X$ followed by an $a$.
  - We have (at least) $\alpha$ already on top of the stack.
  - Thus we need to see next a prefix derived from $\beta a$.

# Note

- The symbol **|** was used before to separate the stack from the rest of input:

  $\alpha$ **|** $\gamma$, where $\alpha$ is the stack and $\gamma$ is the remaining string of terminals.

- In items, **|** is used to mark a prefix of a production RHS:

  $$X \rightarrow \alpha \,|\, \beta, \quad a$$

  – Here $\beta$ might contain non-terminals as well.

- In either case, the stack is on the left of **|**

# Convention

- We add to our grammar a fresh new start symbol $S$ and a production $S \rightarrow E$
  - Where $E$ is the old start symbol.

- The initial parsing context contains:

$$S \rightarrow \mid E \ , \ \$$$

  - Trying to find an $S$ as a string derived from $E\$$
  - The stack is empty.

# LR(1) Items (Cont.)

- In context containing

$$E \rightarrow E + \text{\textbar} ( E ) \quad , +$$

  - If ( follows then we can perform a shift to context containing

$$E \rightarrow E + ( \text{\textbar} E ) \quad , +$$

- In context containing

$$E \rightarrow E + ( E ) \text{\textbar} \quad , +$$

  - We can perform a reduction with $E \rightarrow E + ( E )$
  - But only if a + follows

# LR(1) Items (Cont.)

- Consider the item

$$E \rightarrow E + ( \, \mathbf{I} \, E ) \quad , \, +$$

- We expect a string derived from $E \, ) \, +$
- Our example has two productions for $E$

$$E \rightarrow int \quad and \quad E \rightarrow E + ( E )$$

- We describe this by extending the context with two more items:

$$E \rightarrow \mathbf{I} \, int \qquad \qquad , \, )$$
$$E \rightarrow \mathbf{I} \, E + ( E ) \quad , \, )$$

# The Closure Operation

- The operation of extending the context with items is called the closure operation.

**Closure**(Items) =
  repeat
    for each [X → α **|** Yβ, a] in Items
      for each production Y → γ
        for each b in First(βa)
          add [Y → **|** γ, b] to Items
  until Items is unchanged

# Constructing the Parsing DFA (1)

- Construct the start context:

  $E \rightarrow E + ( E ) \mid int$

  Closure({$S \rightarrow \mid E, \$$})

$$S \rightarrow \mid E \qquad , \$$$
$$E \rightarrow \mid E+(E), \$$$
$$E \rightarrow \mid int \qquad , \$$$
$$E \rightarrow \mid E+(E), +$$
$$E \rightarrow \mid int \qquad , +$$

- We abbreviate as:

$$S \rightarrow \mid E \qquad , \$$$
$$E \rightarrow \mid E+(E) \quad , \$/+$$
$$E \rightarrow \mid int \qquad , \$/+$$

# Constructing the Parsing DFA (2)

- A DFA state is a closed set of LR(1) items.

- The start state contains $[S \rightarrow \textbf{|}\ E\ , \$]$.

- A state that contains $[X \rightarrow \alpha\ \textbf{|}, b]$ is labeled with "reduce with $X \rightarrow \alpha$ on b".

- And now the transitions ...

# The DFA Transitions

- A state "State" that contains $[X \rightarrow \alpha \textcolor{red}{|} \gamma\beta, b]$ has a transition labeled $\gamma$ to a state that contains the items "**Transition**(State, $\gamma$)"
  - $\gamma$ can be a terminal or a non-terminal

**Transition**(State, $\gamma$)

   Items = $\varnothing$

   for each $[X \rightarrow \alpha \textcolor{red}{|} \gamma\beta, b]$ in State

      add $[X \rightarrow \alpha\gamma \textcolor{red}{|} \beta, b]$ to Items

   return Closure(Items)

# Constructing the Parsing DFA: Example



**0**
$S \rightarrow \mathbf{I}\, E \qquad , \$$
$E \rightarrow \mathbf{I}\, E+(E), \$/+$
$E \rightarrow \mathbf{I}\, int \qquad , \$/+$

**1**
$E \rightarrow int\, \mathbf{I}, \$/+$    $E \rightarrow int$ on $\$, +$

int

**2**
$S \rightarrow E\, \mathbf{I} \qquad , \$$
$E \rightarrow E\, \mathbf{I} +(E), \$/+$

accept on $\$$

E

**3**
$E \rightarrow E+ \mathbf{I}\, (E), \$/+$

+

(

**4**
$E \rightarrow E+( \mathbf{I}\, E)\, , \$/+$
$E \rightarrow \mathbf{I}\, E+(E)\, , )/+$
$E \rightarrow \mathbf{I}\, int \qquad , )/+$

E

int

**6**
$E \rightarrow E+(E\, \mathbf{I}\, )\, , \$/+$
$E \rightarrow E\, \mathbf{I} +(E)\, , )/+$

**5**
$E \rightarrow int\, \mathbf{I}\, , )/+$    $E \rightarrow int$ on $), +$

+

)

and so on…

23

# LR Parsing Tables: Notes

- Parsing tables (i.e., the DFA) can be constructed automatically for a CFG.

- But we still need to understand the construction to work with parser generators.
  - E.g., they report errors in terms of sets of items.

- What kind of errors can we expect?

# Shift/Reduce Conflicts

- If a DFA state contains both

    $[X \rightarrow \alpha \mid a\beta, b]$  and  $[Y \rightarrow \gamma \mid, a]$


- Then on input "a" we could either
  - Shift into state $[X \rightarrow \alpha a \mid \beta, b]$, or
  - Reduce with $Y \rightarrow \gamma$


- This is called a *shift-reduce conflict*

# Shift/Reduce Conflicts

- Typically due to ambiguities in the grammar.
- Classic example: the dangling else

    S → if E then S | if E then S else S | OTHER

- We will have a DFA state containing:

    [S → if E then S **I**,            else]
    [S → if E then S **I** else S,   x]

- If else follows then we can shift or reduce.
- Default (yacc, ML-yacc, bison, etc.) is to shift.
  - Default behavior is as needed in this case.

# More Shift/Reduce Conflicts

- Consider the ambiguous grammar:
  $$E \rightarrow E + E \mid E * E \mid int$$

- We will have the states containing:

  [E → E * **|** E,  +]             [E → E * E **|**,  +]

  [E → **|** E + E,  +]   ⇒$^E$  [E → E **|** + E,  +]

  …                                        …

- Again we have a shift/reduce on input +
  - We need to reduce (* binds more tightly than +)
  - Recall solution: declare the precedence of * and +

27

# More Shift/Reduce Conflicts

- In yacc declare precedence and associativity:

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

- Precedence of a rule = that of its last terminal.

  See yacc manual for ways to override this default.

- Resolve shift/reduce conflict with a <u>shift</u> if:
  - no precedence declared for either rule or terminal;
  - input terminal has higher precedence than the rule;
  - the precedences are the same and right associative.

# Using Precedence to Solve S/R Conflicts

- Back to our example:

  $[E \rightarrow E * \mathbf{|} E, \ +]$       $[E \rightarrow E * E \mathbf{|}, \ +]$

  $[E \rightarrow \mathbf{|} E + E, \ +] \Rightarrow^E$    $[E \rightarrow E \mathbf{|} + E, \ +]$

       …                         …

- We will choose reduce because precedence of rule $E \rightarrow E * E$ is higher than of terminal $+$ .

# Using Precedence to Solve S/R Conflicts

- Same grammar as before:
  $$E \rightarrow E + E \mid E * E \mid \text{int}$$

- We will also have the states:

  $[E \rightarrow E + | E, +]$    $[E \rightarrow E + E |, +]$

  $[E \rightarrow | E + E, +]$  $\Rightarrow^E$  $[E \rightarrow E | + E, +]$

  ...                        ...

- Now we also have a shift/reduce on input $+$

  - We will choose reduce because $E \rightarrow E + E$ and $+$ have the same precedence and $+$ is left-associative.

# Using Precedence to Solve S/R Conflicts

- Back to our dangling else example:

    [S → if E then S l,            else]

    [S → if E then S l else S,   x]

- Can eliminate conflict by declaring else having higher precedence than then.

- But this starts to look like "hacking the tables".

- Best to avoid overuse of precedence declarations or we will end with unexpected parse trees.

# Precedence Declarations Revisited

The term "precedence declaration" is misleading!

These declarations do not define precedence; instead, they define conflict resolutions.

  I.e., they instruct shift-reduce parsers to resolve conflicts in certain ways.

  These two are not quite the same!

# Reduce/Reduce Conflicts

- If a DFA state contains both

$$[X \rightarrow \alpha \mathbin{|}, a] \quad \text{and} \quad [Y \rightarrow \beta \mathbin{|}, a]$$

  – Then on input "a" we do not know which production to reduce.

- This is called a *reduce/reduce conflict*

# Reduce/Reduce Conflicts

- Usually due to gross ambiguity in the grammar.
- Ex. A grammar for a sequence of identifiers:

$$S \rightarrow \varepsilon \mid id \mid id\ S$$

- There are two parse trees for the string id

$$S \rightarrow id$$

$$S \rightarrow id\ S \rightarrow id$$

- How does this confuse the parser?
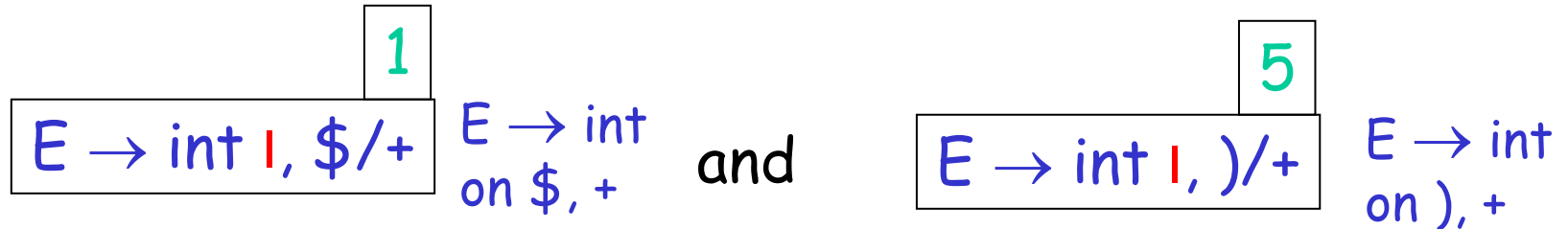
# More on Reduce/Reduce Conflicts

- Consider the states

$$[S' \rightarrow | S, \quad \$]$$
$$[S \rightarrow |, \qquad \$]$$
$$[S \rightarrow | \text{ id}, \quad \$]$$
$$[S \rightarrow | \text{ id } S, \quad \$]$$

$$\Rightarrow^{\text{id}}$$

$$[S \rightarrow \text{id } |, \qquad \$]$$
$$[S \rightarrow \text{id } | S, \quad \$]$$
$$[S \rightarrow |, \qquad \$]$$
$$[S \rightarrow | \text{ id}, \qquad \$]$$
$$[S \rightarrow | \text{ id } S, \quad \$]$$

- Reduce/reduce conflict on input $\$$

$$S' \rightarrow S \rightarrow \text{id}$$
$$S' \rightarrow S \rightarrow \text{id } S \rightarrow \text{id}$$

- Better to rewrite the grammar as: $S \rightarrow \varepsilon \mid \text{id } S$

# Using Parser Generators

- Parser generators automatically construct the parsing DFA given a CFG.

  - Use precedence declarations and default conventions to resolve conflicts.

  - The parser algorithm is the same for all grammars (and is provided as a library function).

- But most parser generators do not construct the DFA as described before.

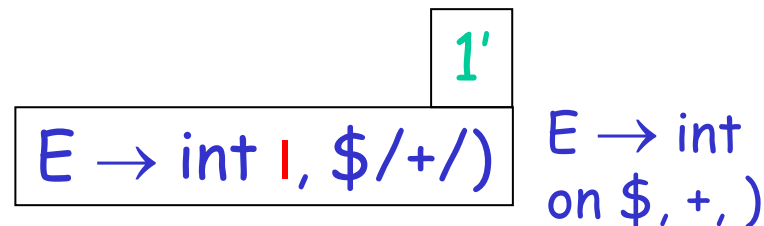  - Because the LR(1) parsing DFA has 1000s of states even for a simple language.

# LR(1) Parsing Tables are Big

- But many states are similar, e.g.



$$E \rightarrow int \;|\,, \$/+$$   $E \rightarrow int$ on $\$, +$   **and**   $E \rightarrow int \;|\,, )/+$   $E \rightarrow int$ on $), +$

- <u>Idea</u>: merge the DFA states whose items differ only in the lookahead tokens
  - We say that such states have the same core

- We obtain

$$E \rightarrow int \;|\,, \$/+/)$$   $E \rightarrow int$ on $\$, +, )$

# The Core of a Set of LR Items

**Definition**: The core of a set of LR items is the set of first components
- Without the lookahead terminals

- Example: the core of
$$\{[X \rightarrow \alpha \mid \beta, b], [Y \rightarrow \gamma \mid \delta, d]\}$$
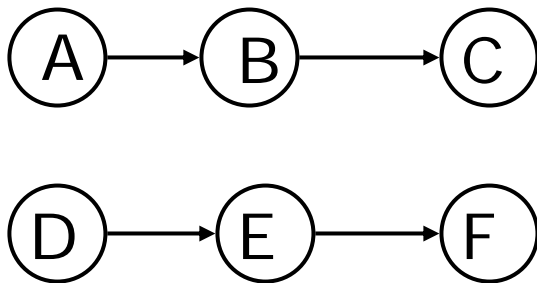is
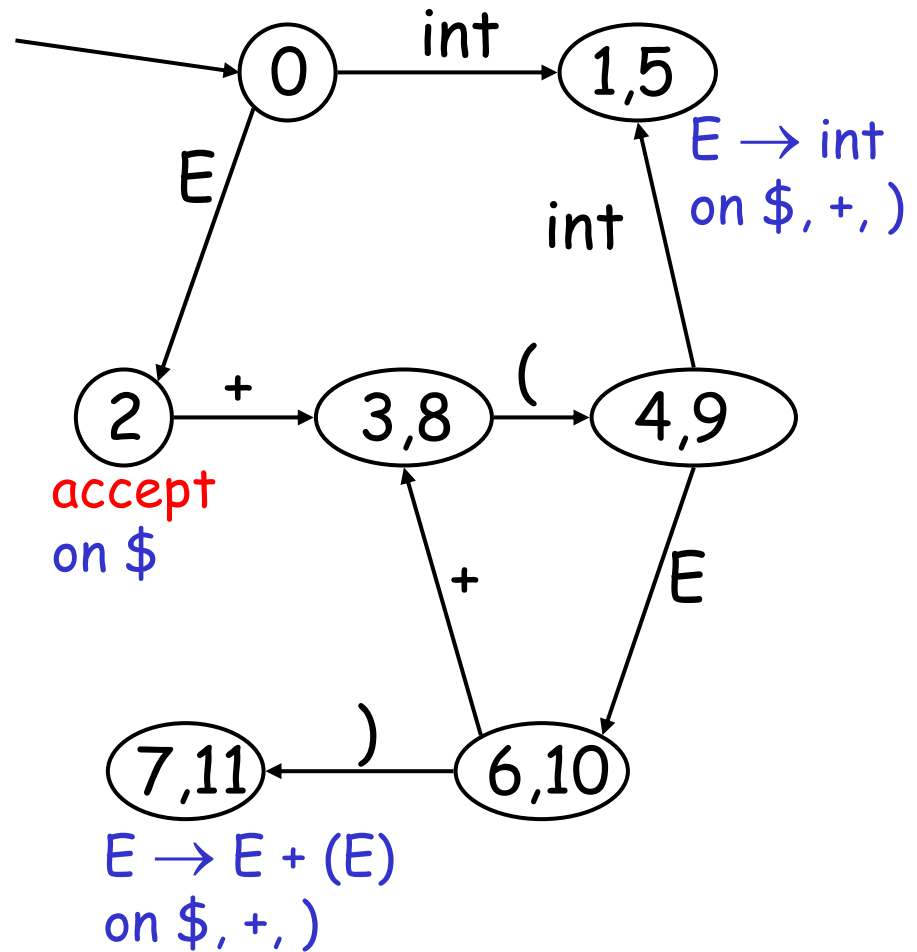$$\{X \rightarrow \alpha \mid \beta, \ Y \rightarrow \gamma \mid \delta\}$$

# LALR States

- Consider for example the LR(1) states
$$\{[X \rightarrow \alpha \mathbin{\text{\textbar}}, a], [Y \rightarrow \beta \mathbin{\text{\textbar}}, c]\}$$
$$\{[X \rightarrow \alpha \mathbin{\text{\textbar}}, b], [Y \rightarrow \beta \mathbin{\text{\textbar}}, d]\}$$

- They have the same core and can be merged
- The merged state contains:
$$\{[X \rightarrow \alpha \mathbin{\text{\textbar}}, a/b], [Y \rightarrow \beta \mathbin{\text{\textbar}}, c/d]\}$$

- These are called LALR(1) states
  - Stands for LookAhead LR
  - Typically 10 times fewer LALR(1) states than LR(1)

# A LALR(1) DFA

- Repeat until all states have distinct core
  - Choose two distinct states with same core
  - Merge the states by creating a new one with the union of all the items
  - Point edges from predecessors to new state
  - New state points to all the previous successors
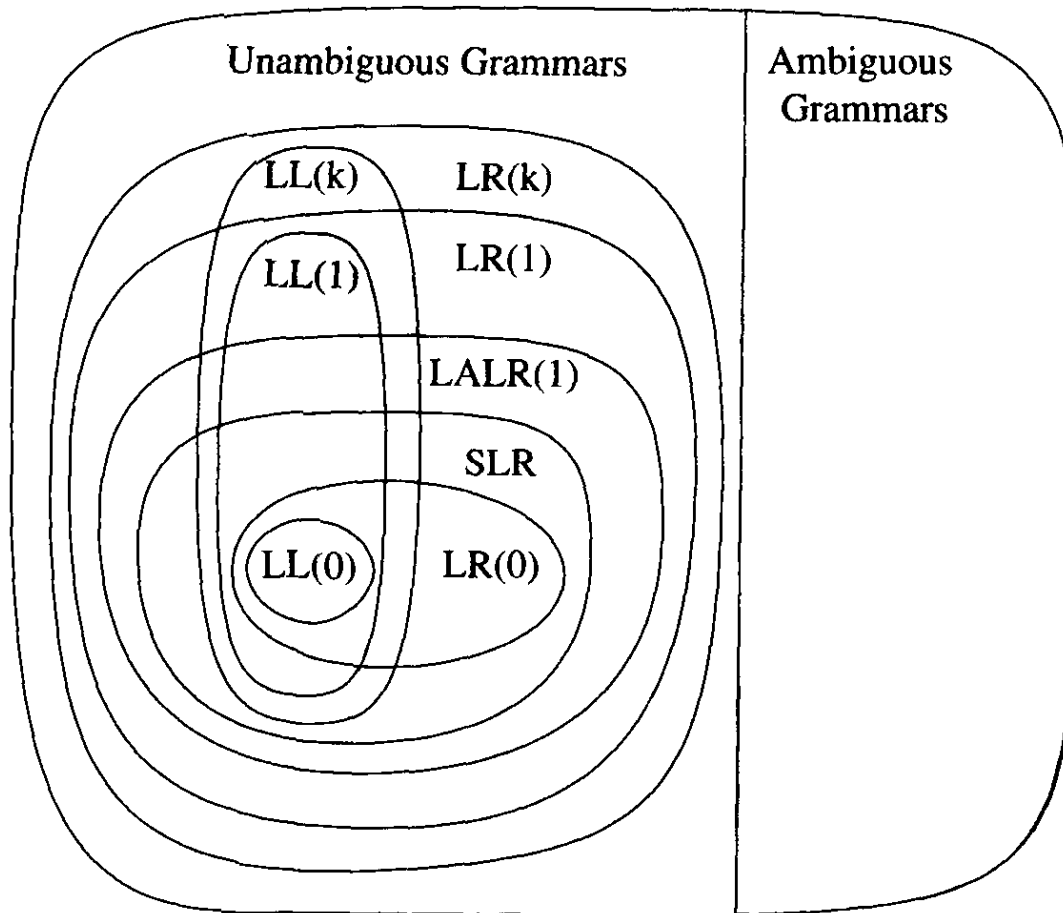
# Conversion LR(1) to LALR(1): Example.

# The LALR Parser Can Have Conflicts

- Consider for example the LR(1) states:

$$\{[X \rightarrow \alpha \cdot, a], [Y \rightarrow \beta \cdot, b]\}$$
$$\{[X \rightarrow \alpha \cdot, b], [Y \rightarrow \beta \cdot, a]\}$$

- And the merged LALR(1) state:

$$\{[X \rightarrow \alpha \cdot, a/b], [Y \rightarrow \beta \cdot, a/b]\}$$

- Has a <u>new</u> reduce/reduce conflict!

- In practice, such cases are rare.

# LALR vs. LR Parsing: Things to keep in mind

- LALR languages are not natural.
  - They are an "efficiency hack" on LR languages.

- Any reasonable programming language has a LALR(1) grammar.

- LALR(1) parsing has become a standard for programming languages and parser generators.

# A Hierarchy of Grammar Classes



From Andrew Appel, "Modern Compiler Implementation in ML"