November 29

# EECS 483:
# COMPILER CONSTRUCTION

# LR GRAMMARS

# Bottom-up Parsing  (LR Parsers)

- LR(k) parser:
  - <u>L</u>eft-to-right scanning
  - <u>R</u>ightmost derivation
  - k lookahead symbols

- LR grammars are more expressive than LL
  - Can handle left-recursive (and right recursive) grammars; virtually all programming languages
  - Easier to express programming language syntax (no left factoring)

- Technique:  "Shift-Reduce" parsers
  - Work bottom up instead of top down
  - Construct right-most derivation of a program in the grammar
  - Used by many parser generators (e.g. yacc, ocamlyacc, lalrpop, etc.)
  - Better error detection/recovery

# Top-down vs. Bottom up
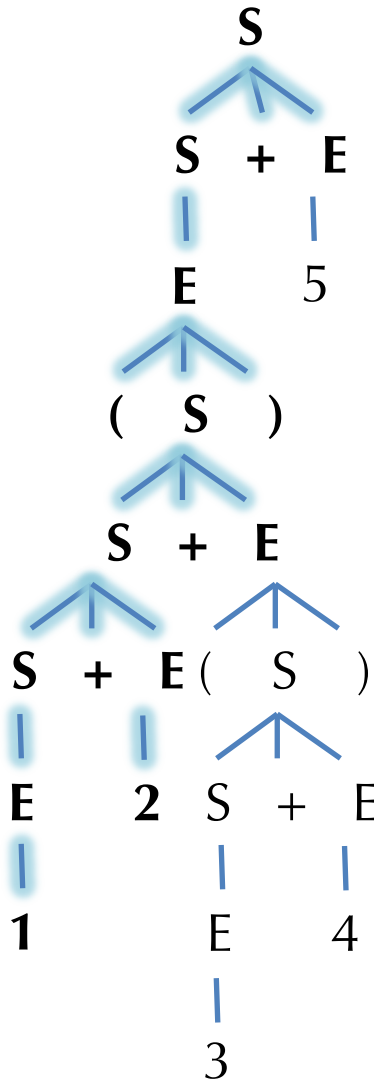
- Consider the left-recursive grammar:

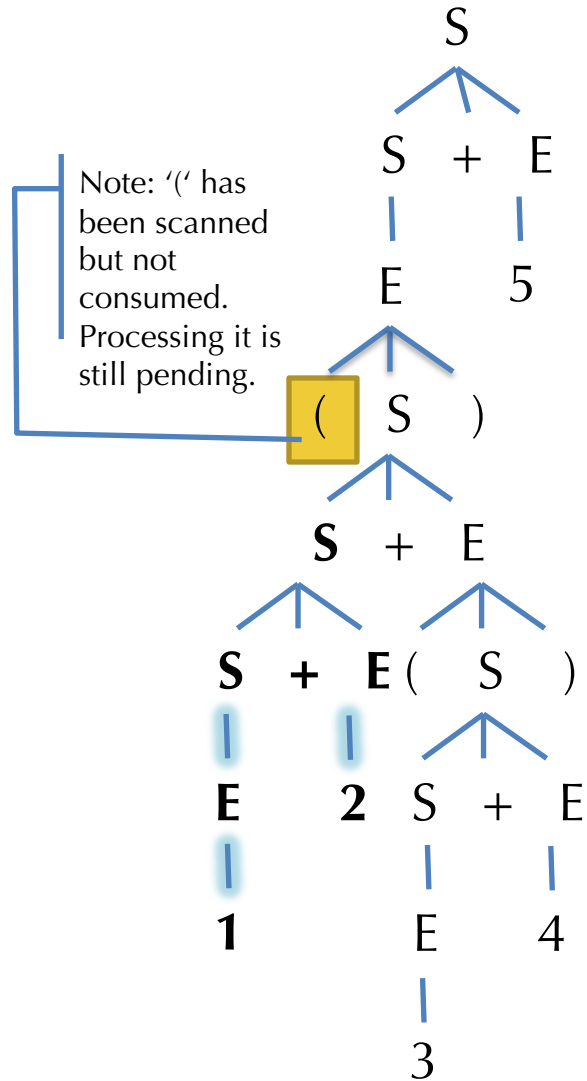  $S \longmapsto S + E \mid E$
  $E \longmapsto$ number $\mid ( S )$

- $(1 + 2 + (3 + 4)) + 5$

- What part of the tree must we know after scanning just "(1 + 2" ?

- In top-down, must be able to guess which productions to use…

Note: '(' has been scanned but not consumed. Processing it is still pending.

Top-down

Bottom-up

# Progress of Bottom-up Parsing

| Reductions | Scanned | Input Remaining |
|---|---|---|
| (1 + 2 + (3 + 4)) + 5 ↩ | | (1 + 2 + (3 + 4)) + 5 |
| (**E** + 2 + (3 + 4)) + 5 ↩ | ( | 1 + 2 + (3 + 4)) + 5 |
| (**S** + 2 + (3 + 4)) + 5 ↩ | (1 | + 2 + (3 + 4)) + 5 |
| (S + **E** + (3 + 4)) + 5 ↩ | (1 + 2 | + (3 + 4)) + 5 |
| (**S** + (3 + 4)) + 5 ↩ | (1 + 2 | + (3 + 4)) + 5 |
| (S + (**E** + 4)) + 5 ↩ | (1 + 2 + (3 | + 4)) + 5 |
| (S + (**S** + 4)) + 5 ↩ | (1 + 2 + (3 | + 4)) + 5 |
| (S + (S + **E**)) + 5 ↩ | (1 + 2 + (3 + 4 | )) + 5 |
| (S + (**S**)) + 5 ↩ | (1 + 2 + (3 + 4 | )) + 5 |
| (S + **E**) + 5 ↩ | (1 + 2 + (3 + 4) | ) + 5 |
| (**S**) + 5 ↩ | (1 + 2 + (3 + 4) | ) + 5 |
| **E** + 5 ↩ | (1 + 2 + (3 + 4)) | + 5 |
| **S** + 5 ↩ | (1 + 2 + (3 + 4)) | + 5 |
| S + **E** ↩ | (1 + 2 + (3 + 4)) + 5 | |
| S | | |

*Rightmost derivation*

S ↦ S + E | E
E ↦ number | ( S )

# Shift/Reduce Parsing

$$S \longmapsto S + E \;|\; E$$
$$E \longmapsto number \;|\; ( \, S \, )$$

- Parser state:
  - Stack of terminals and nonterminals.
  - Unconsumed input is a string of terminals
  - Current derivation step is    stack + input
- Parsing is a sequence of *shift* and *reduce* operations:
- Shift: move look-ahead token to the stack
- Reduce: Replace symbols $\gamma$ at top of stack with nonterminal X such that X $\longmapsto \gamma$ is a production.  (pop $\gamma$, push X)

| Stack | Input | Action |
|---|---|---|
|  | (1 + 2 + (3 + 4)) + 5 | shift ( |
| ( | 1 + 2 + (3 + 4)) + 5 | shift 1 |
| (1 | + 2 + (3 + 4)) + 5 | reduce: E $\longmapsto$ number |
| (E | + 2 + (3 + 4)) + 5 | reduce: S $\longmapsto$ E |
| (S | + 2 + (3 + 4)) + 5 | shift + |
| (S + | 2 + (3 + 4)) + 5 | shift 2 |
| (S + 2 | + (3 + 4)) + 5 | reduce: E $\longmapsto$ number |
| (S + E | + (3 + 4)) + 5 | reduce: S $\longmapsto$ S + E |
| (S | + (3 + 4)) + 5 | shift + |

# Shift/Reduce Parsing

- Parser state:
  - Stack of terminals and nonterminals.
  - Unconsumed input is a string of terminals
  - Current derivation step is     stack + input

$$S \longmapsto S + E \mid E$$
$$E \longmapsto number \mid ( \, S \, )$$

- Invariant: Stack plus input is a step in building the Rightmost derivation in reverse

| Stack | Input | Derivation steps |
|---|---|---|
| | $(1 + 2 + (3 + 4)) + 5$ | $(1 + 2 + (3 + 4)) + 5$ |
| ( | $1 + 2 + (3 + 4)) + 5$ | |
| (1 | $+ 2 + (3 + 4)) + 5$ | |
| (E | $+ 2 + (3 + 4)) + 5$ | $(\underline{E} + 2 + (3 + 4)) + 5$ |
| (S | $+ 2 + (3 + 4)) + 5$ | $(\underline{S} + 2 + (3 + 4)) + 5$ |
| (S + | $2 + (3 + 4)) + 5$ | |
| (S + 2 | $+ (3 + 4)) + 5$ | |
| (S + E | $+ (3 + 4)) + 5$ | $(S + \underline{E} + (3 + 4)) + 5$ |
| (S | $+ (3 + 4)) + 5$ | $(\underline{S} + (3 + 4)) + 5$ |

Rightmost derivation →

7

Simple LR parsing with no look ahead.

# LR(0) GRAMMARS

# LR Parser States

- Goal: know what set of reductions are legal at any given point.
- Idea: Summarize all possible stack prefixes α as a finite parser state.
  - Parser state is computed by a DFA that reads the stack σ.
  - Accept states of the DFA correspond to unique reductions that apply.

- Example: LR(0) parsing
  - **L**eft-to-right scanning, **R**ight-most derivation, **zero** look-ahead tokens
  - Too weak to handle many language grammars (e.g. the "sum" grammar)
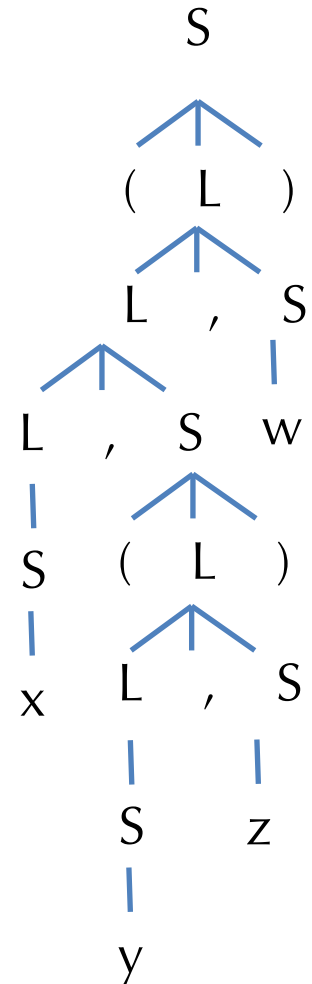  - But, helpful for understanding how the shift-reduce parser works.

# Example LR(0) Grammar: Tuples

- Example grammar for non-empty tuples and identifiers:

$$S \longmapsto ( L ) \;\mid\; id$$
$$L \longmapsto S \;\mid\; L , S$$

- Example strings:
  - x
  - (x,y)
  - ((((x))))
  - (x, (y, z), w)
  - (x, (y, (z, w)))

Parse tree for:
(x, (y, z), w)

# Shift/Reduce Parsing

$$S \longmapsto ( \, L \, ) \quad | \quad id$$
$$L \longmapsto S \quad | \quad L \, , \, S$$

- Parser state:
  - Stack of terminals and nonterminals.
  - Unconsumed input is a string of terminals
  - Current derivation step is      stack + input
- Parsing is a sequence of *shift* and *reduce* operations:
- Shift: move look-ahead token to the stack: e.g.

| Stack | Input | Action |
|-------|-------|--------|
|       | (x,  (y, z), w) | shift ( |
| (     | x,  (y, z), w) | shift x |

- Reduce: Replace symbols γ at top of stack with nonterminal X such that $X \longmapsto \gamma$ is a production.  (pop γ, push X): e.g.

| Stack | Input | Action |
|-------|-------|--------|
| (x    | ,  (y, z), w) | reduce $S \longmapsto id$ |
| (S    | ,  (y, z), w) | reduce $L \longmapsto S$ |

# Example Run

| Stack | Input | Action |
|---|---|---|
| | (x,  (y, z), w) | shift ( |
| ( | x,  (y, z), w | shift x |
| (x | ,  (y, z), w | reduce S $\mapsto$ id |
| (S | ,  (y, z), w | reduce L $\mapsto$ S |
| (L | ,  (y, z), w | shift , |
| (L, | (y, z), w | shift ( |
| (L, ( | y, z), w | shift y |
| (L, (y | , z), w | reduce S $\mapsto$ id |
| (L, (S | , z), w | reduce L $\mapsto$ S |
| (L, (L | , z), w | shift , |
| (L, (L, | z), w | shift z |
| (L, (L, z | ), w | reduce S $\mapsto$ id |
| (L, (L, S | ), w | reduce L $\mapsto$ L, S |
| (L, (L | ), w | shift ) |
| (L, (L) | , w | reduce S $\mapsto$ ( L ) |
| (L, S | , w | reduce L $\mapsto$ L, S |
| (L | , w | shift , |
| (L, | w) | shift w |
| (L, w | ) | reduce S $\mapsto$ id |
| (L, S | ) | reduce L $\mapsto$ L, S |
| (L | ) | shift ) |
| (L) | | reduce S $\mapsto$ ( L ) |
| S | | |

$$S \mapsto ( L ) \mid id$$
$$L \mapsto S \mid L , S$$

# Action Selection Problem

- Given a stack σ and a look-ahead symbol b, should the parser:
  - Shift b onto the stack (new stack is σb)
  - Reduce a production X ⟼ γ, assuming that σ = αγ  (new stack is αX)?

- Sometimes the parser can reduce but shouldn't
  - For example, X ⟼ ε can *always* be reduced
- Sometimes the stack can be reduced in different ways

- Main idea:  decide what to do based on a *prefix* α of the stack plus the look-ahead symbol.
  - The prefix α is different for different possible reductions since in productions X ⟼ γ and Y ⟼ β, γ and β might have different lengths.

- Main goal: know what set of reductions are legal at any point.
  - How do we keep track?

# LR(0) States

- An LR(0) *state* is a *set* of *items* keeping track of progress on possible upcoming reductions.
- An LR(0) *item* is a production from the language with an extra separator "." somewhere in the right-hand-side

$$S \longmapsto ( L ) \mid id$$
$$L \longmapsto S \mid L , S$$

- Example items:   $S \longmapsto .( L )$   or   $S \longmapsto (. L)$   or   $L \longmapsto S.$
- Intuition:
  - Stuff before the '.' is already on the stack (beginnings of possible γ's to be reduced)
  - Stuff after the '.' is what might be seen next
  - The prefixes α are represented by the state itself

# Constructing the DFA: Start state & Closure

- First step:  Add a new production
  $S' \mapsto S\$$  to the grammar

$$S' \mapsto S\$$$
$$S \mapsto (\ L\ )\ \ |\ \ id$$
$$L \mapsto S\ \ \ |\ \ \ L\ ,\ S$$

- Start state of the DFA =  empty stack, so it contains the item:
  $S' \mapsto .S\$$

- Closure of a state:
  - Adds items for all productions whose LHS nonterminal occurs in an item in the state just after the '.'
  - The added items have the '.' located at the beginning (no symbols for those items have been added to the stack yet)
  - Note that newly added items may cause yet more items to be added to the state… keep iterating until a *fixed point* is reached.

- Example:  CLOSURE($\{S' \mapsto .S\$\}$)  =  $\{S' \mapsto .S\$, S \mapsto .(L), S \mapsto .id\}$

- Resulting "closed state" contains the set of all possible productions that might be reduced next.

# Example: Constructing the DFA

$$S' \mapsto S\$$$
$$S \mapsto (\,L\,)\ \mid\ id$$
$$L \mapsto S\ \mid\ L\,,\,S$$

$$S' \mapsto .S\$$$

- First, we construct a state with the initial item $S' \mapsto .S\$$

# Example: Constructing the DFA

S' $\mapsto$ S$
S $\mapsto$ ( L )  |  id
L $\mapsto$ S  |  L , S

S' $\mapsto$ .S$
S $\mapsto$ .( L )
S $\mapsto$ .id

- Next, we take the closure of that state:
  CLOSURE({S' $\mapsto$ .S$}) = {S' $\mapsto$ .S$, S $\mapsto$ .( L ), S $\mapsto$ .id}

- In the set of items, the nonterminal S appears after the '.'
- So we add items for each S production in the grammar
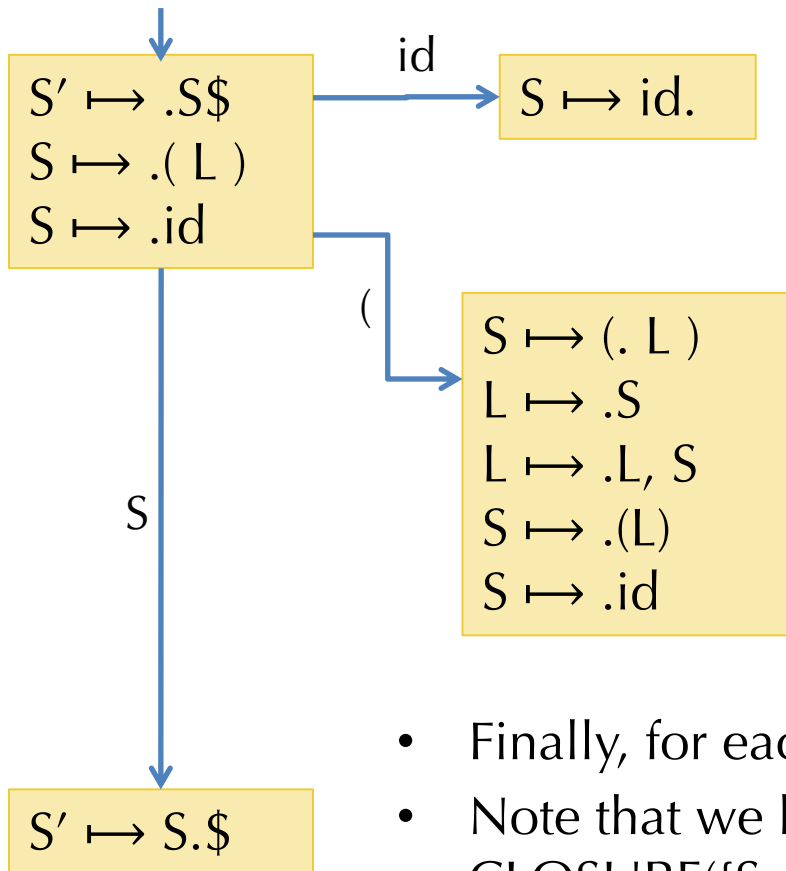
# Example: Constructing the DFA

$$S' \mapsto S\$$$
$$S \mapsto ( \ L \ ) \ | \ id$$
$$L \mapsto S \ | \ L \ , \ S$$

$S' \mapsto .S\$$
$S \mapsto .( \ L \ )$
$S \mapsto .id$

→ (id) → $S \mapsto id.$

→ (( ) → $S \mapsto (. \ L \ )$

→ (S) → $S' \mapsto S.\$$

- Next we add the transitions:
- First, we see what terminals and nonterminals can appear after the '.' in the source state.
  - Outgoing edges have those label.
- The target state (initially) includes all items from the source state that have the edge-label symbol after the '.', but we advance the '.'  (to simulate shifting the item onto the stack)
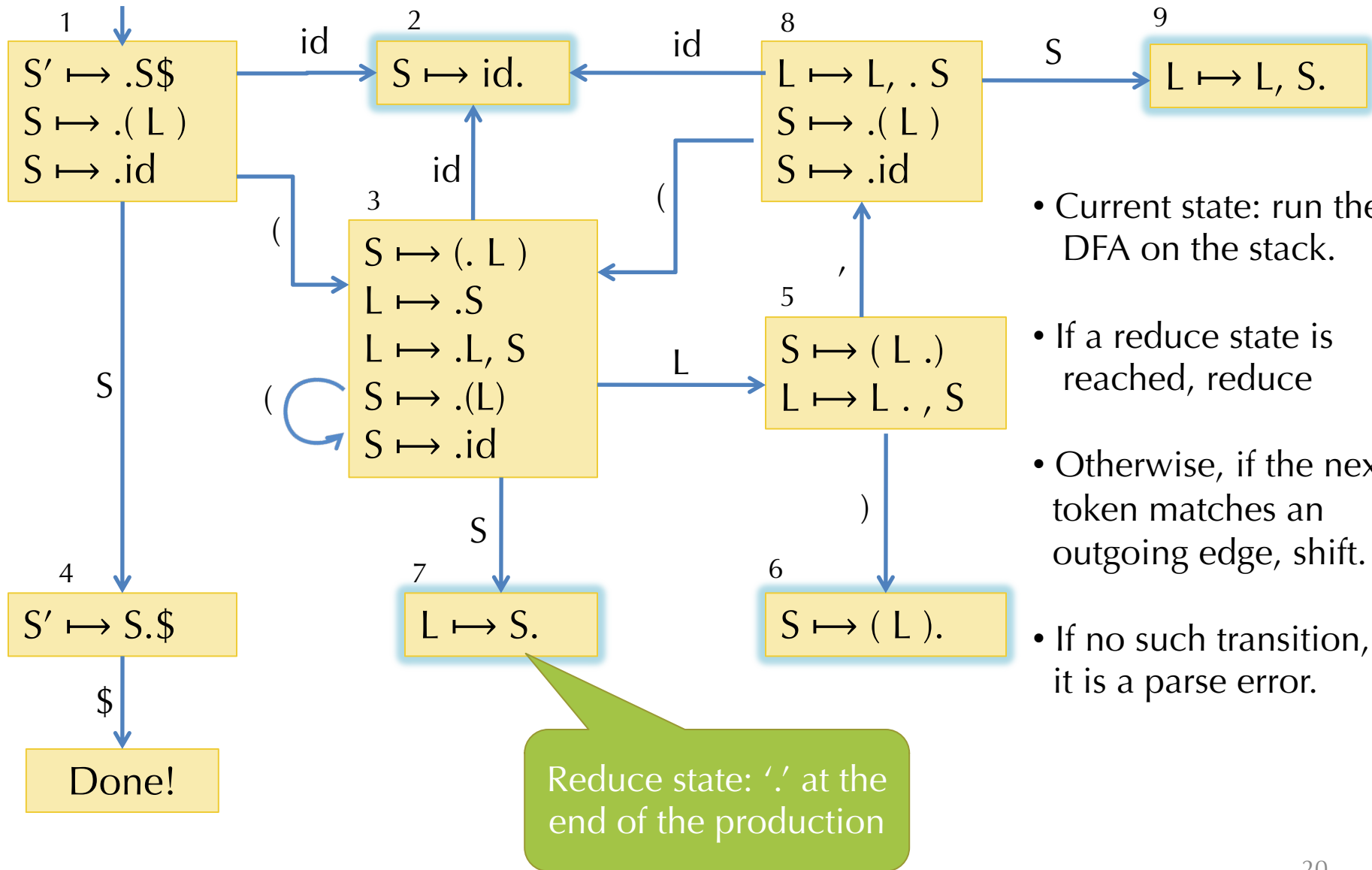
# Example: Constructing the DFA

$S' \mapsto S\$$
$S \mapsto ( L ) \mid id$
$L \mapsto S \mid L , S$

$S' \mapsto .S\$$
$S \mapsto .( L )$
$S \mapsto .id$

id →

$S \mapsto id.$

( →

$S \mapsto (. L )$
$L \mapsto .S$
$L \mapsto .L, S$
$S \mapsto .(L)$
$S \mapsto .id$

S ↓

$S' \mapsto S.\$$

- Finally, for each new state, we take the closure.
- Note that we have to perform two iterations to compute CLOSURE($\{S \mapsto ( . L )\}$)
  - First iteration adds $L \mapsto .S$ and $L \mapsto .L, S$
  - Second iteration adds $S \mapsto .(L)$ and $S \mapsto .id$

19

# Full DFA for the Example

**1**
$S' \mapsto .S\$$
$S \mapsto .( L )$
$S \mapsto .id$

→ id →

**2**
$S \mapsto id.$

← id ←

**8**
$L \mapsto L, . S$
$S \mapsto .( L )$
$S \mapsto .id$

→ S →

**9**
$L \mapsto L, S.$

id ↑

**3**
$S \mapsto (. L )$
$L \mapsto .S$
$L \mapsto .L, S$
$S \mapsto .(L)$
$S \mapsto .id$

(

**5**
$S \mapsto ( L .)$
$L \mapsto L . , S$

→ L →

,

**4**
$S' \mapsto S.\$$

S

$

**7**
$L \mapsto S.$

**6**
$S \mapsto ( L ).$

)

Done!

- Current state: run the DFA on the stack.

- If a reduce state is reached, reduce

- Otherwise, if the next token matches an outgoing edge, shift.

- If no such transition, it is a parse error.

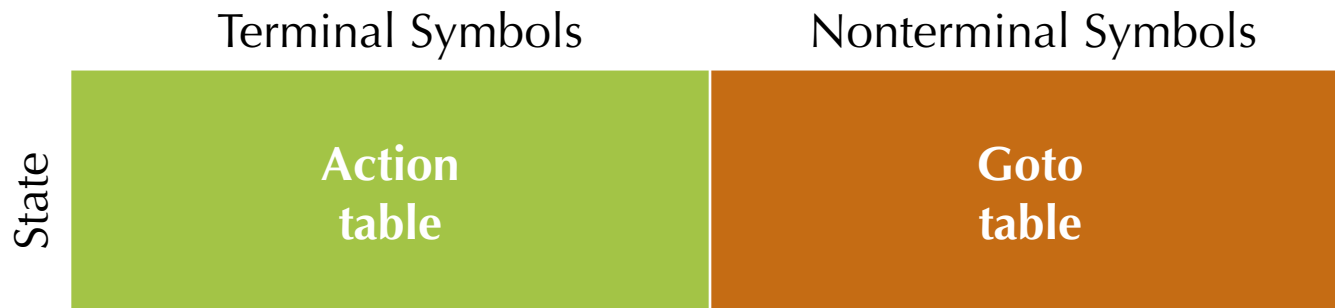Reduce state: '.' at the end of the production

20

# Using the DFA

- Run the parser stack through the DFA.
- The resulting state tells us which productions might be reduced next.
  - If not in a reduce state, then shift the next symbol and transition according to DFA.
  - If in a reduce state, $X \longmapsto \gamma$ with stack $\alpha\gamma$, pop $\gamma$ and push $X$.

- Optimization: No need to re-run the DFA from beginning every step
  - Store the state with each symbol on the stack: e.g. $_1(_3(_3 L_5)_6$
  - On a reduction $X \longmapsto \gamma$, pop stack to reveal the state too:
    e.g.   From stack $_1(_3(_3 L_5)_6$ reduce $S \longmapsto ( L )$ to reach stack $_1(_3$
  - Next, push the reduction symbol: e.g. to reach stack $_1(_3 S$
  - Then take just one step in the DFA to find next state: $_1(_3 S_7$

# Implementing the Parsing Table

Represent the DFA as a table of shape:

state * (terminals + nonterminals)

- Entries for the "action table" specify two kinds of actions:
  - Shift and goto state n
  - Reduce using reduction $X \mapsto \gamma$
    - First pop $\gamma$ off the stack to reveal the state
    - Look up X in the "goto table" and goto that state

|   | Terminal Symbols | Nonterminal Symbols |
|---|---|---|
| State | **Action table** | **Goto table** |

# Example Parse Table

| | ( | ) | id | , | $ | S | L |
|---|---|---|---|---|---|---|---|
| 1 | s3 | | s2 | | | g4 | |
| 2 | S↦id | S↦id | S↦id | S↦id | S↦id | | |
| 3 | s3 | | s2 | | | g7 | g5 |
| 4 | | | | | DONE | | |
| 5 | | s6 | | | s8 | | |
| 6 | S ↦ (L) | S ↦ (L) | S ↦ (L) | S ↦ (L) | S ↦ (L) | | |
| 7 | L ↦ S | L ↦ S | L ↦ S | L ↦ S | L ↦ S | | |
| 8 | s3 | | s2 | | | g9 | |
| 9 | L ↦ L,S | L ↦ L,S | L ↦ L,S | L ↦ L,S | L ↦ L,S | | |

sx = shift and goto state x
gx = goto state x

# Example

- Parse the token stream:  (x, (y, z), w)$

| Stack | Stream | Action (according to table) |
|---|---|---|
| $\varepsilon_1$ | (x, (y, z), w)$ | s3 |
| $\varepsilon_1(_3$ | x, (y, z), w)$ | s2 |
| $\varepsilon_1(_3 x_2$ | , (y, z), w)$ | Reduce: S$\longmapsto$id |
| $\varepsilon_1(_3 S$ | , (y, z), w)$ | g7   (from state 3 follow S) |
| $\varepsilon_1(_3 S_7$ | , (y, z), w)$ | Reduce: L$\longmapsto$S |
| $\varepsilon_1(_3 L$ | , (y, z), w)$ | g5   (from state 3 follow L) |
| $\varepsilon_1(_3 L_5$ | , (y, z), w)$ | s8 |
| $\varepsilon_1(_3 L_{5,8}$ | (y, z), w)$ | s3 |
| $\varepsilon_1(_3 L_{5,8}(_3$ | y, z), w)$ | s2 |

# LR(0) Limitations

- An LR(0) machine only works if states with reduce actions have a *single* reduce action.
  - In such states, the machine *always* reduces (ignoring lookahead)
- With more complex grammars, the DFA construction will yield states with shift/reduce and reduce/reduce conflicts:

| OK | shift/reduce | reduce/reduce |
|----|--------------|---------------|
| $S \mapsto ( L ).$ | $S \mapsto ( L ).$ <br> $L \mapsto .L , S$ | $S \mapsto L ,S.$ <br> $S \mapsto ,S.$ |

- Such conflicts can often be resolved by using a look-ahead symbol: SLR(1) or LR(1)

# Examples

- Consider the left associative and right associative "sum" grammars:

left

$$S \mapsto S + E \mid E$$
$$E \mapsto number \mid ( S )$$

right

$$S \mapsto E + S \mid E$$
$$E \mapsto number \mid ( S )$$

- One is LR(0) the other isn't…  which is which and why?

- What kind of conflict do you get?  Shift/reduce or Reduce/reduce?

- Ambiguities in associativity/precedence usually lead to shift/reduce conflicts.

# SLR(1) ("simple" LR) Parsers

- What conflicts are there in LR(0) parsing?
  - reduce/reduce conflict:  an LR(0) state has two reduce actions
  - shift/reduce conflict: an LR(0) state mixes reduce and shift actions
- Can we use lookahead to disambiguate?

- SLR(1) – uses the same DFA construction as LR(0)
  - modifies the actions based on lookahead

- Suppose reducing an A nonterminal is possible in some state:
  - compute Follow(A) for the given grammar
  - if the lookahead symbol is in Follow(A), then reduce, otherwise shift
  - can disambiguate between reduce/reduce conflicts if the follow sets are disjoint

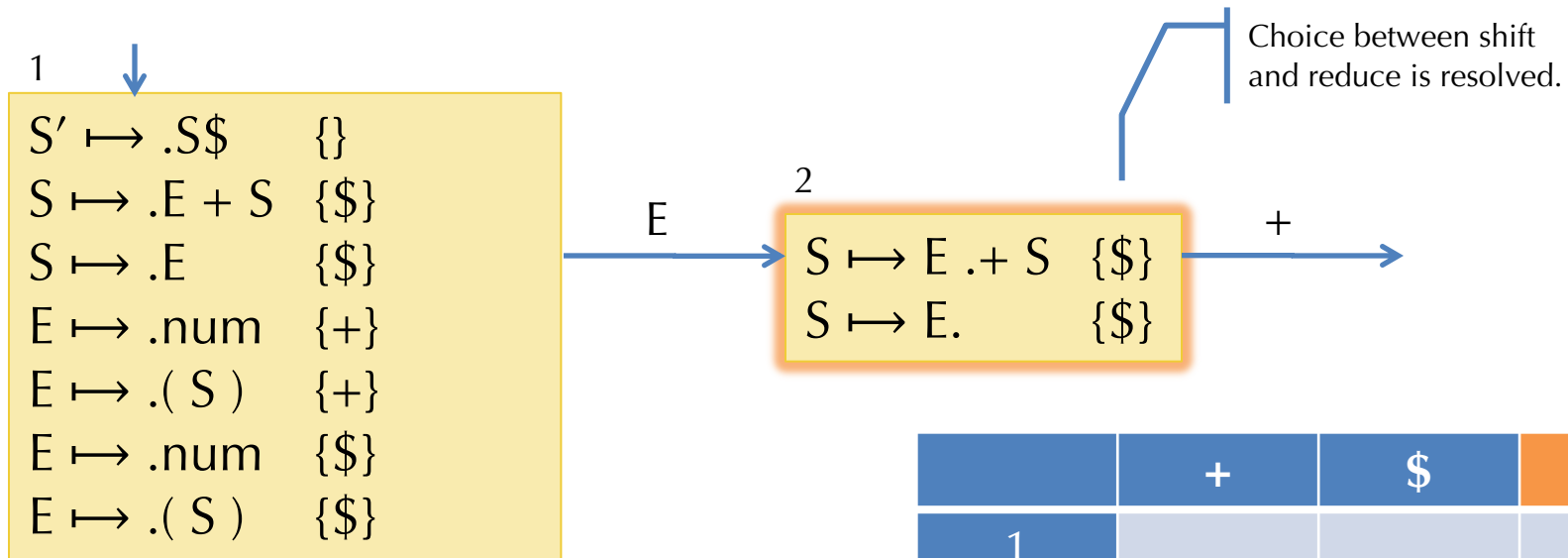# LR(1) Parsing

- Algorithm is similar to LR(0) DFA construction:
    - LR(1) state = set of LR(1) items
    - An LR(1) item is an LR(0) item + a set of look-ahead symbols:
        $$A \mapsto \alpha.\beta \ , \ \mathcal{L}$$

- LR(1) closure is a little more complex:
- Form the set of items just as for LR(0) algorithm.
- Whenever a new item $C \mapsto .\gamma$ is added because $A \mapsto \beta.C\delta \ , \ \mathcal{L}$ is already in the set, we need to compute its look-ahead set $\mathcal{M}$:
    1. The look-ahead set $\mathcal{M}$ includes $FIRST(\delta)$
       (the set of terminals that may start strings derived from $\delta$)
    2. If $\delta$ is itself $\varepsilon$ or can derive $\varepsilon$ (i.e. it is nullable), then the look-ahead $\mathcal{M}$ also contains $\mathcal{L}$

# Example Closure

$$S' \longmapsto S\$$$
$$S \longmapsto E + S \ | \ E$$
$$E \longmapsto number \ | \ ( S )$$

- Start item:     $S' \longmapsto .S\$$  ,  {}

- Since S is to the right of a '.', add:
    $S \longmapsto .E + S$  ,  {\$}                    Note: {\$} is FIRST(\$)
    $S \longmapsto .E$        ,  {\$}

- Need to keep closing, since E appears to the right of a '.' in '.E + S':
    $E \longmapsto .number$ ,  {+}                    Note: + added for reason 1
    $E \longmapsto .( S )$     ,  {+}                    FIRST(+ S) = {+}

- Because E also appears to the right of '.' in '.E' we get:
    $E \longmapsto .number$ ,  {\$}                    Note: \$ added for reason 2
    $E \longmapsto .( S )$     ,  {\$}                    $\delta$ is $\varepsilon$

- All items are distinct, so we're done

# Using the DFA

1

$S' \mapsto .S\$ \qquad \{\}$
$S \mapsto .E + S \quad \{\$\}$
$S \mapsto .E \qquad\quad \{\$\}$
$E \mapsto .num \qquad \{+\}$
$E \mapsto .( S ) \qquad \{+\}$
$E \mapsto .num \qquad \{\$\}$
$E \mapsto .( S ) \qquad \{\$\}$

E →

2

$S \mapsto E .+ S \quad \{\$\}$
$S \mapsto E. \qquad\quad \{\$\}$

+ →

Choice between shift and reduce is resolved.

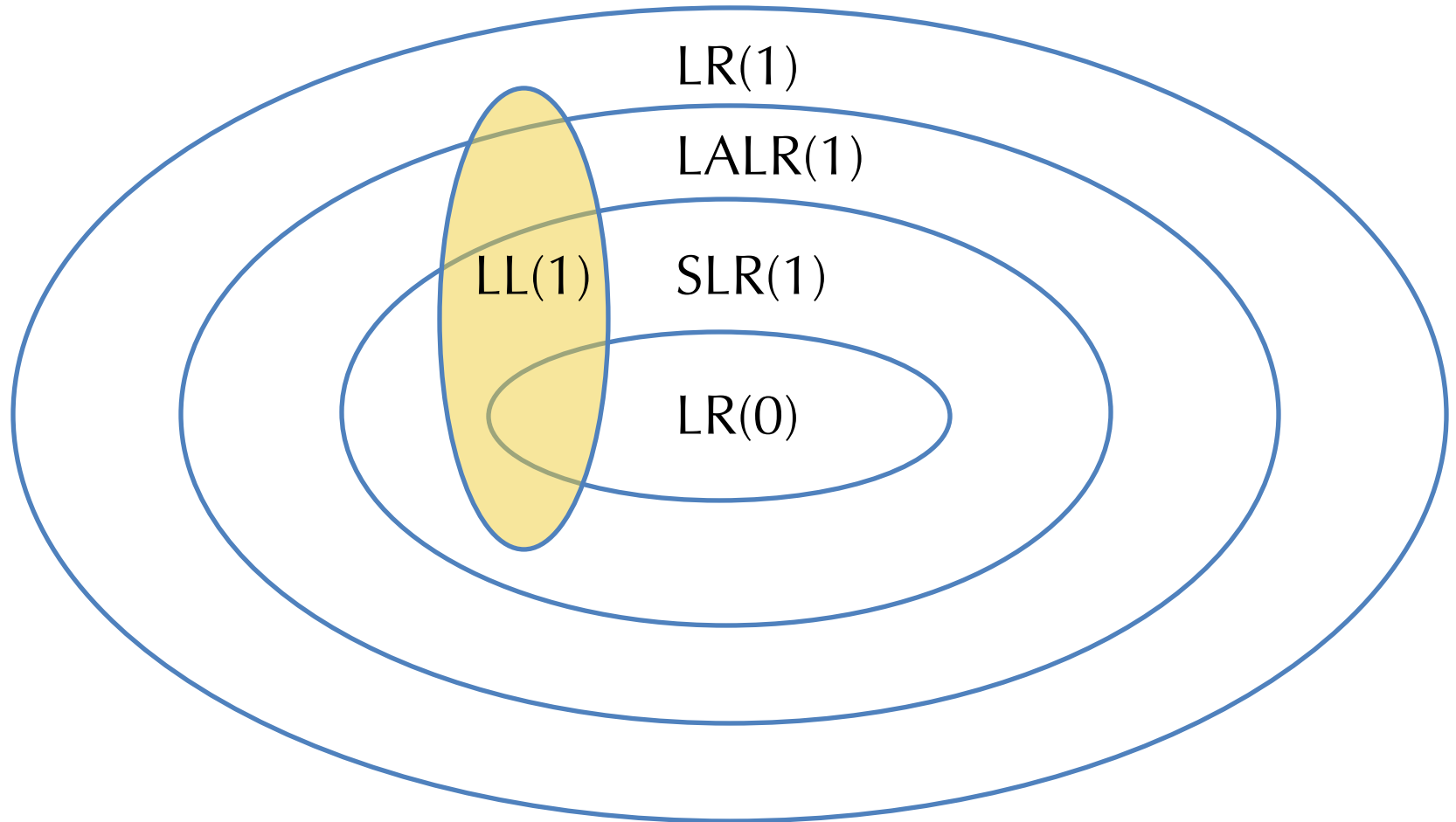| | + | $ | E |
|---|---|---|---|
| 1 | | | g2 |
| 2 | s3 | $S \mapsto E$ | |

Fragment of the Action & Goto tables

- The behavior is determined if:
  - There is no overlap among the look-ahead sets for each reduce item, and
  - None of the look-ahead symbols appear to the right of a '.'

# LR variants

- LR(1) gives maximal power out of a 1 look-ahead symbol parsing table
  - DFA + stack is a push-down automaton
- In practice, LR(1) tables are big.
  - Modern implementations (*e.g.,* menhir) directly generate code

- LALR(1)  = "Look-ahead LR"
  - Merge any two LR(1) states whose items are identical except for the look-ahead sets:

| | |
|---|---|
| S' ⟼ .S\$ | {} |
| S ⟼ .E + S | {\$} |
| S ⟼ .E | {\$} |
| E ⟼ .num | {+} |
| E ⟼ .( S ) | {+} |
| E ⟼ .num | {\$} |
| E ⟼ .( S ) | {\$} |

| | |
|---|---|
| S' ⟼ .S\$ | {} |
| S ⟼ .E + S | {\$} |
| S ⟼ .E | {\$} |
| E ⟼ .num | {+,\$} |
| E ⟼ .( S ) | {+,\$} |

  - Such merging can lead to nondeterminism (*e.g.,* reduce/reduce conflicts), but
  - Results in a much smaller parse table and works well in practice
  - This is the usual technology for automatic parser generators: yacc, ocamlyacc
- GLR = "Generalized LR" parsing
  - Efficiently compute the set of *all* parses for a given input
  - Later passes should disambiguate based on other context

# Classification of Grammars

Debugging parser conflicts.

Disambiguating grammars.

# LALRPOP DEMO