

Lecture 11

# **CIS 341: COMPILERS**

# Announcements

- Reminder: HW3 LLVM backend
  - Due: Monday, Feb 23<sup>rd</sup> at 11:59:59pm
- Midterm Exam: March 5<sup>th</sup> in class!
- HW4: Parsing & basic code generation
  - Available next week
  - Due: After break

Searching for derivations.

# LL & LR PARSING

# CFGs Mathematically

- A Context-free Grammar (CFG) consists of
  - A set of *terminals* (e.g., a token or  $\epsilon$ )
  - A set of *nonterminals* (e.g.,  $S$  and other syntactic variables)
  - A designated nonterminal called the *start symbol*
  - A set of productions:  $\text{LHS} \mapsto \text{RHS}$ 
    - LHS is a nonterminal
    - RHS is a *string* of terminals and nonterminals

- Example: The balanced parentheses language:

$$S \mapsto (S)S$$

$$S \mapsto \epsilon$$

- How many terminals? How many nonterminals? Productions?

# Consider finding left-most derivations

- Look at only one input symbol at a time.

$$\begin{array}{l} S \mapsto E + S \mid E \\ E \mapsto \text{number} \mid ( S ) \end{array}$$

Partly-derived String	Look-ahead	<b>Parsed</b> /Unparsed Input
<u>S</u>	(	(1 + 2 + (3 + 4)) + 5
$\mapsto$ <u>E</u> + S	(	(1 + 2 + (3 + 4)) + 5
$\mapsto$ ( <u>S</u> ) + S	1	(1 + 2 + (3 + 4)) + 5
$\mapsto$ ( <u>E</u> + S) + S	1	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + <u>S</u> ) + S	2	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + <u>E</u> + S) + S	2	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + 2 + <u>S</u> ) + S	(	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + 2 + <u>E</u> ) + S	(	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + 2 + ( <u>S</u> )) + S	3	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + 2 + ( <u>E</u> + S)) + S	3	(1 + 2 + (3 + 4)) + 5
$\mapsto$ ...		

# There is a problem

- We want to decide which production to apply based on the look-ahead symbol.
- But, there is a choice:

$$\begin{array}{l} S \mapsto E + S \mid E \\ E \mapsto \text{number} \mid ( S ) \end{array}$$

(1)  $S \mapsto E \mapsto (S) \mapsto (E) \mapsto (1)$

vs.

(1) + 2  $S \mapsto E + S \mapsto (S) + S \mapsto (E) + S \mapsto (1) + S \mapsto (1) + E$   
 $\mapsto (1) + 2$

- Given the look-ahead symbol: '(' it isn't clear whether to pick  $S \mapsto E$  or  $S \mapsto E + S$  first.

# LL(1) GRAMMARS

# Grammar is the problem

- Not all grammars can be parsed “top-down” with only a single lookahead symbol.
- *Top-down*: starting from the start symbol (root of the parse tree) and going down
- LL(1) means
  - Left-to-right scanning
  - Left-most derivation,
  - 1 lookahead symbol
- This language isn’t “LL(1)”
- Is it LL(k) for some k?
- What can we do?

$$\begin{aligned} S &\mapsto E + S \mid E \\ E &\mapsto \text{number} \mid ( S ) \end{aligned}$$



# Making a grammar LL(1)

- *Problem:* We can't decide which S production to apply until we see the symbol after the first expression.
- *Solution:* "Left-factor" the grammar. There is a common S prefix for each choice, so add a new non-terminal  $S'$  at the decision point:

$$\begin{array}{l} S \mapsto E + S \mid E \\ E \mapsto \text{number} \mid ( S ) \end{array}$$

$$\begin{array}{l} S \mapsto ES' \\ S' \mapsto \epsilon \\ S' \mapsto + S \\ E \mapsto \text{number} \mid ( S ) \end{array}$$

- Also need to eliminate left-recursion somehow. Why?
- Consider:

$$\begin{array}{l} S \mapsto S + E \mid E \\ E \mapsto \text{number} \mid ( S ) \end{array}$$

# LL(1) Parse of the input string

- Look at only one input symbol at a time.

$$\begin{aligned} S &\mapsto ES' \\ S' &\mapsto \varepsilon \\ S' &\mapsto + S \\ E &\mapsto \text{number} \mid ( S ) \end{aligned}$$

Partly-derived String	Look-ahead	Parsed/Unparsed Input
<u>S</u>	(	(1 + 2 + (3 + 4)) + 5
$\mapsto$ <u>E</u> S'	(	(1 + 2 + (3 + 4)) + 5
$\mapsto$ ( <u>S</u> ) S'	1	(1 + 2 + (3 + 4)) + 5
$\mapsto$ ( <u>E</u> S') S'	1	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 <u>S'</u> ) S'	+	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + <u>S</u> ) S'	2	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + <u>E</u> S') S'	2	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + 2 <u>S'</u> ) S'	+	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + 2 + <u>S</u> ) S'	(	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + 2 + <u>E</u> S') S'	(	(1 + 2 + (3 + 4)) + 5
$\mapsto$ (1 + 2 + ( <u>S</u> )S') S'	3	(1 + 2 + (3 + 4)) + 5

# Predictive Parsing

- Given an LL(1) grammar:
  - For a given nonterminal, the lookahead symbol uniquely determines the production to apply.
  - Top-down parsing = predictive parsing
  - Driven by a predictive parsing table:  
nonterminal \* input token  $\rightarrow$  production

$T \mapsto S\$$   
 $S \mapsto ES'$   
 $S' \mapsto \epsilon$   
 $S' \mapsto + S$   
 $E \mapsto \text{number} \mid ( S )$

	number	+	(	)	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \epsilon$	$\mapsto \epsilon$
E	$\mapsto \text{num.}$		$\mapsto ( S )$		

- Note: it is convenient to add a special *end-of-file* token \$ and a start symbol T (top-level) that requires \$.

# How do we construct the parse table?

- Consider a given production:  $A \rightarrow \gamma$
- Construct the set of all input tokens that may appear *first* in strings that can be derived from  $\gamma$ 
  - Add the production  $\rightarrow \gamma$  to the entry (A,token) for each such token.
- If  $\gamma$  can derive  $\epsilon$  (the empty string), then we construct the set of all input tokens that may *follow* the nonterminal A in the grammar.
  - Add the production  $\rightarrow \gamma$  to the entry (A, token) for each such token.
- Note: if there are two different productions for a given entry, the grammar is not LL(1)

# Example

- $\text{First}(T) = \text{First}(S)$
- $\text{First}(S) = \text{First}(E)$
- $\text{First}(S') = \{ + \}$
- $\text{First}(E) = \{ \text{number}, '(' \}$

$T \mapsto S\$$   
 $S \mapsto ES'$   
 $S' \mapsto \epsilon$   
 $S' \mapsto + S$   
 $E \mapsto \text{number} \mid ( S )$

- $\text{Follow}(S') = \text{Follow}(S)$
- $\text{Follow}(S) = \{ \$, ')' \} \cup \text{Follow}(S')$

**Note:** we want the *least* solution to this system of set equations... a *fixpoint* computation. More on these later in the course.

	number	+	(	)	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \epsilon$	$\mapsto \epsilon$
E	$\mapsto \text{num.}$		$\mapsto ( S )$		

# Converting the table to code

- Define  $n$  mutually recursive functions
  - one for each nonterminal  $A$ : `parse_A`
  - The type of `parse_A` is `unit -> ast` if  $A$  is *not* an auxiliary nonterminal
  - Parse functions for auxiliary nonterminals (e.g.  $S'$ ) take extra `ast`'s as inputs, one for each nonterminal in the “factored” prefix.
- Each function “peeks” at the lookahead token and then follows the production rule in the corresponding entry.
  - Consume terminal tokens from the input stream
  - Call `parse_X` to create sub-tree for nonterminal  $X$
  - If the rule ends in an auxiliary nonterminal, call it with appropriate `ast`'s. (The auxiliary rule is responsible for creating the `ast` after looking at more input.)
  - Otherwise, this function builds the `ast` tree itself and returns it.

	number	+	(	)	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \epsilon$	$\mapsto \epsilon$
E	$\mapsto \text{num.}$		$\mapsto ( S )$		

Hand-generated LL(1) code for the table above.

## DEMO: PARSER.ML

# LL(1) Summary

- Top-down parsing that finds the leftmost derivation.
- Language Grammar  $\Rightarrow$  LL(1) grammar  $\Rightarrow$  prediction table  $\Rightarrow$  recursive-descent parser
- Problems:
  - Grammar must be LL(1)
  - Can extend to LL(k) (it just makes the table bigger)
  - Grammar cannot be left recursive (parser functions will loop!)
- Is there a better way?



# LR GRAMMARS

# Bottom-up Parsing (LR Parsers)

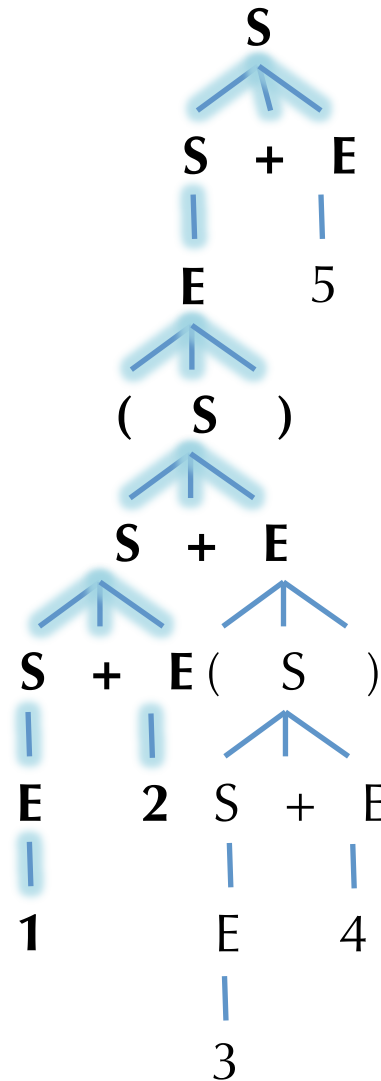
- LR(k) parser:
  - Left-to-right scanning
  - Rightmost derivation
  - k lookahead symbols
- LR grammars are more expressive than LL
  - Can handle left-recursive (and right recursive) grammars; virtually all programming languages
  - Easier to express programming language syntax (no left factoring)
- 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, CUP, ocamllyacc, merlin, etc.)
  - Better error detection/recovery

# Top-down vs. Bottom up

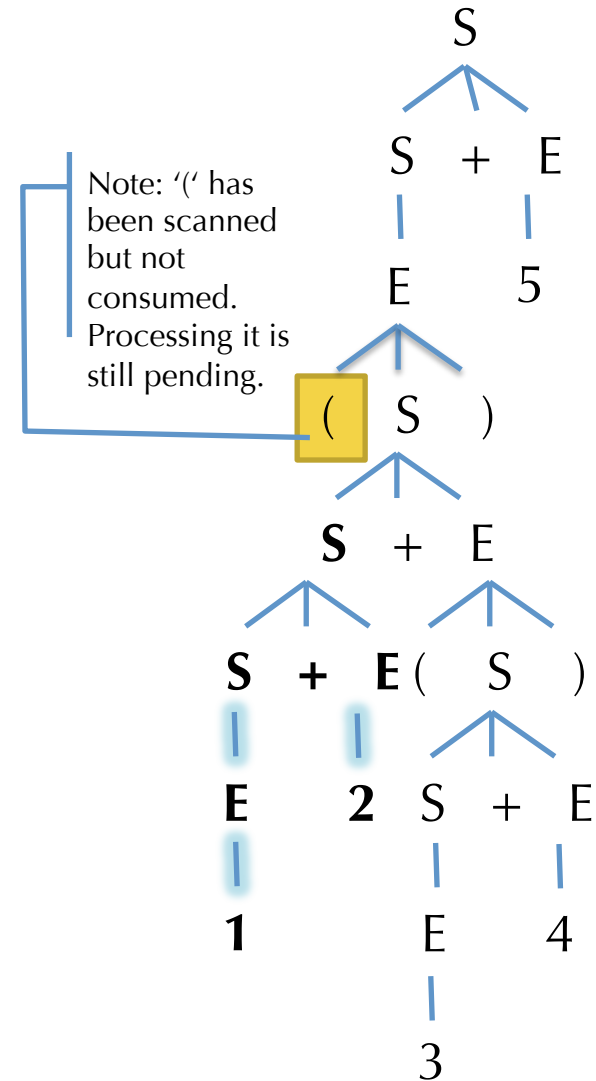
- Consider the left-recursive grammar:

$S \mapsto S + E \mid E$   
 $E \mapsto \text{number} \mid ( S )$

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



Top-down



Bottom-up

# Progress of Bottom-up Parsing

	Reductions	Scanned	Input Remaining
	$(1 + 2 + (3 + 4)) + 5 \leftarrow$		$(1 + 2 + (3 + 4)) + 5$
	$(\underline{E} + 2 + (3 + 4)) + 5 \leftarrow$	$($	$+ 2 + (3 + 4)) + 5$
	$(\underline{S} + 2 + (3 + 4)) + 5 \leftarrow$	$(1$	$+ 2 + (3 + 4)) + 5$
	$(S + \underline{E} + (3 + 4)) + 5 \leftarrow$	$(1 + 2$	$+ (3 + 4)) + 5$
	$(\underline{S} + (3 + 4)) + 5 \leftarrow$	$(1 + 2$	$+ (3 + 4)) + 5$
	$(S + (\underline{E} + 4)) + 5 \leftarrow$	$(1 + 2 + (3$	$+ 4)) + 5$
	$(S + (\underline{S} + 4)) + 5 \leftarrow$	$(1 + 2 + (3$	$+ 4)) + 5$
	$(S + (S + \underline{E})) + 5 \leftarrow$	$(1 + 2 + (3 + 4$	$)) + 5$
	$(S + (\underline{S})) + 5 \leftarrow$	$(1 + 2 + (3 + 4$	$)) + 5$
	$(S + \underline{E}) + 5 \leftarrow$	$(1 + 2 + (3 + 4)$	$) + 5$
	$(\underline{S}) + 5 \leftarrow$	$(1 + 2 + (3 + 4)$	$) + 5$
	$\underline{E} + 5 \leftarrow$	$(1 + 2 + (3 + 4))$	$+ 5$
	$\underline{S} + 5 \leftarrow$	$(1 + 2 + (3 + 4))$	$+ 5$
	$S + \underline{E} \leftarrow$	$(1 + 2 + (3 + 4)) + 5$	
	$S$		

Rightmost derivation

$S \mapsto S + E \mid E$   
 $E \mapsto \text{number} \mid ( S )$

# Shift/Reduce Parsing

- Parser state:
  - Stack of terminals and nonterminals.
  - Unconsumed input is a string of terminals
  - Current derivation step is  $\text{stack} + \text{input}$
- 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  $X \mapsto \gamma$  is a production. (pop  $\gamma$ , push  $X$ )

$$S \mapsto S + E \mid E$$

$$E \mapsto \text{number} \mid ( S )$$

Stack	Input	Action
	(1 + 2 + (3 + 4)) + 5	shift (
(	1 + 2 + (3 + 4)) + 5	shift 1
(1	+ 2 + (3 + 4)) + 5	reduce: $E \mapsto \text{number}$
(E	+ 2 + (3 + 4)) + 5	reduce: $S \mapsto E$
(S	+ 2 + (3 + 4)) + 5	shift +
(S +	2 + (3 + 4)) + 5	shift 2
(S + 2	+ (3 + 4)) + 5	reduce: $E \mapsto \text{number}$

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  $\alpha$  as a finite parser state.
  - Parser state is computed by a DFA that reads the stack  $\sigma$ .
  - Accept states of the DFA correspond to unique reductions that apply.
- Example: LR(0) parsing
  - Left-to-right scanning, Right-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.

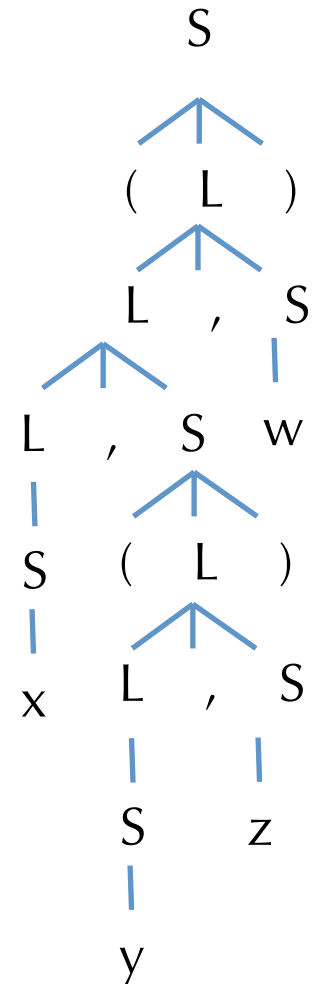
# Example LR(0) Grammar: Tuples

- Example grammar for non-empty tuples and identifiers:

$$\begin{array}{lcl} S & \mapsto & ( L ) \mid \text{id} \\ L & \mapsto & S \mid L, S \end{array}$$

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

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





# Shift/Reduce Parsing

- 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.

$$\begin{array}{l} S \mapsto ( L ) \mid id \\ L \mapsto S \mid L , S \end{array}$$

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

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

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

# Example Run

Stack	Input	Action
	(x, (y, z), w)	shift (
(	x, (y, z), w)	shift x
(x	, (y, z), w)	reduce $S \mapsto \text{id}$
(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 $S \mapsto \text{id}$
(L, (S	, z), w)	reduce $L \mapsto S$
(L, (L	, z), w)	shift ,
(L, (L,	z), w)	shift z
(L, (L, z	), w)	reduce $S \mapsto \text{id}$
(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 ,

$S \mapsto ( L ) \mid \text{id}$   
 $L \mapsto S \mid L, S$

# Action Selection Problem

- Given a stack  $\sigma$  and a look-ahead symbol  $b$ , should the parser:
  - Shift  $b$  onto the stack (new stack is  $\sigma b$ )
  - Reduce a production  $X \mapsto \gamma$ , assuming that  $\sigma = \alpha\gamma$  (new stack is  $\alpha X$ )?
- Sometimes the parser can reduce but shouldn't
  - For example,  $X \mapsto \epsilon$  can *always* be reduced
- Sometimes the stack can be reduced in different ways
- Main idea: decide what to do based on a *prefix*  $\alpha$  of the stack plus the look-ahead symbol.
  - The prefix  $\alpha$  is different for different possible reductions since in productions  $X \mapsto \gamma$  and  $Y \mapsto \beta$ ,  $\gamma$  and  $\beta$  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

$$\begin{array}{l} S \mapsto ( L ) \mid id \\ L \mapsto S \mid L , S \end{array}$$

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

# Constructing the DFA: Start state & Closure

- First step: Add a new production  $S' \mapsto S\$$  to the grammar
- 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:  $\text{CLOSURE}(\{S' \mapsto .S\$\}) = \{S' \mapsto .S\$, S \mapsto .(L), S \mapsto .id\}$
- Resulting “closed state” contains the set of all possible productions that might be reduced next.

$$\begin{array}{l} S' \mapsto S\$ \\ S \mapsto ( L ) \mid id \\ L \mapsto S \mid L , S \end{array}$$

# Example: Constructing the DFA

$S' \mapsto .S\$$

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

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

# Example: Constructing the DFA

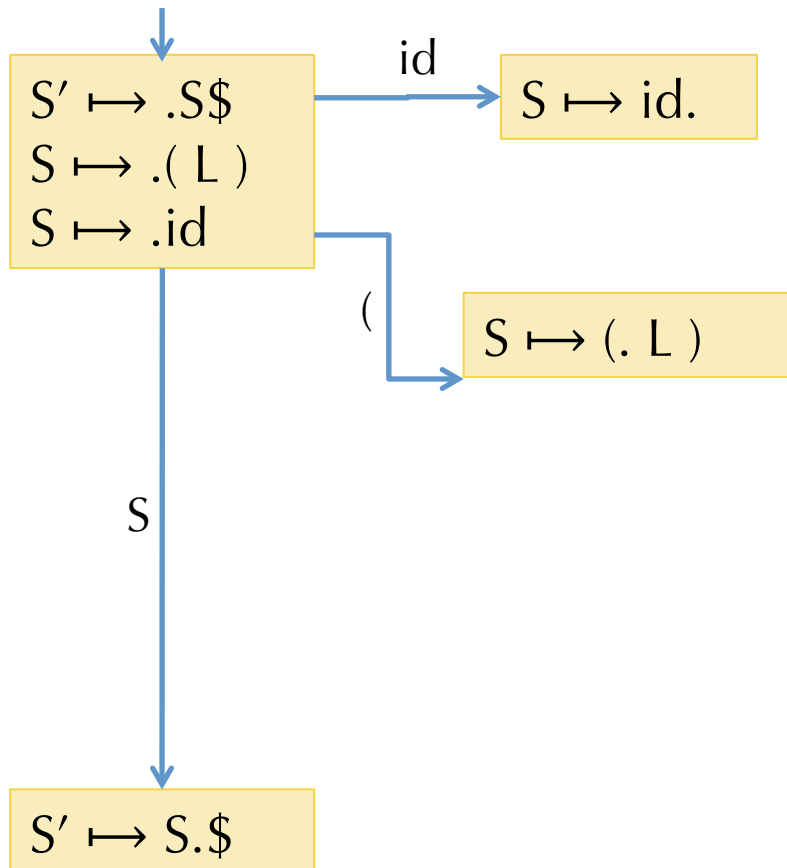
↓  
 $S' \mapsto .S\$$   
 $S \mapsto .( L )$   
 $S \mapsto .id$

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

- 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

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

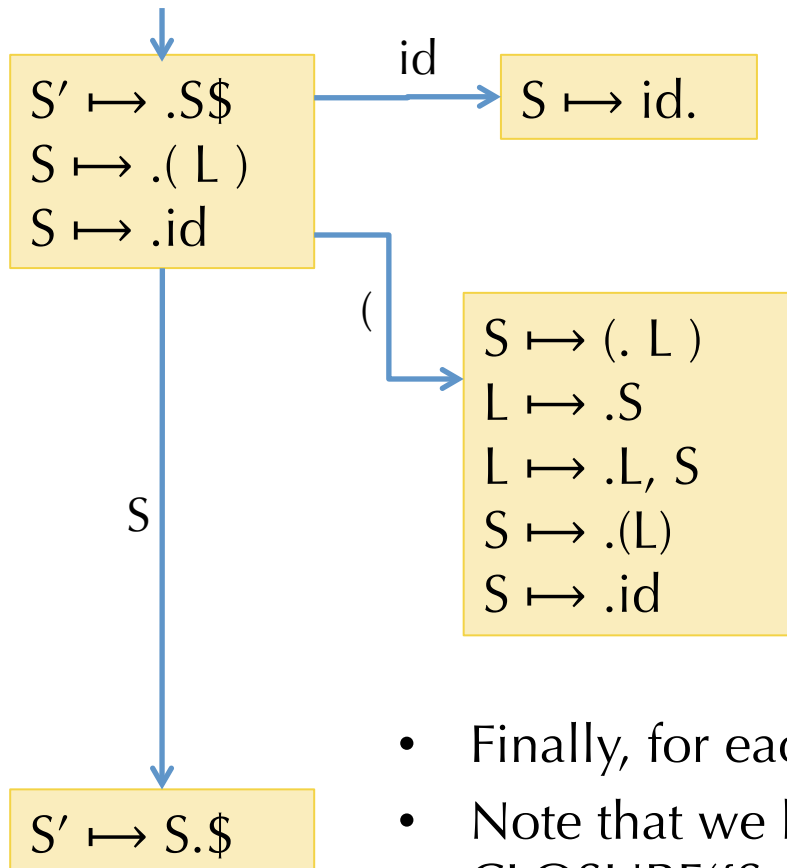


- Next we add the transitions:
- First, we see what terminals and nonterminals can appear after the  $'.'$  in the source state.
  - Outgoing edges have those label.
- The target state (initially) includes all items from the source state that have the edge-label symbol after the  $'.'$ , but we advance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



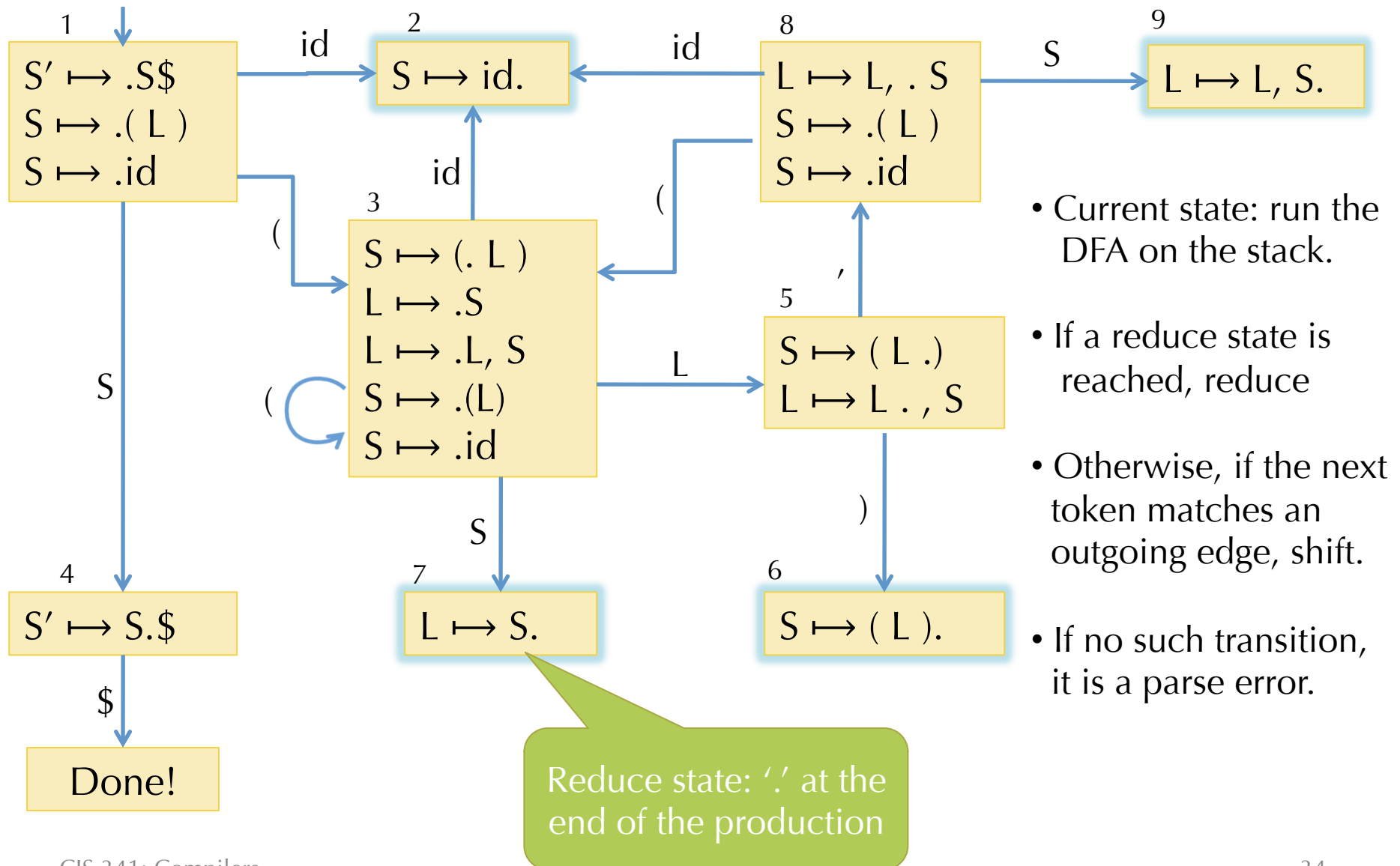
# Example: Constructing the DFA

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



- 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  $S \mapsto .id$

# Full DFA for the Example



- 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.

## Using the DFA

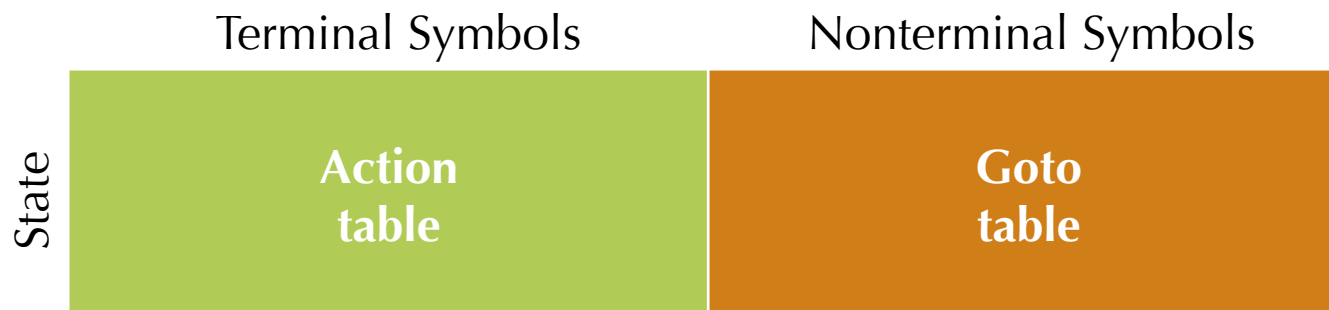
- Run the parser stack through the DFA.
- The resulting state tells us which productions might be reduced next.
  - If not in a reduce state, then shift the next symbol and transition according to DFA.
  - If in a reduce state,  $X \mapsto \gamma$  with stack  $\alpha\gamma$ , pop  $\gamma$  and push  $X$ .
- Optimization: No need to re-run the DFA from beginning every step
  - Store the state with each symbol on the stack: e.g.  $_1(_3(_3L_5)_6$
  - On a reduction  $X \mapsto \gamma$ , pop stack to reveal the state too:  
e.g. From stack  $_1(_3(_3L_5)_6$  reduce  $S \mapsto ( L )$  to reach stack  $_1(_3$
  - Next, push the reduction symbol: e.g. to reach stack  $_1(_3S$
  - Then take just one step in the DFA to find next state:  $_1(_3S_7$

# 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



# Example Parse Table

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

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)
$\epsilon_1$	$(x, (y, z), w)\$$	s3
$\epsilon_1($	$x, (y, z), w)\$$	s2
$\epsilon_1($	$, (y, z), w)\$$	Reduce: $S \mapsto id$
$\epsilon_1($	$, (y, z), w)\$$	g7 (from state 3 follow S)
$\epsilon_1($	$, (y, z), w)\$$	Reduce: $L \mapsto S$
$\epsilon_1($	$, (y, z), w)\$$	g5 (from state 3 follow L)
$\epsilon_1($	$, (y, z), w)\$$	s8
$\epsilon_1($	$(y, z), w)\$$	s3
$\epsilon_1($	$y, z), w)\$$	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

$S \mapsto ( L ).$

shift/reduce

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

reduce/reduce

$S \mapsto L , S.$   
 $S \mapsto , S.$

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

# Examples

- Consider the left associative and right associative “sum” grammars:

left

$$\begin{array}{l} S \mapsto S + E \mid E \\ E \mapsto \text{number} \mid ( S ) \end{array}$$

right

$$\begin{array}{l} S \mapsto E + S \mid E \\ E \mapsto \text{number} \mid ( S ) \end{array}$$

- 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.



# LR(1) Parsing

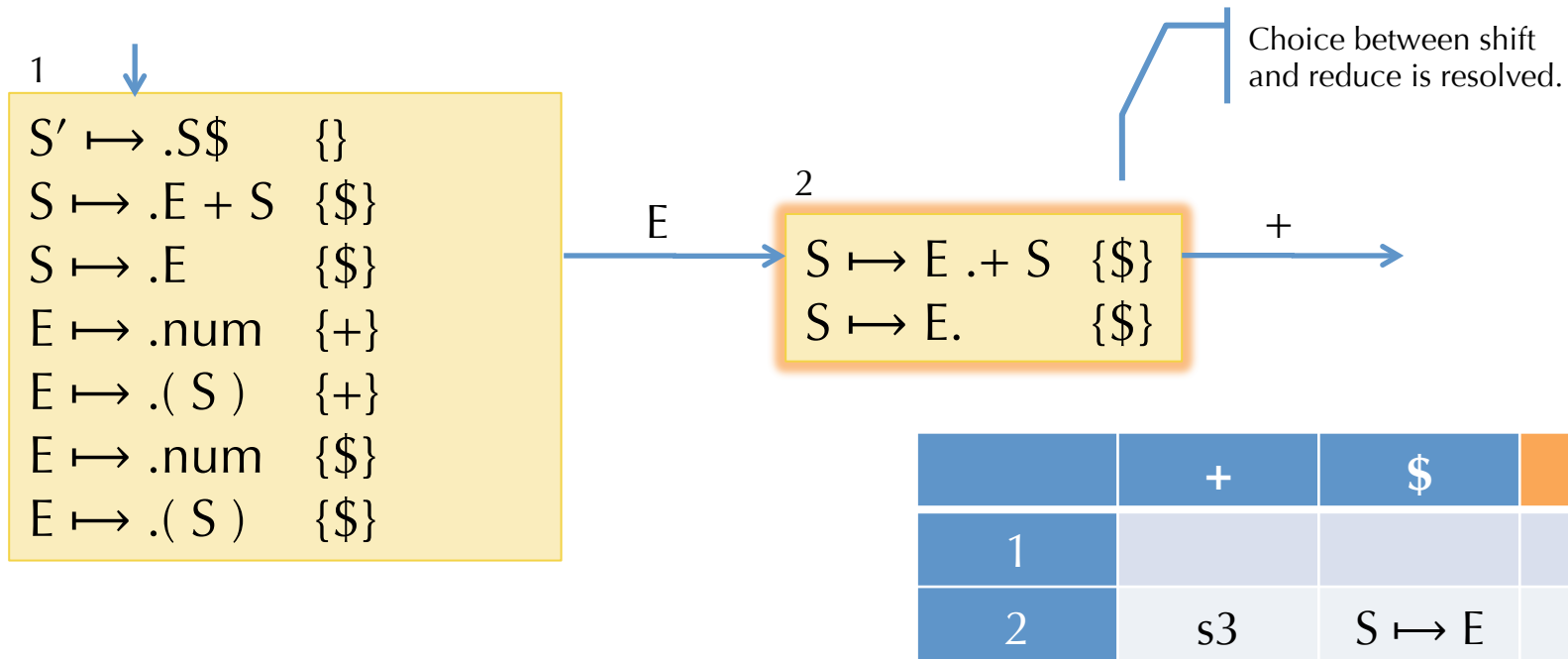
- Algorithm is similar to LR(0) DFA construction:
  - LR(1) state = set of LR(1) items
  - An LR(1) item is an LR(0) item + a set of look-ahead symbols:  
 $A \mapsto \alpha.\beta, \mathcal{L}$
- LR(1) closure is a little more complex:
- Form the set of items just as for LR(0) algorithm.
- Whenever a 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  $\text{FIRST}(\delta)$   
(the set of terminals that may start strings derived from  $\delta$ )
  2. If  $\delta$  can derive  $\epsilon$  (it is nullable), then the look-ahead  $\mathcal{M}$  also contains  $\mathcal{L}$

# Example Closure

$$\begin{aligned} S' &\mapsto S\$ \\ S &\mapsto E + S \mid E \\ E &\mapsto \text{number} \mid ( S ) \end{aligned}$$

- Start item:  $S' \mapsto .S\$$  ,  $\{\}$
- Since S is to the right of a '.', add:  
 $S \mapsto .E + S$  ,  $\{\$\}$       Note:  $\{\$\}$  is FIRST( $\$$ )  
 $S \mapsto .E$  ,  $\{\$\}$
- Need to keep closing, since E appears to the right of a '.' in ' $.E + S$ ':  
 $E \mapsto .\text{number}$  ,  $\{+\}$       Note: + added for reason 1  
 $E \mapsto .( S )$  ,  $\{+\}$
- Because E also appears to the right of '.' in ' $.E$ ' we get:  
 $E \mapsto .\text{number}$  ,  $\{\$\}$       Note: \$ added for reason 2  
 $E \mapsto .( S )$  ,  $\{\$\}$
- All items are distinct, so we're done

# Using the DFA



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 (recall 262)
- 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:

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

$S' \mapsto .S\$$	$\{\}$
$S \mapsto .E + S$	$\{\$ \}$
$S \mapsto .E$	$\{\$ \}$
$E \mapsto .num$	$\{+, \$ \}$
$E \mapsto .( S )$	$\{+, \$ \}$
  - 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, ocaml yacc
- 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

