

November 27

# **EECS 483: COMPILER CONSTRUCTION**



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: LHS  $\mapsto$  RHS
    - LHS is a nonterminal
    - RHS is a *string* of terminals and nonterminals
- Example: The balanced parentheses language:

$$S \mapsto (S)S$$

$$S \mapsto \epsilon$$

# Derivations in CFGs

- Example: derive  $(1 + 2 + (3 + 4)) + 5$

- $\underline{S} \mapsto \underline{E} + S$

$$\mapsto (\underline{S}) + S$$

$$\mapsto (\underline{E} + S) + S$$

$$\mapsto (1 + \underline{S}) + S$$

$$\mapsto (1 + \underline{E} + S) + S$$

$$\mapsto (1 + 2 + \underline{S}) + S$$

$$\mapsto (1 + 2 + \underline{E}) + S$$

$$\mapsto (1 + 2 + (\underline{S})) + S$$

$$\mapsto (1 + 2 + (\underline{E} + S)) + S$$

$$\mapsto (1 + 2 + (3 + \underline{S})) + S$$

$$\mapsto (1 + 2 + (3 + \underline{E})) + S$$

$$\mapsto (1 + 2 + (3 + 4)) + \underline{S}$$

$$\mapsto (1 + 2 + (3 + 4)) + \underline{E}$$

$$\mapsto (1 + 2 + (3 + 4)) + 5$$

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

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

For arbitrary strings  $\alpha, \beta, \gamma$  and production rule  $A \mapsto \beta$  a single step of the derivation is:

$$\square \quad \square \quad \alpha A \gamma \mapsto \alpha \beta \gamma$$

( *substitute*  $\beta$  for an occurrence of  $A$  )

In general, there are many possible derivations for a given string

Note: Underline indicates symbol being expanded.

# Example: Left- and rightmost derivations

- Leftmost derivation:

- $\underline{S} \mapsto \underline{E} + S$
- $\mapsto (\underline{S}) + S$
- $\mapsto (\underline{E} + S) + S$
- $\mapsto (1 + \underline{S}) + S$
- $\mapsto (1 + \underline{E} + S) + S$
- $\mapsto (1 + 2 + \underline{S}) + S$
- $\mapsto (1 + 2 + \underline{E}) + S$
- $\mapsto (1 + 2 + (\underline{S})) + S$
- $\mapsto (1 + 2 + (\underline{E} + S)) + S$
- $\mapsto (1 + 2 + (3 + \underline{S})) + S$
- $\mapsto (1 + 2 + (3 + \underline{E})) + S$
- $\mapsto (1 + 2 + (3 + 4)) + \underline{S}$
- $\mapsto (1 + 2 + (3 + 4)) + \underline{E}$
- $\mapsto (1 + 2 + (3 + 4)) + 5$

- Rightmost derivation:

- $\underline{S} \mapsto E + \underline{S}$
- $\mapsto E + \underline{E}$
- $\mapsto \underline{E} + 5$
- $\mapsto (\underline{S}) + 5$
- $\mapsto (E + \underline{S}) + 5$
- $\mapsto (E + E + \underline{S}) + 5$
- $\mapsto (E + E + \underline{E}) + 5$
- $\mapsto (E + E + (\underline{S})) + 5$
- $\mapsto (E + E + (E + \underline{S})) + 5$
- $\mapsto (E + E + (E + \underline{E})) + 5$
- $\mapsto (E + E + (\underline{E} + 4)) + 5$
- $\mapsto (E + \underline{E} + (3 + 4)) + 5$
- $\mapsto (\underline{E} + 2 + (3 + 4)) + 5$
- $\mapsto (1 + 2 + (3 + 4)) + 5$

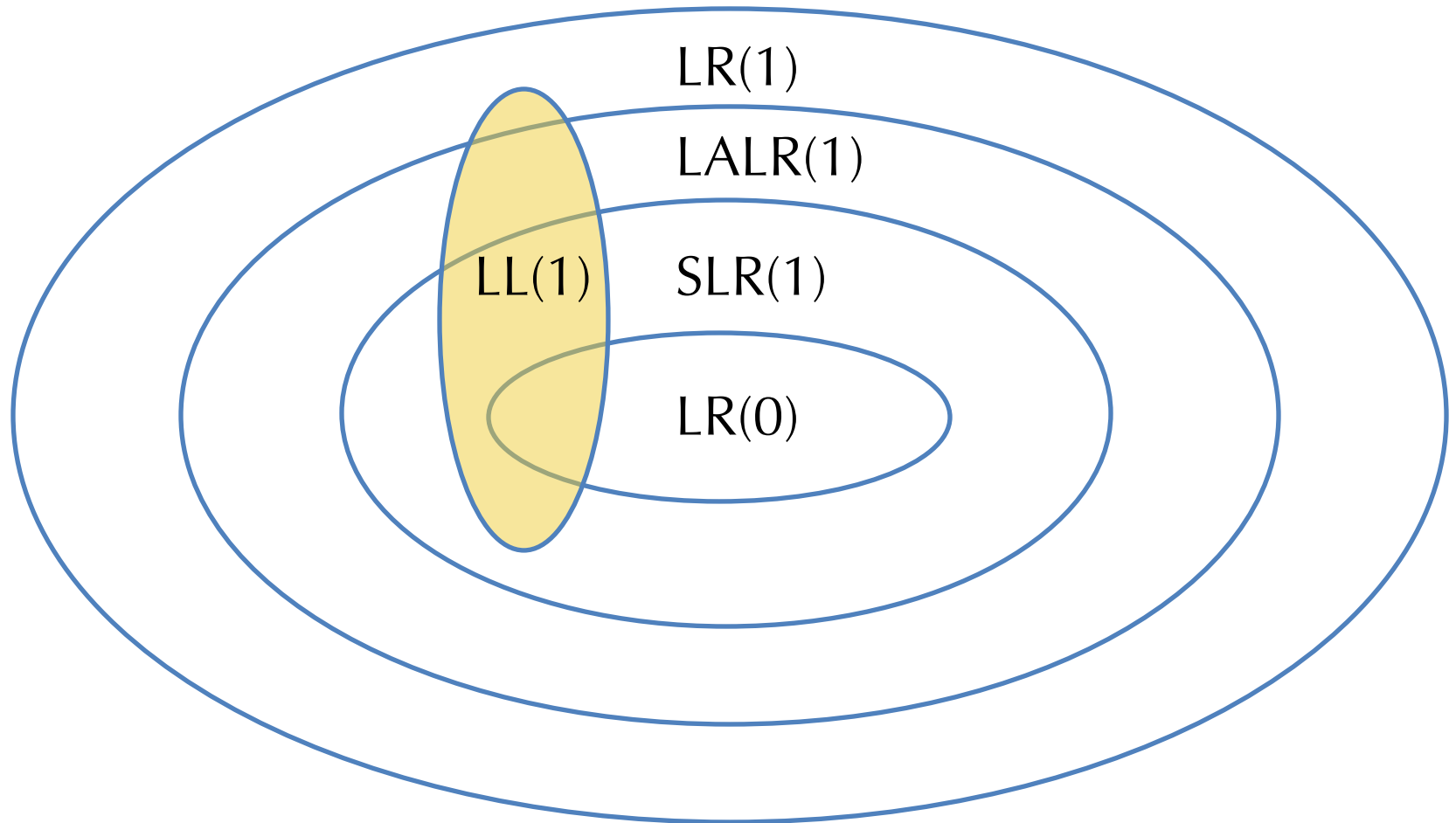
$$S \mapsto E + S \mid E$$

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

# CFGs In Practice

- Context-free Grammars are elegant, *declarative* specifications, generalizing regular expressions
- A parser for a CFG amounts to a *search procedure* for derivations
- Unlike regular expressions, which are easily compiled to linear time recognizers, practical algorithms for parsing *general* CFGs are  $O(n^3)$  in input string length
  - Compromise: add restrictions to the CFGs
    - Benefit: Linear time
    - Drawback: have to rewrite the grammar to make it fit the restrictions

# Classification of Grammars



# LL(1) GRAMMARS



# Consider finding left-most derivations

- Look at only one input symbol at a time.

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

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	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.

# 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:

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



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

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

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

# 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  $(A, \text{token})$  for each such token.
- If  $\gamma$  can derive  $\varepsilon$  (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  $(A, \text{token})$  for each such token.
  
- Note: The grammar is LL(1) *if and only if* all entries have at most one production

# 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. Just like in program analysis!

	number	+	(	)	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto ES'$		$\mapsto ES'$		
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  $() \rightarrow \text{ast}$  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: HANDPARSER.RS

# 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?



Next time  
**LR GRAMMARS**