

Lecture 11

# COMPILER DESIGN

# Announcements

- **HW 3:** Compiling LLVMlite
  - Familiarize yourself with (a subset of) the LLVM IR
  - Implement a translation down to (inefficient) X86lite



# RECAP

# Parsing

Source Code

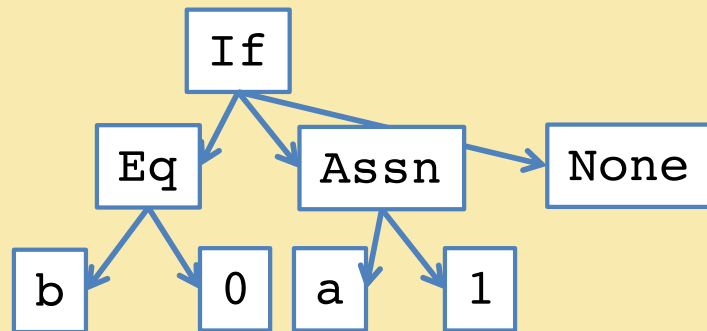
(Character stream)

```
if (b == 0) { a = 1; }
```

Token stream:

if	(	b	==	0	)	{	a	=	1	;	}
----	---	---	----	---	---	---	---	---	---	---	---

Abstract Syntax Tree:



Intermediate code:

```
11: %cnd = icmp eq i64 %b, 0
    br i1 %cnd, label %12, label %13
12: store i64* %a, 1
    br label %13
13:
```

Assembly Code

```
11: cmpq %eax, $0
    jeq 12
    jmp 13
12:
...
```

Lexical Analysis

Parsing

Analysis & Transformation

Backend

# Context-free Grammars (CFGs)

- A Context-free Grammar (CFG) consists of
  - A set of *terminals* (e.g., a lexical token, but [how about  \$\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  $(1 + 2 + (3 + 4)) + 5$
- $\underline{S} \mapsto \underline{E} + S$ 
  - $\mapsto (\underline{S}) + S$
  - $\mapsto (\underline{E} + S) + S$
  - $\mapsto (1 + \underline{S}) + S$
  - $\mapsto (1 + \underline{E} + S) + S$
  - $\mapsto (1 + 2 + \underline{S}) + S$
  - $\mapsto (1 + 2 + \underline{E}) + S$
  - $\mapsto (1 + 2 + (\underline{S})) + S$
  - $\mapsto (1 + 2 + (\underline{E} + S)) + S$
  - $\mapsto (1 + 2 + (3 + \underline{S})) + S$
  - $\mapsto (1 + 2 + (3 + \underline{E})) + S$
  - $\mapsto (1 + 2 + (3 + 4)) + \underline{S}$
  - $\mapsto (1 + 2 + (3 + 4)) + \underline{E}$
  - $\mapsto (1 + 2 + (3 + 4)) + 5$

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

# 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

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

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

Searching for derivations

# LL & LR PARSING

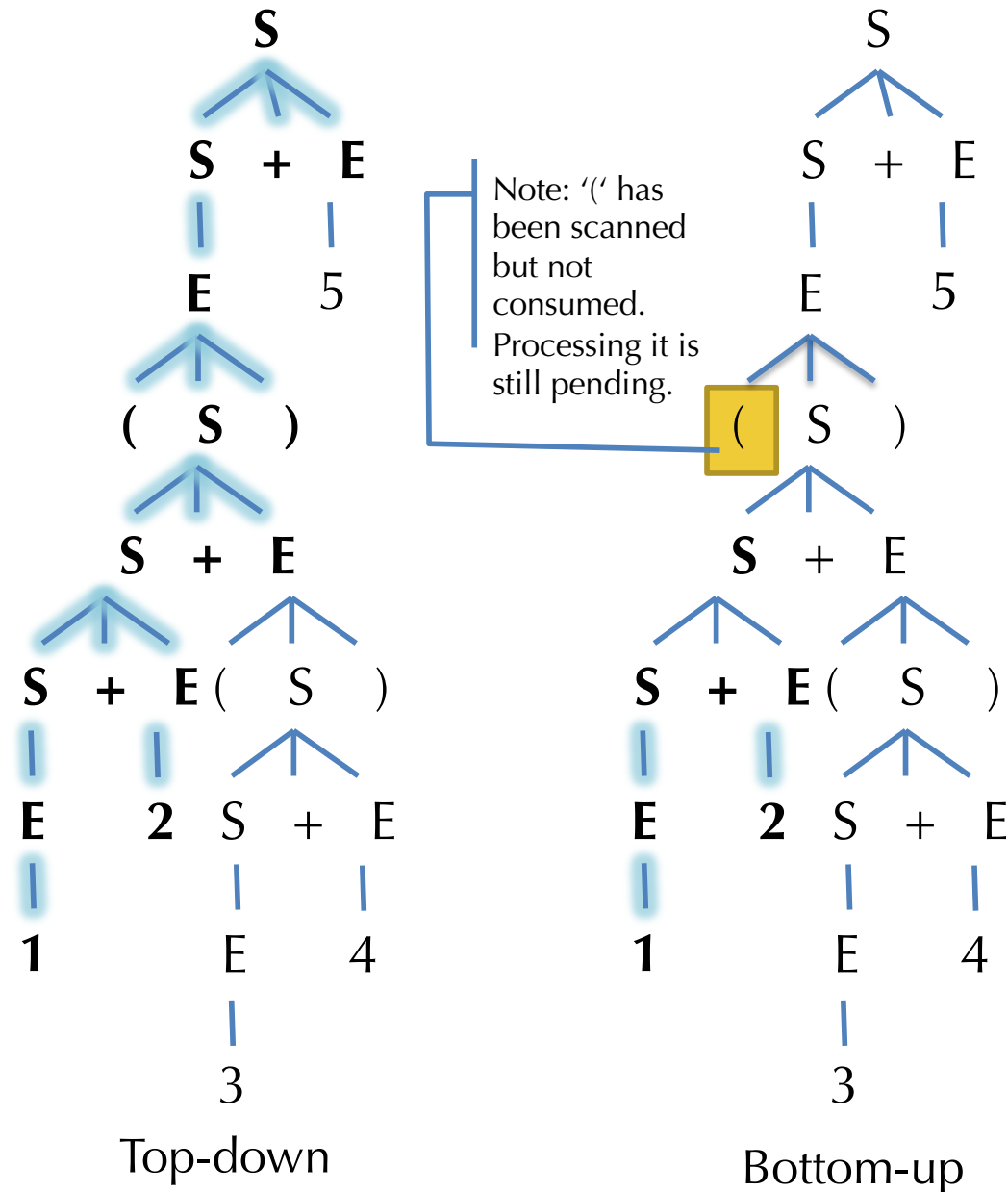


# Top-down vs. Bottom up

- 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



# Consider finding left-most derivations

Look at only **one** input symbol at a time

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

<u>Partly-derived String</u>	<u>Look-ahead</u>	<u>Parsed/Unparsed Input</u>
<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

- Want to decide *which production to apply* based on the *look-ahead symbol*

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

- But, consider

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

(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 is unclear whether to pick first
  - $S \mapsto E$                       or
  - $S \mapsto E + S$

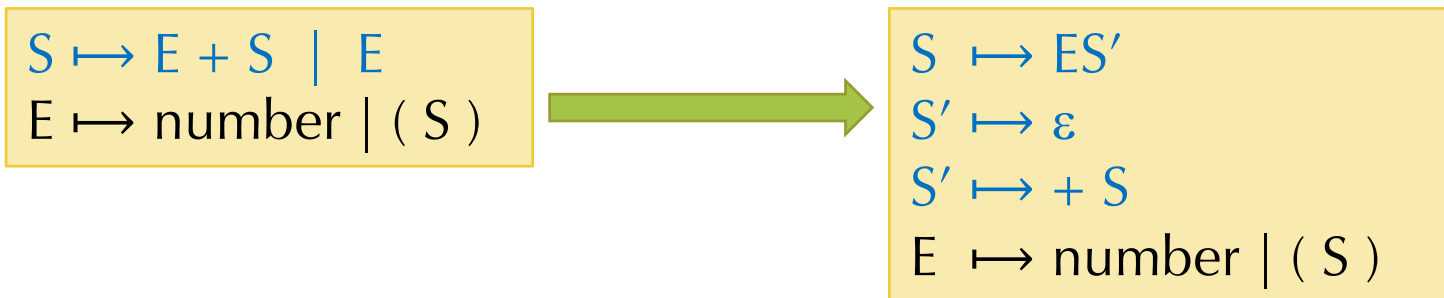
# LL(1) GRAMMARS

# Grammar is the problem

- Not all grammars can be parsed “top-down” with a single lookahead
  - *Top-down*
    - Start from the start symbol (root of the parse tree), and
    - Go 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?
- $$S \mapsto E + S \mid E$$
$$E \mapsto \text{number} \mid ( S )$$
- What can we do?

# 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



# Making a grammar LL(1)

- **Problem:** We can't decide which S production to apply until we see the symbol after the first expression
- Also need to **eliminate left-recursion** somehow. Why?
- Consider

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

# How to Remove Left Recursion?

- In general

$$S \rightarrow S \alpha_1 \mid \dots \mid S \alpha_n \mid \beta_1 \mid \dots \mid \beta_m$$

- Rewrite as

$$S \rightarrow \beta_1 S' \mid \dots \mid \beta_m S'$$

$$S' \rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \varepsilon$$

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



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



# LL(1) parse of the input string

Look at only **one** input symbol at a time

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

<u>Partly-derived String</u>	<u>Look-ahead</u>	<u>Parsed/Unparsed Input</u>
<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  $\times$  input token  $\rightarrow$  production**

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

	number	+	(	)	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \varepsilon$	$\mapsto \varepsilon$
E	$\mapsto \text{number}$		$\mapsto ( S )$		

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

# LL(1) parse of the input string

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

Partly-derived String

**T**  
 $\mapsto \underline{S}\$$   
 $\mapsto \underline{E} S'\$$   
 $\mapsto (\underline{S}) S'\$$   
 $\mapsto (\underline{E} S') S'\$$   
 $\mapsto (1 \underline{S}') S'\$$   
 $\mapsto (1 + \underline{S}) S'\$$   
 $\mapsto (1 + \underline{E} S') S'\$$   
 $\mapsto (1 + 2 \underline{S}') S'\$$   
 $\mapsto (1 + 2 + \underline{S}) S'\$$   
 $\mapsto (1 + 2 + \underline{E} S') S'\$$   
 $\mapsto (1 + 2 + (\underline{S})S') S'\$$

Look-ahead

(  
(  
(  
1  
1  
+  
2  
2  
+  
(  
(  
3

Parsed/Unparsed Input

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

# How do we construct the parse table?

- Consider a given production:  $A \rightarrow \gamma$
- **(Case 1)** 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
- **(Case 2)** 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
- **Important:** if there are two different productions for a given entry, the grammar is not LL(1)

# Example

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

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

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

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

# Converting the table to code

```
T  $\mapsto$  S$  
S  $\mapsto$  ES'  
S'  $\mapsto$   $\epsilon$   
S'  $\mapsto$  + S  
E  $\mapsto$  number | ( S )
```

- Define **N** mutually recursive functions
  - One for each nonterminal **A**: **parse\_A**
  - **parse\_A** of type **unit**  $\rightarrow$  **ast** if A is *not* an auxiliary nonterminal
  - Otherwise, **parse\_A** takes extra ast's as inputs, one for each nonterminal in the “factored” prefix

# Converting the table to code

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

- Each function **parse\_A**
  - “Peeks” at the lookahead token
  - 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 creates the ast after looking at more input)
  - Otherwise, this function builds the ast tree itself and returns it

# LL(1) Summary

- Top-down parsing that finds the leftmost derivation
- Language grammar  $\Rightarrow$   
LL(1) grammar  $\Rightarrow$   
prediction table  $\Rightarrow$   
recursive-descent parser
- Parser generator based on LL: ANTLR
- 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/alternative way?



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

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

## DEMO: PARSER.ML