Lecture 11 CIS 341: COMPILERS

#### Announcements

- HW3: LLVM lite
  - Available on the course web pages.
  - Due: Weds., February 23<sup>rd</sup> at 11:59:59pm

it is officially too late to **START EARLY!!** 

- Midterm: March 3<sup>rd</sup>
  - In class
  - One-page, letter-sized, double-sided "cheat sheet" of notes permitted
  - Coverage: interpreters / program transformers / x86 / calling conventions / IRs / LLVM / Lexing / Parsing
  - See examples of previous exams on the web pages

#### Parsing



# **Parsing: Finding Syntactic Structure**



# **CONTEXT FREE GRAMMARS**

Zdancewic CIS 341: Compilers

#### **Context-free Grammars**

• Here is a specification of the language of balanced parens:

$$S \mapsto (S)S$$
  
 $S \mapsto \varepsilon$ 

Note: Once again we have to take care to distinguish meta-language elements (e.g. "S" and " $\mapsto$ ") from object-language elements (e.g. "(").\*

- The definition is *recursive* S mentions itself.
- Idea: "derive" a string in the language by starting with S and rewriting according to the rules:

- Example:  $S \mapsto (S)S \mapsto ((S)S)S \mapsto ((\epsilon)S)S \mapsto ((\epsilon)S)\epsilon \mapsto ((\epsilon)\epsilon)\epsilon = (())$ 

- You can replace the "nonterminal" S by one of its definitions anywhere
- A context-free grammar accepts a string iff there is a derivation from the start symbol

## **CFGs Mathematically**

- A Context-free Grammar (CFG) consists of
  - A set of *terminals* (e.g., a lexical token or  $\varepsilon$ )
  - 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:

$$S \longmapsto (S)S$$
$$S \longmapsto \varepsilon$$

• How many terminals? How many nonterminals? Productions?

#### **Another Example: Sum Grammar**

• A grammar that accepts parenthesized sums of numbers:

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

e.g.: (1 + 2 + (3 + 4)) + 5

• Note the vertical bar '|' is shorthand for multiple productions:

```
S \mapsto E + SS \mapsto EE \mapsto numberE \mapsto (S)
```

4 productions 2 nonterminals: S, E 4 terminals: (, ), +, number Start symbol: S

#### **Derivations in CFGs**

- Example: derive (1 + 2 + (3 + 4)) + 5
- $\underline{\mathbf{S}} \mapsto \underline{\mathbf{E}} + \mathbf{S}$ 
  - $\longmapsto (\underline{\mathbf{S}}) + S$
  - $\longmapsto (\underline{\textbf{E}} + S) + S$
  - $\mapsto (1 + \underline{S}) + S$
  - $\longmapsto (1 + \underline{\mathbf{E}} + \mathbf{S}) + \mathbf{S}$
  - $\mapsto (1 + 2 + \mathbf{S}) + \mathbf{S}$
  - $\mapsto$  (1 + 2 + **E**) + S
  - $\mapsto (1 + 2 + (\mathbf{S})) + \mathbf{S}$
  - $\longmapsto (1 + 2 + (\underline{\mathbf{E}} + S)) + S$
  - $\mapsto (1 + 2 + (3 + \underline{\mathbf{S}})) + S$  $\mapsto (1 + 2 + (3 + \mathbf{E})) + S$

 $\mapsto$  (1 + 2 + (3 + 4)) + **E** 

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

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

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

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

 $\alpha A \gamma \mapsto \alpha \beta \gamma$ 

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

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

Note: Underline indicates symbol being expanded.

#### **From Derivations to Parse Trees**

- Tree representation of the derivation
- Leaves of the tree are terminals
  - In-order traversal yields the input sequence of tokens
- Internal nodes: nonterminals
- No information about the *order* of the derivation steps

• 
$$(1 + 2 + (3 + 4)) + 5$$



CIS 341: Compilers

#### **From Parse Trees to Abstract Syntax**



### **Derivation Orders**

- Productions of the grammar can be applied in any order.
- There are two standard orders:
  - *Leftmost derivation*: Find the left-most nonterminal and apply a production to it.
  - *Rightmost derivation*: Find the right-most nonterminal and apply a production there.
- Note that both strategies (and any other) yield the same parse tree!
  - Parse tree doesn't contain the information about what order the productions were applied.

#### **Example: Left- and rightmost derivations**

- Leftmost derivation:
- $\mathbf{S} \mapsto \mathbf{E} + \mathbf{S}$  $\mapsto$  (**S**) + S  $\mapsto$  (**E** + S) + S  $\mapsto$  (1 + **S**) + S  $\mapsto$  (1 + **E** + S) + S  $\mapsto$  (1 + 2 + **S**) + S  $\mapsto$  (1 + 2 + **E**) + S  $\mapsto$  (1 + 2 + (**S**)) + S  $\mapsto$  (1 + 2 + (**E** + S)) + S  $\mapsto$  (1 + 2 + (3 + **S**)) + S  $\mapsto$  (1 + 2 + (3 + **E**)) + S  $\mapsto$  (1 + 2 + (3 + 4)) + **S**  $\mapsto$  (1 + 2 + (3 + 4)) + **E**  $\mapsto$  (1 + 2 + (3 + 4)) + 5

Rightmost derivation:

 $\mathbf{S} \mapsto \mathbf{E} + \mathbf{S}$  $S \mapsto E + S \mid E$ ⊷ E + **E**  $E \mapsto number \mid (S)$  $\mapsto$  **E** + 5  $\mapsto$  (**S**) + 5  $\mapsto$  (E + **S**) + 5  $\mapsto$  (E + E + **S**) + 5  $\mapsto$  (E + E + **E**) + 5  $\mapsto$  (E + E + (**S**)) + 5  $\mapsto$  (E + E + (E + **S**)) + 5  $\mapsto$  (E + E + (E + E)) + 5  $\mapsto (\mathsf{E} + \mathsf{E} + (\mathbf{E} + 4)) + 5$  $\mapsto$  (E + E + (3 + 4)) + 5  $\mapsto$  (**E** + 2 + (3 + 4)) + 5  $\mapsto$  (1 + 2 + (3 + 4)) + 5

# **Loops and Termination**

- Some care is needed when defining CFGs
- Consider:



- This grammar has nonterminal definitions that are "nonproductive".
   (i.e. they don't mention any terminal symbols)
- There is no finite derivation starting from S, so the language is empty.
- Consider:  $S \mapsto (S)$ 
  - This grammar is productive, but again there is no finite derivation starting from S, so the language is empty
- Easily generalize these examples to a "cycle" of many nonterminals, which can be harder to find in a large grammar
- Upshot: be aware of "vacuously empty" CFG grammars.
  - Every nonterminal should eventually rewrite to an alternative that contains only terminal symbols.

Associativity, ambiguity, and precedence.

# GRAMMARS FOR PROGRAMMING LANGUAGES

#### Associativity

Consider the input: 1 + 2 + 3

Leftmost derivation: Rightmost derivation:

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

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



 $S \mapsto E + S \mid E$ 

 $E \mapsto number \mid (S)$ 



CIS 341: Compilers

#### Associativity

- This grammar makes '+' *right associative*...
  - i.e., the abstract syntax tree is the same for both
    - 1 + 2 + 3 and 1 + (2 + 3)
- Note that the grammar is *right recursive*...



S refers to itself on the right of +

- How would you make '+' left associative?
- What are the trees for "1 + 2 + 3"?



• Consider this grammar:

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

- Claim: it accepts the *same* set of strings as the previous one.
- What's the difference?
- Consider these *two* leftmost derivations:

$$- \underline{\mathbf{S}} \mapsto \underline{\mathbf{S}} + \mathbf{S} \mapsto \mathbf{1} + \underline{\mathbf{S}} \mapsto \mathbf{1} + \underline{\mathbf{S}} + \mathbf{S} \mapsto \mathbf{1} + \mathbf{2} + \underline{\mathbf{S}} \mapsto \mathbf{1} + \mathbf{2} + \mathbf{3}$$

$$- \underline{\mathbf{S}} \mapsto \underline{\mathbf{S}} + \mathbf{S} \mapsto \underline{\mathbf{S}} + \mathbf{S} + \mathbf{S} \mapsto \mathbf{1} + \underline{\mathbf{S}} + \mathbf{S} \mapsto \mathbf{1} + \mathbf{2} + \underline{\mathbf{S}} \mapsto \mathbf{1} + \mathbf{2} + \mathbf{3}$$

- One derivation gives left associativity, the other gives right associativity to '+'
  - Which is which?



## Why do we care about ambiguity?

- The '+' operation is associative, so it doesn't matter which tree we pick. Mathematically, x + (y + z) = (x + y) + z
  - But, some operations aren't associative. Examples?
  - Some operations are only left (or right) associative. Examples?
- Moreover, if there are multiple operations, ambiguity in the grammar leads to ambiguity in their *precedence*
- Consider:

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

- Input: 1 + 2 \* 3
  - One parse = (1 + 2) \* 3 = 9
  - The other = 1 + (2 \* 3) = 7



# **Eliminating Ambiguity**

- We can often eliminate ambiguity by adding nonterminals and allowing recursion only on the left (or right) .
- Higher-precedence operators go *farther* from the start symbol.
- Example:

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

- To disambiguate:
  - Decide (following math) to make '\*' higher precedence than '+'
  - Make '+' left associative
  - Make '\*' right associative
- Note:
  - S<sub>2</sub> corresponds to 'atomic' expressions

$$S_0 \mapsto S_0 + S_1 | S_1$$
  

$$S_1 \mapsto S_2 * S_1 | S_2$$
  

$$S_2 \mapsto \text{number} | (S_0)$$

#### **Context Free Grammars: Summary**

- Context-free grammars allow concise specifications of programming languages.
  - An unambiguous CFG specifies how to parse: convert a token stream to a (parse tree)
  - Ambiguity can (often) be removed by encoding precedence and associativity in the grammar.
- Even with an unambiguous CFG, there may be more than one derivation
  - Though all derivations correspond to the same abstract syntax tree.
- Still to come: finding a derivation
  - But first: menhir

parser.mly, lexer.mll, range.ml, ast.ml, main.ml

# **DEMO: BOOLEAN LOGIC**

Zdancewic CIS 341: Compilers

Searching for derivations.

# LL & LR PARSING

Zdancewic CIS 341: Compilers

## **CFGs Mathematically**

- A Context-free Grammar (CFG) consists of
  - A set of *terminals* (e.g., a token or  $\varepsilon$ )
  - 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:

$$S \mapsto (S)S$$
$$S \mapsto \varepsilon$$

• How many terminals? How many nonterminals? Productions?

## **Consider finding left-most derivations**

• Look at only one input symbol at a time. S

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

Partly-derived String	Look-ahead	Parsed/Unparsed Input
<u>S</u>	(	(1 + 2 + (3 + 4)) + 5
$\mapsto \underline{\mathbf{E}} + \mathbf{S}$	(	(1 + 2 + (3 + 4)) + 5
$\mapsto (\underline{\mathbf{S}}) + \mathbf{S}$	1	(1 + 2 + (3 + 4)) + 5
$\longmapsto (\underline{\mathbf{E}} + \mathbf{S}) + \mathbf{S}$	1	(1 + 2 + (3 + 4)) + 5
$\longmapsto (1 + \underline{\mathbf{S}}) + \mathbf{S}$	2	(1 + 2 + (3 + 4)) + 5
$\longmapsto (1 + \underline{\mathbf{E}} + \mathbf{S}) + \mathbf{S}$	2	(1 + 2 + (3 + 4)) + 5
$\mapsto (1 + 2 + \underline{\mathbf{S}}) + \mathbf{S}$	(	(1 + 2 + (3 + 4)) + 5
$\mapsto (1 + 2 + \underline{\mathbf{E}}) + \mathbf{S}$	(	(1 + 2 + (3 + 4)) + 5
$\mapsto (1 + 2 + (\underline{\mathbf{S}})) + \mathbf{S}$	3	(1 + 2 + (3 + 4)) + 5
$\mapsto (1 + 2 + (\underline{\mathbf{E}} + \mathbf{S})) + \mathbf{E} + \mathbf{E}$	S 3	(1 + 2 + (3 + 4)) + 5
$\mapsto \dots$		

### There is a problem

 $S \mapsto E + S \mid E$ 

 $E \mapsto number \mid (S)$ 

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

(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

Zdancewic CIS 341: Compilers

#### **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
  - <u>L</u>eft-most derivation,
  - <u>1</u> lookahead symbol
- This language isn't "LL(1)"
- Is it LL(k) for some k?

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

• 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:

$$\begin{array}{c|c} S \mapsto E + S & | & E \\ E \mapsto number \mid (S) \end{array} \xrightarrow{} & S \mapsto ES' \\ S' \mapsto \varepsilon \\ S' \mapsto + S \\ E \mapsto number \mid (S) \end{array}$$

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

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

# LL(1) Parse of the input string

• Look at only one input symbol at a time.

$$S \mapsto ES'$$

$$S' \mapsto \varepsilon$$

$$S' \mapsto + S$$

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

Partly-derived String	Look-ahead	Parsed/Unparsed Input
<u>S</u>	(	(1 + 2 + (3 + 4)) + 5
→ <u>E</u> S′	(	(1 + 2 + (3 + 4)) + 5
$\mapsto (\underline{\mathbf{S}}) S'$	1	(1 + 2 + (3 + 4)) + 5
$\longmapsto (\underline{\mathbf{E}} S') S'$	1	(1 + 2 + (3 + 4)) + 5
→ (1 <u><b>S'</b></u> ) S'	+	(1 + 2 + (3 + 4)) + 5
$\mapsto (1 + \underline{\mathbf{S}}) \mathbf{S'}$	2	(1 + 2 + (3 + 4)) + 5
$\longmapsto (1 + \underline{\mathbf{E}} S') S'$	2	(1 + 2 + (3 + 4)) + 5
$\mapsto (1 + 2 \mathbf{\underline{S'}}) \mathbf{S'}$	+	(1 + 2 + (3 + 4)) + 5
$\mapsto (1 + 2 + \underline{\mathbf{S}}) \mathbf{S'}$	(	(1 + 2 + (3 + 4)) + 5
$\longmapsto (1 + 2 + \underline{\mathbf{E}} S') S'$	(	(1 + 2 + (3 + 4)) + 5
$\longmapsto (1 + 2 + (\underline{\mathbf{S}})\mathbf{S'}) \mathbf{S'}$	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

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

	number	+	(	)	\$ (EOF)
Т	$\mapsto$ S\$		⊷S\$		
S	$\mapsto E S'$		⊷E S′		
S'		$\mapsto + S$		$\mapsto \epsilon$	$\mapsto \epsilon$
E	⊢ 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  $\varepsilon$  (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, token) for each such token.

• Note: if there are two different productions for a given entry, the grammar is not LL(1)

#### Example

First(T) = First(S)ullet $T \mapsto S$ \$ First(S) = First(E)۲  $S \mapsto ES'$  $First(S') = \{ + \}$ ٠  $S' \mapsto \varepsilon$  $First(E) = \{ number, '(') \}$ ۲  $S' \mapsto + S$  $E \mapsto number \mid (S)$ Follow(S') = Follow(S)٠ **Note:** we want the *least* Follow(S) = { \$, ')' } U Follow(S') solution to this system of set • equations... a *fixpoint* computation. More on these later in the course. number \$ (EOF) +  $\mapsto$  S\$ →S\$ Τ  $\mapsto E S'$  $\mapsto E S'$ S  $\mapsto$  + S **S'**  $\mapsto \epsilon$  $\mapsto \epsilon$  $\mapsto$  (S) E  $\mapsto$  num.

### **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)
Т	$\mapsto$ S\$		⊢→S\$		
S	$\mapsto E S'$		⊢→E S′		
S'		$\mapsto$ + S		$\mapsto \epsilon$	$\mapsto \epsilon$
E	⊢ 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 ⇒ LL(1) grammar ⇒ prediction table ⇒ recursivedescent 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**

Zdancewic CIS 341: Compilers

### **Bottom-up Parsing (LR Parsers)**

- LR(k) parser:
  - Left-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, CUP, ocamlyacc, merlin, etc.)
  - Better error detection/recovery

#### **Top-down vs. Bottom up**

• Consider the leftrecursive grammar:

> $S \mapsto S + E \mid E$ E \low number | (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...



#### **Progress of Bottom-up Parsing**

Reductions	Scanned	Input Remaining
$(1 + 2 + (3 + 4)) + 5 \leftarrow 1$		(1 + 2 + (3 + 4)) + 5
$(\underline{\mathbf{E}} + 2 + (3 + 4)) + 5 \longleftarrow$	(	1 + 2 + (3 + 4)) + 5
$(\underline{\mathbf{S}} + 2 + (3 + 4)) + 5 \longleftarrow$	(1	+2+(3+4))+5
$(\mathbf{S} + \mathbf{\underline{E}} + (3 + 4)) + 5 \longleftarrow$	(1 + 2	+(3+4))+5
$(\underline{\mathbf{S}} + (3 + 4)) + 5 \longleftarrow$	(1 + 2	+(3+4))+5
$(S + (\underline{E} + 4)) + 5 \longleftarrow$	(1 + 2 + (3 + 3))	(+ 4)) + 5
$(S + (\underline{S} + 4)) + 5 \longleftarrow$	(1 + 2 + (3 + (3 + (3 + (3 + (3 + (3 + (3	(+ 4)) + 5
$(S + (S + \underline{E})) + 5 \longleftarrow$	(1 + 2 + (3 + 4))	)) + 5
$(S + (\underline{S})) + 5 \longleftarrow$	(1 + 2 + (3 + 4))	)) + 5
$(S + \underline{E}) + 5 \longleftarrow$	(1 + 2 + (3 + 4))	) + 5
( <u><b>S</b></u> ) + 5 ↔	(1 + 2 + (3 + 4))	) + 5
<u><b>E</b></u> + 5 ↔	(1 + 2 + (3 + 4))	+ 5
<u><b>S</b></u> + 5 ↔	(1 + 2 + (3 + 4))	+ 5
S + <u>E</u> ←	(1 + 2 + (3 + 4)) + 5	
S		

# **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
- Reduce: Replace symbols  $\gamma$  at top of stack with nonterminal X such that X  $\mapsto \gamma$  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: $E \mapsto number$
(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: $E \mapsto number$



Simple LR parsing with no look ahead.

# LR(0) GRAMMARS

Zdancewic CIS 341: Compilers

#### **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  $\boldsymbol{\sigma}.$
  - Accept states of the DFA correspond to unique reductions that apply.
- Example: LR(0) parsing
  - <u>L</u>eft-to-right scanning, <u>R</u>ight-most derivation, <u>zero</u> 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:

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

• Example strings:

x (x,y) ((((x)))) (x, (y, z), w) (x, (y, (z, w)))

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



# **Shift/Reduce Parsing**

- 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.

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

• 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 $L \mapsto S$

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

### **Example Run**

 $S \mapsto id$ 

 $L \mapsto S$ 

 $S \mapsto id$ 

 $\mathsf{L} \mapsto \mathsf{S}$ 

 $S \mapsto id$ 

 $L \mapsto L, S$ 

 $S \longmapsto (\ L \ )$ 

 $L \mapsto L, S$ 

Stack	Input	Action
	(x, (y, z), w)	shift (
(	x, (y, z), w)	shift x
(x	, (y, z), w)	reduce
(S	, (y, z), w)	reduce
(L	, (y, z), w)	shift ,
(L,	(y, z), w)	shift (
(L, (	y, z), w)	shift y
(L, (y	, z), w)	reduce
(L, (S	, z), w)	reduce
(L, (L	, z), w)	shift ,
(L, (L,	z), w)	shift z
(L, (L, z	), W)	reduce
(L, (L, S	), W)	reduce
(L, (L	), W)	shift )
(L, (L)	, w)	reduce
(L, S	, w)	reduce
CIS 34 (Compilers	, W)	shift ,

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

### **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,  $X \mapsto \varepsilon$  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

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

- 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 γ'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:

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

- 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.



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



- 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**



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

- 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**



 $\begin{array}{l} S' \longmapsto S \\ S \longmapsto (L) &| id \\ L \longmapsto S &| L, S \end{array}$ 

- Finally, for each new state, we take the closure.
- Note that we have to perform two iterations to compute CLOSURE({S → (.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**



9

 $L \mapsto L, S.$ 

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

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
$\varepsilon_1(_3$	x, (y, z), w)\$	s2
$\varepsilon_1(_3X_2$	, (y, z), w)\$	Reduce: S⊷id
$\epsilon_1(_3S)$	, (y, z), w)\$	g7 (from state 3 follow S)
$\epsilon_1(_3S_7$	, (y, z), w)\$	Reduce: L→S
$\epsilon_1(_3L)$	, (y, z), w)\$	g5 (from state 3 follow L)
$\varepsilon_1(_3L_5$	, (y, z), w)\$	s8
$\epsilon_1(_3L_{5\prime 8})$	(y, z), w)\$	s3
$\varepsilon_1(_3L_{5,8}(_3$	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:

OKshift/reducereduce/reduce
$$S \mapsto (L).$$
 $S \mapsto (L).$  $S \mapsto L, S.$  $L \mapsto .L, S$  $S \mapsto ,S.$ 

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

#### **Examples**

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



- 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 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  ${\cal M}$  also contains  ${\cal L}$

### **Example Closure**

```
S' \mapsto S

S \mapsto E + S \mid E

E \mapsto number \mid (S)
```

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

# **Using the DFA**



- 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 '.'

Fragment of the Action & Goto tables

#### 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 lookahead sets:  $s' \mapsto ss = 0$



- 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**

