## CSC488/2107 Winter 2019 – Compilers & Interpreters

https://www.cs.toronto.edu/~csc488h/

Peter McCormick pdm@cs.toronto.edu

## Agenda

- Recognize, Analyze, Transform
- Lexical analysis
- Building lexical analyzers

### Recognize Analyze Transform



## Syntax Analysis

- The syntax of a language defines the rules by which a sequence of tokens can be recognized as a legal construction in that language.
- Recognize and distinguish legal and illegal sentences of the language

## Addition expressions

# Addition expressions tokens as regular expressions

Digits = '0' ... '9' <u>Plus</u> = '+' <u>Literal</u> = Digits+

Token = Plus | Literal

### Addition expressions as regular expressions

### Digits = $('0' ... '9')^+$ Expr = $(Digits '+')^*$ Digits

### Digits = ( '0' ... '9' )+ Expr = ( Digits '+' )\* Digits



### $Expr = (('0' ... '9')^{+} '+')^{*} ('0' ... '9')^{+}$

### **Expressions with parentheses**

488 (400+88) (400+(44+44)) ((400+4)+(42+42))



Abbreviations in regular expressions are just syntactic sugar...

They add no additional expressive power

### Recursion adds expressive power

## **Recursive abbreviations?**

- Simplifies regular expressions
  - Alternation within expressions is no longer required
  - Alternation pipe | is no longer required
  - Kleene closure is no longer required

## **Recursive abbreviations**



## **Recursive abbreviations**



## **Recursive abbreviations**







This simple but powerful notation of recursive abbreviations is referred to as *context-free grammars* (CFG)

## Context-free grammars

- CFGs can describe richer languages than regular expressions
  - Thus need something more powerful than finite automata to recognize them (key is having *memory*)
- A *language* is a set of strings over an alphabet  $\Sigma$
- Alphabet  $\Sigma$  ranges over *tokens*, not characters
  - Example:  $\Sigma = \{ IF, IDENT, PLUS, LBRACE, EQ, ... \}$
- A CFG consists of a set of *productions*

## Context-free grammars

A *production* is of the form:

### symbol $\longrightarrow$ symbol symbol ... symbol

- Right hand side has 0 or more symbols (0 means  $\lambda$ )
- A *symbol* is either:
  - Terminal: a token from alphabet  $\Sigma$
  - Non-terminal: appears on the left hand side of a production
- Left hand side symbol is always a non-terminal
- Distinguished start production (typically the first)

- $S \longrightarrow S$  ';' S $S \longrightarrow id$  '=' E
- $S \longrightarrow$  'print' E
- $E \longrightarrow id$
- $E \longrightarrow num$
- $E \longrightarrow E'+'E$
- E → '(' E ')'

### id = num ; id = num ; print id + ( id + num )

x = 400; y = 42; print x + (y + 46)

## Derivations

- Derive a sentence from the grammar to show that it is in the language
- Begin with the start symbol, and repeatedly replace right hand side non-terminals with symbols from productions
- Many possible derivation orders
  - Left-most derivation: always expand the left-most nonterminal first, working towards the right
  - Right-most derivation



- $S \longrightarrow S$  ';' S  $S \longrightarrow id$  '=' E
- $S \longrightarrow `print' E$
- $E \longrightarrow id$
- $E \longrightarrow num$
- $E \longrightarrow E' + E'$
- E → '(' E ')'



num

## Ambiguous grammars

Two possible parse trees for the same sentence



<u>(1 + 2) + 3</u>

<u>1 + (2 + 3)</u>

## Removing ambiguity

We want left-associativity



T for terms because they are added together

### Removing ambiguity 1 + 2 + 3 $S \longrightarrow E$ $E \longrightarrow T$ S $E \longrightarrow E + T$ Ε Ε $T \longrightarrow id$ Ε num $T \longrightarrow num$ num $T \longrightarrow (E)$

num

## More ambiguity

We want multiplication & division to be higher-precedence, or to bind more tightly



F for factors because they are multiplied together

## Impossible to derive

- $S \longrightarrow E$
- $E \longrightarrow T$
- $E \longrightarrow E + T$
- $E \longrightarrow E T$
- $T \longrightarrow T * F$  $T \longrightarrow T / F$  $T \longrightarrow F$
- $F \longrightarrow id$
- $F \longrightarrow num$
- $\mathsf{F} \longrightarrow$  (  $\mathsf{E}$  )



### **Backus-Naur Form (BNF)**

### Backus-Naur Form (BNF) – ALGOL 60 Report

<expression></expression>	= : :   	<term> <expression> "+" <term> <expression> "-" <term></term></expression></term></expression></term>
<term></term>	: :=   	<factor> <term> "*" <factor> <term> "/" <factor></factor></term></factor></term></factor>
<factor></factor>	::=	<pre><identifier> <number> "(" <expression> ")"</expression></number></identifier></pre>

## PLY (Python Lex-Yacc)

expression		term expression PLUS term expression MINUS term
term	•	factor term STAR factor term SLASH factor
factor	•	identifier number LPAREN expression RPAREN

## Parsing Top Down vs Bottom Up

## **Top Down Parsing**

- Starting from the start symbol S, find the correct sequence of production expansions that will transform S into the input token stream (if possible)
  - Called a *derivation*
- Build the parse tree from the root (S) to the leaves (terminals)
- Top down techniques: LL(k), Recursive Descent

## **Bottom Up Parsing**

- Given a stream of input tokens, find the correct sequence of production contractions that take the input back to the start symbol
  - Called a *reduction* or *reverse derivation*
- Build the parse tree from the leaves (terminals) to the root (S)
- Bottom up techniques: LR(k), LALR(k), SLR(k)

### LL(1) Grammars

#### Peter McCormick

January 24, 2018

(ロ)、(型)、(E)、(E)、 E) の(()

#### Context-free Grammars

• Define a *context-free grammar*  $G = (N, \Sigma, S, P)$  where

- Σ is the set of *terminals* (the *alphabet*), each represented as lower case Roman letters (t,x,y,z)
- ► *N* is the set of *non-terminals*, each represented by capitalized Roman letters (*A*, *B*, *X*, *Y*)
- $S \in N$  is the distinguished *start* symbol
- P is a finite set of productions
- ►  $V = \Sigma \cup N$  is the vocabulary of the grammar. Strings ranging over V, such as A w X t, are each represented by Greek letters  $(\alpha, \beta, \gamma)$

- X ⇒<sup>+</sup> α means that there is a valid sequence of derivations starting from the non-terminal X to the string α
- Let \$ be a special *end-of-input* marker

#### Productions

A production is a *rewriting rule* of the form:  $X \rightarrow \alpha_1 \alpha_2 \cdots \alpha_m$ where  $X \in N$ ,  $\alpha_i \in V^*$  for each *i*, and  $m \ge 0$ . If m = 0, then  $X \rightarrow$ That is, X can be replaced by  $\lambda$ , the empty string.

- *Rewriting* is the act of replacing a non-terminal X with the right hand of a production for X, so  $\alpha_1 \cdots \alpha_m$  replacing X
- $X \Rightarrow^+ \alpha$  means that there is a valid sequence of derivations starting from the non-terminal X to the vocabulary string  $\alpha$

#### Definitions

A vocabulary string is in *sentential form* if it is in the set of all strings that can be derived from the start symbol:

$$\{w: S \Rightarrow^* w\}$$

A *language* is the set of all terminal strings that can be derived from the start symbol:

$$\{w: S \Rightarrow^* w\} \cap \Sigma^*$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

#### Definitions

Given a production  $X \rightarrow \gamma$  and a sentenial form  $\alpha X \beta$ , then

 $\begin{array}{ll} \alpha X\beta & \Rightarrow & \alpha \gamma \beta \text{ is a derivation in one step} \\ \alpha X\beta & \Rightarrow^* & \alpha \gamma \beta \text{ is a derivation in zero or more steps} \\ \alpha X\beta & \Rightarrow^+ & \alpha \gamma \beta \text{ is a derivation in one or more steps} \end{array}$ 

• A *left-most derivation* expands non-terminals left to right, while a *right-most derivation* expands right to left

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

### LL(1) Grammars

• LL is a set of all languages that can be parser by an LL parser

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- An *LL* parser consumes its input Left-to-right, producing a Leftmost-derivation
- LL(k) means LL with k tokens worth of input look ahead
- LL(1) means 1 input token lookahead

#### Nullable

 $\alpha \in V^+$  is *nullable* iff  $\alpha \Rightarrow^* \lambda$ 



#### Nullable

Given the productions:  $X \rightarrow Y Z W$  $Y \rightarrow \lambda$  $Z \rightarrow z | Y$  $W \rightarrow w | YZ$ Then X is nullable since:  $X \Rightarrow YZW$  $\Rightarrow \underline{Z}W$  (since Y is nullable)  $\Rightarrow YW$  $\Rightarrow$  <u>W</u> (again since Y is nullable)  $\Rightarrow YZ$  $\Rightarrow$  <u>Z</u> (again)  $\Rightarrow \underline{Y} \Rightarrow \lambda$ 

### First Sets

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = ● ● ●

#### Follow Sets

### $Follow(X) = \left\{ t \in \Sigma : S \Rightarrow^+ \alpha X \, t \, \beta \right\}$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

#### Follow Set Construction Rules

Starting from initially empty follow sets, iteratively apply these rules until the sets no longer change. #1. If S is the goal symbol Add {\$} to Follow(S) #2. If  $S \Rightarrow^+ \alpha X$  Add {\$} to Follow(X) #3. If  $S \Rightarrow^+ \alpha X \mathbf{t}\beta$  Add {t} to Follow(X) #4. If  $S \Rightarrow^+ \alpha X \mathbf{t}\beta$  If Y is not nullable: · Add First(Y) to Follow(X) Else if Y is nullable: · Add First(Y)  $\cup$  First( $\beta$ ) to Follow(X) #5. If  $X \rightarrow \alpha \mathbf{Y}$  Add Follow(X) to Follow(Y)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

### Predict Sets for LL(1) Grammars

Given a non-terminal X defined with several alternate productions:  $\begin{array}{cccc} X & \rightarrow & \gamma_1 \\ & \rightarrow & \gamma_2 \\ & & \cdots \\ & \rightarrow & \gamma_m \end{array}$ The predict set for each production  $X \rightarrow \gamma_i$  is defined as  $\begin{array}{cccc} Predict(X \rightarrow \gamma_i) &= & First(\gamma_i) & (\text{if } \gamma_i \text{ is not nullable}) \\ &= & First(\gamma_i) \cup Follow(X) & (\text{if } \gamma_i \text{ is nullable}) \end{array}$ 

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

### Predict Sets for LL(1) Grammars

• To be an *LL*(1) grammar, the Predict sets for all productions for a given non-terminal *X* must be mutually disjoint

$$X \rightarrow Y \alpha$$
  
 $\rightarrow Y \beta$ 

is not LL(1) because the Predict sets for the productions of X are not mutually disjoint:

$$\begin{aligned} Predict (X \to Y \alpha) &= First (Y \alpha) \text{ (if } Y \text{ is } not \text{ nullable}) \\ &= First (Y) \\ &= First (Y \beta) \\ &= Predict (X \to Y \beta) \end{aligned}$$
$$\begin{aligned} Predict (X \to Y \alpha) &= First (Y \alpha) \cup Follow (X) \text{ (if } Y \text{ is nullable}) \\ Predict (X \to Y \beta) &= First (Y \beta) \cup Follow (X) \text{ (if } Y \text{ is nullable}) \end{aligned}$$

### Factoring Out Common Prefix



#### Left Recursion



#### Factoring Out Left Recursion



▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ 三重 - 釣A(?)

## Next Week

- More on practical parser construction
- Bottom up parsing
- LR(1) grammars