### CS345H: Programming Languages

#### Lecture 5: Introduction to Parsing

Thomas Dillig

Limitations of Regular Languages

Limitations of Regular Languages

Parser Overview

- Limitations of Regular Languages
- Parser Overview
- Context-free Grammars (CFGs)

- Limitations of Regular Languages
- Parser Overview
- Context-free Grammars (CFGs)
- Derivations

- Limitations of Regular Languages
- Parser Overview
- Context-free Grammars (CFGs)
- Derivations
- Ambiguity

 Last time, we saw that regular languages are very useful for partitioning input into tokens

- Last time, we saw that regular languages are very useful for partitioning input into tokens
- But regular languages are not expressive enough to turn a stream of tokens into structure

## **Regular Languages**

- Last time, we saw that regular languages are very useful for partitioning input into tokens
- But regular languages are not expressive enough to turn a stream of tokens into structure
- ▶ For this, we need a more expressive formal language

## Beyond Regular Languages

Many languages are not regular

# Beyond Regular Languages

- Many languages are not regular
- Classic Example: Strings of balanced parenthesis:

 $\{(^i)^j \mid i \ge 0\}$ 

# Beyond Regular Languages

- Many languages are not regular
- Classic Example: Strings of balanced parenthesis:

$$\{(^i)^j \mid i \ge 0\}$$

Question: Why is there no automata that can recognize this language?

What Can Regular Languages Express?

Languages requiring counting modulo a fixed integer

# What Can Regular Languages Express?

- Languages requiring counting modulo a fixed integer
- Intuition: A finite automaton that runs long enough must repeat states

# What Can Regular Languages Express?

- Languages requiring counting modulo a fixed integer
- Intuition: A finite automaton that runs long enough must repeat states
- Finite automaton cannot remember the number of times it has visited a particular state

## Side Note: Comments in L

#### Recall: Comments in L start with (\*, end with \*) and can be nested

# Side Note: Comments in L

- Recall: Comments in L start with (\*, end with \*) and can be nested
- ► Also Recall: Comments are removed during lexing

### Side Note: Comments in L

- Recall: Comments in L start with (\*, end with \*) and can be nested
- ► Also Recall: Comments are removed during lexing
- Question: Are comments in L a regular language?

The Functionality of the Parser

Input: sequence of tokens from the lexer

### The Functionality of the Parser

- Input: sequence of tokens from the lexer
- Output: parse tree of the program

### Example

Consider the following L expression: if x<>y then 1 else 2

#### Example

- Consider the following L expression: if x<>y then 1 else 2
- Parse Input: TOKEN\_IF TOKEN\_ID("x") TOKEN\_NEQ TOKEN\_ID("y") TOKEN\_THEN TOKEN\_INT(1) TOKEN\_ELSE TOKEN\_INT(2)

#### Example

- Consider the following L expression: if x<>y then 1 else 2
- Parse Input: TOKEN\_IF TOKEN\_ID("x") TOKEN\_NEQ TOKEN\_ID("y") TOKEN\_THEN TOKEN\_INT(1) TOKEN\_ELSE TOKEN\_INT(2)



Phase	Input	Output
Lexer		
Parser		

Phase	Input	Output
Lexer	String of characters	
Parser		

Phase	Input	Output
Lexer	String of characters	String of tokens
Parser		

Phase	Input	Output
Lexer	String of characters	String of tokens
Parser	String of tokens	

Phase	Input	Output
Lexer	String of characters	String of tokens
Parser	String of tokens	Parse tree

Not all strings of tokens are programs ....

- Not all strings of tokens are programs ....
- Parser must distinguish between valid and invalid strings of tokens

- Not all strings of tokens are programs ....
- Parser must distinguish between valid and invalid strings of tokens
- ► We need:

- Not all strings of tokens are programs ....
- Parser must distinguish between valid and invalid strings of tokens
- We need:
  - A language for describing valid strings of tokens

- Not all strings of tokens are programs ....
- Parser must distinguish between valid and invalid strings of tokens
- We need:
  - A language for describing valid strings of tokens
  - A method for recognizing if a string of tokens is in this language or not

# Context-free Grammars (CFGs)

Programming language constructs have recursive structure

# Context-free Grammars (CFGs)

- Programming language constructs have recursive structure
- Example: An L expression is expression + expression, if expression then expression else expression,

. . .

# Context-free Grammars (CFGs)

- Programming language constructs have recursive structure
- Example: An L expression is expression + expression, if expression then expression else expression, ...
- Context free grammars are a natural notation for this recursive structure
• A CFG consists of:

• A CFG consists of:

• A set of terminals T

• A CFG consists of:

- $\blacktriangleright$  A set of terminals T
- A set of non-terminals N

- A CFG consists of:
  - $\blacktriangleright$  A set of terminals T
  - $\blacktriangleright$  A set of non-terminals N
  - ► A start symbol *S* (non-terminal)

- A CFG consists of:
  - $\blacktriangleright$  A set of terminals T
  - $\blacktriangleright$  A set of non-terminals N
  - ► A start symbol *S* (non-terminal)
  - A set of productions

$$X \to Y_1 Y_2 \dots Y_n$$

where  $X \in N$  and  $Y_i \in (T \cup N \cup \{\varepsilon\})$ 

# Notational Conventions in this Class

Non-terminals are always written upper-case

# Notational Conventions in this Class

- Non-terminals are always written upper-case
- Terminals are written lower-case

# Notational Conventions in this Class

- Non-terminals are always written upper-case
- Terminals are written lower-case
- The start symbol is the left-hand side of the first production

# **CFG Examples**

A fragment of L

# **CFG Examples**

▶ A fragment of L EXPR  $\rightarrow$  if EXPR then EXPR else EXPR | EXPR + EXPR | id

# CFG Examples continued

Simple arithmetic expressions:

# CFG Examples continued

► Simple arithmetic expressions: EXPR → E \* E | E + E | (E) | id

• Recall production rules:  $X \to Y_1 \dots Y_n$ 

- Recall production rules:  $X \to Y_1 \dots Y_n$
- Means that X can be replaced by  $Y_1 \dots Y_N$

- Recall production rules:  $X \to Y_1 \dots Y_n$
- Means that X can be replaced by  $Y_1 \dots Y_N$
- More specifically:

- Recall production rules:  $X \to Y_1 \dots Y_n$
- Means that X can be replaced by  $Y_1 \dots Y_N$
- More specifically:
  - 1. Begin with string consisting of the start symbol "S"

- Recall production rules:  $X \to Y_1 \dots Y_n$
- Means that X can be replaced by  $Y_1 \dots Y_N$
- More specifically:
  - 1. Begin with string consisting of the start symbol "S"
  - 2. Replace any non-terminal X in string with the right-hand side of some production

$$X \to Y_1 \dots Y_n$$

- Recall production rules:  $X \to Y_1 \dots Y_n$
- Means that X can be replaced by  $Y_1 \ldots Y_N$
- More specifically:
  - 1. Begin with string consisting of the start symbol "S"
  - 2. Replace any non-terminal X in string with the right-hand side of some production

$$X \to Y_1 \dots Y_n$$

3. Repeat (2) until there are no non-terminals in the string

# The Language of a CFG continued

More formally, write

$$X_1 \dots X_i \dots X_n \to X_1 \dots X_{i-1} Y_1 \dots Y_m X_{i+1} \dots X_n$$

if there is a production

$$X_i \to Y_1 \dots Y_m$$

### The Language of a CFG continued

More formally, write

$$X_1 \dots X_i \dots X_n \to X_1 \dots X_{i-1} \, Y_1 \dots \, Y_m X_{i+1} \dots X_n$$
 if there is a production

$$X_i \to Y_1 \dots Y_m$$

► Abbreviation: Write  $X_1 
dots X_n \rightarrow^* Y_1 
dots Y_m$  if  $X_1 
dots X_n \rightarrow dots \rightarrow Y_1 
dots Y_m$  in 0 or more steps

## The Language of a CFG continued

▶ Now, let G be a context-free grammar with start symbol S. Then the language of G is:

 $\{a_1 \dots a_n | S \to^* a_1 \dots a_n \text{ and every } a_i \text{ is a terminal}\}$ 

 Terminals are called "terminals" because there are no rules for replacing them

- Terminals are called "terminals" because there are no rules for replacing them
- Once generated, terminals are permanent

- Terminals are called "terminals" because there are no rules for replacing them
- Once generated, terminals are permanent
- Question: What should terminals be when parsing a programming language?

- Terminals are called "terminals" because there are no rules for replacing them
- Once generated, terminals are permanent
- Question: What should terminals be when parsing a programming language?
- Answer: Tokens

• L(G) is the language of CFG G

- L(G) is the language of CFG G
- Strings of balanced parentheses:

 $\{(^i)^j | i \ge 0\}$ 

- L(G) is the language of CFG G
- Strings of balanced parentheses:

 $\{(^i)^j | i \ge 0\}$ 

CFG:

$$\begin{array}{c} S \to (S) \\ S \to \varepsilon \end{array}$$

- L(G) is the language of CFG G
- Strings of balanced parentheses:

 $\{(^i)^j | i \ge 0\}$ 

CFG:

$$\begin{array}{c} S \to (S) \\ S \to \varepsilon \end{array}$$

or equivalently

 $S \to (S) \mid \varepsilon$ 

Recall the earlier fragment of L:

► Recall the earlier fragment of L: EXPR → if EXPR then EXPR else EXPR | EXPR + EXPR | id

► Recall the earlier fragment of L: EXPR → if EXPR then EXPR else EXPR | EXPR + EXPR | id

Some strings in this language:

► Recall the earlier fragment of L: EXPR → if EXPR then EXPR else EXPR | EXPR + EXPR | id

Some strings in this language:

```
► ID
IF ID THEN ID ELSE ID
ID + ID
IF ID THEN ID+ID ELSE ID
IF IF ID THEN ID ELSE IF THEN ID ELSE ID
```

Recall simple arithmetic expressions:

► Recall simple arithmetic expressions: EXPR → E \* E

```
| E + E
| (E)
| id
```

- ► Recall simple arithmetic expressions: EXPR → E \* E | E + E | (E) | id
- Some strings in this language:
### Examples

- ► Recall simple arithmetic expressions: EXPR → E \* E
  - | E + E | (E) | id
- Some strings in this language:

id
(id)
(id)\*id
id+id
id\*id
id\*(id) ...

#### • The idea of a CFG is a big step towards parsing tokens.

- The idea of a CFG is a big step towards parsing tokens.
- But we don't just want to know if a string of tokens is in a language, we also need parse tree of input tokens

- The idea of a CFG is a big step towards parsing tokens.
- But we don't just want to know if a string of tokens is in a language, we also need parse tree of input tokens
- Must also handle errors gracefully

- The idea of a CFG is a big step towards parsing tokens.
- But we don't just want to know if a string of tokens is in a language, we also need parse tree of input tokens
- Must also handle errors gracefully
- Need an implementation of CFGs (e.g., bison)

A derivation is a sequence of productions

 $S \to \ldots \to \ldots \to \ldots$ 

• A derivation is a sequence of productions

 $S \to \ldots \to \ldots \to \ldots$ 

A derivation can be drawn as a tree

A derivation is a sequence of productions

 $S \to \ldots \to \ldots \to \ldots$ 

A derivation can be drawn as a tree

Start symbol is the tree's root

A derivation is a sequence of productions

 $S \to \ldots \to \ldots \to \ldots$ 

A derivation can be drawn as a tree

- Start symbol is the tree's root
- ▶ For a production  $X \to Y_1 \dots Y_n$  add children  $Y_1 \dots Y_n$  to node X

## Derivation Example

- Ε
- $\rightarrow$  E+E
- $\rightarrow$  E\*E+E
- ightarrow id\*E+E
- $\rightarrow$  id\*id + E
- $\rightarrow$  id\*id + id



Е

Е

 $E \rightarrow E+E$ 



- Е
- $\rightarrow$  E+E
- $\rightarrow$  E\*E+E



- Е
- $\rightarrow$  E+E
- $\rightarrow$  E\*E+E
- $\rightarrow$  id\*E+E



- Е
- $\rightarrow$  E+E
- $\rightarrow$  E\*E+E
- ightarrow id\*E+E
- $\rightarrow$  id\*id + E



- Ε
- $\rightarrow$  E+E
- $\rightarrow$  E\*E+E
- ightarrow id\*E+E
- $\rightarrow$  id\*id + E
- $\rightarrow$  id\*id + id



A parse tree has terminals at the leaves and non-terminals at the interior nodes

- A parse tree has terminals at the leaves and non-terminals at the interior nodes
- > An in-order traversal of the leaves is the original input

- A parse tree has terminals at the leaves and non-terminals at the interior nodes
- > An in-order traversal of the leaves is the original input
- The parse tree shows the associativity of operations, the input token string does not

- A parse tree has terminals at the leaves and non-terminals at the interior nodes
- An in-order traversal of the leaves is the original input
- The parse tree shows the associativity of operations, the input token string does not
- Example: The parse tree from the last slide encodes that times has higher precedence than plus

Left-most and Right-most Derivations

The example we looked at is a left-most derivation

### Left-most and Right-most Derivations

- The example we looked at is a left-most derivation
- This means: At each step, we replace the left-most non-terminal

### Left-most and Right-most Derivations

- The example we looked at is a left-most derivation
- This means: At each step, we replace the left-most non-terminal
- There is also an equivalent notion of right-most derivation

Е

Е

 $E \rightarrow E+E$ 



 $E \rightarrow E+E$ 

 $\rightarrow$  E+id



- Ε
- $\rightarrow$  E+E
- $\rightarrow$  E+id
- $\rightarrow$  E\*E+id



- Е
- $\rightarrow$  E+E
- $\rightarrow$  E+id
- $\rightarrow$  E\*E+id
- $\rightarrow$  E\*id + id



- Ε
- $\rightarrow$  E+E
- $\rightarrow$  E+id
- $\rightarrow$  E\*E+id
- $\rightarrow$  E\*id + id
- $\rightarrow$  id\*id + id



 Observe that left-most and right-most derivations have the same parse tree

- Observe that left-most and right-most derivations have the same parse tree
- > The only difference is the order in which branches are added

- Observe that left-most and right-most derivations have the same parse tree
- ► The only difference is the order in which branches are added
- But when parsing tokens, we only care about the final parse tree, which may have many different derivations

- Observe that left-most and right-most derivations have the same parse tree
- ► The only difference is the order in which branches are added
- But when parsing tokens, we only care about the final parse tree, which may have many different derivations
- Left-most and right-most derivations are important in parser implementations

## Ambiguity

Recall our example grammar:

# Ambiguity

► Recall our example grammar: EXPR → E \* E | E + E | (E) | id

## Ambiguity

#### ► Recall our example grammar: EXPR → E \* E | E + E | (E) | id

Now, consider the string id\*id+id
## Ambiguity continued

This string has two parse trees!

# Ambiguity continued

This string has two parse trees!



# Ambiguity continued

This string has two parse trees!



 A grammar is ambiguous if it has more than one parse tree for some string

- A grammar is ambiguous if it has more than one parse tree for some string
- Equivalently: There is more than one left-most or right-most derivation for some string

- A grammar is ambiguous if it has more than one parse tree for some string
- Equivalently: There is more than one left-most or right-most derivation for some string
- Ambiguity is bad!

- A grammar is ambiguous if it has more than one parse tree for some string
- Equivalently: There is more than one left-most or right-most derivation for some string
- Ambiguity is bad!
- Leaves meaning of programs ill-defined

First method: Rewrite grammar unambiguously

First method: Rewrite grammar unambiguously

Question: How can we write simple arithmetic expressions unambiguously?

First method: Rewrite grammar unambiguously

- Question: How can we write simple arithmetic expressions unambiguously?
- Solution: Enforce precedence of times over plus by generating all pluses fist:

First method: Rewrite grammar unambiguously

- Question: How can we write simple arithmetic expressions unambiguously?
- Solution: Enforce precedence of times over plus by generating all pluses fist:

$$\begin{array}{rrrr} S & \rightarrow & E+S \mid E \\ E & \rightarrow & \mathsf{id} \ast E \mid \mathsf{id} \mid (S) \ast E \mid (S) \end{array}$$

 However, converting grammars to unambiguous form can be very difficult

- However, converting grammars to unambiguous form can be very difficult
- It also often results in horrible, unintuitive grammars with many non-terminals

- However, converting grammars to unambiguous form can be very difficult
- It also often results in horrible, unintuitive grammars with many non-terminals
- It is also fundamentally impossible to transform an ambiguous grammar into a unambiguous grammar

- However, converting grammars to unambiguous form can be very difficult
- It also often results in horrible, unintuitive grammars with many non-terminals
- It is also fundamentally impossible to transform an ambiguous grammar into a unambiguous grammar
- For this reason, tools such as bison include disambiguation mechanisms

Instead of rewriting the grammar:

#### Instead of rewriting the grammar:

Use the more natural ambiguous grammar

#### Instead of rewriting the grammar:

- Use the more natural ambiguous grammar
- Along with disambiguating declarations

- Instead of rewriting the grammar:
  - Use the more natural ambiguous grammar
  - Along with disambiguating declarations
- The parser tool bison allows you to declare precedence and associativity for this

• Consider the grammar  $E \rightarrow E + E \mid id$ 

- Consider the grammar  $E \to E + E \mid \text{id}$
- Ambiguous: Two parse trees of input id + id + id

- Consider the grammar  $E \rightarrow E + E \mid id$
- Ambiguous: Two parse trees of input id + id + id



- Consider the grammar  $E \to E + E \mid \text{id}$
- Ambiguous: Two parse trees of input id + id + id



- Consider the grammar  $E \rightarrow E + E \mid id$
- Ambiguous: Two parse trees of input id + id + id



Declare left associativity of plus as: %left +

- Consider the grammar  $E \rightarrow E + E \mid id$
- Ambiguous: Two parse trees of input id + id + id



Declare left associativity of plus as: %left +

▶ Consider the grammar  $E \to E + E \mid \text{id}$  and input id + id \* id

▶ Consider the grammar  $E \rightarrow E + E \mid \text{id and input}$ id + id \* id



▶ Consider the grammar  $E \rightarrow E + E \mid \text{id and input}$ id + id \* id

E

id

id



▶ Consider the grammar  $E \rightarrow E + E \mid \text{id and input}$ id + id \* id





- Precedence Declaration:
  - %left +
  - %left \*

► Consider the grammar E → E + E | id and input id + id \* id





- Precedence Declaration:
  - %left +
  - %left \*

 We have seem how to specify programming language syntax with CFGs

- We have seem how to specify programming language syntax with CFGs
- We built parse trees that express the high-level syntactic structure

- We have seem how to specify programming language syntax with CFGs
- We built parse trees that express the high-level syntactic structure
- Parse trees of programs are known as abstract syntax trees

- We have seem how to specify programming language syntax with CFGs
- We built parse trees that express the high-level syntactic structure
- Parse trees of programs are known as abstract syntax trees
- We discussed ambiguity of CFGs