# CS 412
## Introduction to Compilers

Ross Tate
Cornell University
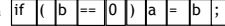
Lecture 4: Syntactic Analysis

---

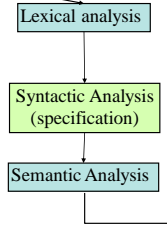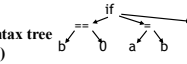# Where we are

Source code
(character stream)          if (b == 0) a = b;

                                              Lexical analysis

Token stream    if  (  b  ==  0  )  a  =  b  ;

                                              Syntactic Analysis
                                              (specification)

Abstract syntax tree
(AST)

                                              Semantic Analysis
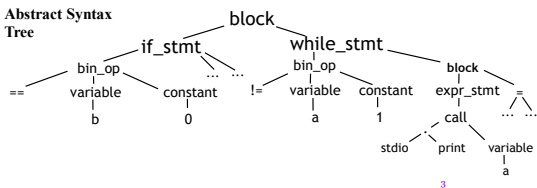
2

---

# What is Syntactic Analysis?

Source code
(token stream)
```
{
    if (b == (0)) a = b;
    while (a != 1) {
        stdio.print(a);
        a = a - 1;
    }
}
```
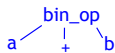
Abstract Syntax Tree

block
├── if_stmt
│   ├── bin_op ... ...
│   │   ├── ==
│   │   ├── variable — b
│   │   └── constant — 0
└── while_stmt
    ├── bin_op
    │   ├── !=
    │   ├── variable — a
    │   └── constant — 1
    └── block
        ├── expr_stmt
        │   └── call
        │       ├── stdio
        │       ├── print
        │       └── variable — a
        └── = ... ...

3

---

# Parsing

- Parsing: recognizing whether a program (or sentence) is grammatically well-formed & identifying the function of each component.

"I gave him the book"

sentence
├── subject: I
├── verb: gave
├── indirect object: him
└── object
    └── noun phrase
        ├── article: the
        └── noun: book

4

---

# Overview of Syntactic Analysis

- Input: stream of tokens

- Output: abstract syntax tree
  – Abstract syntax tree removes extra syntax
  $a + b \approx (a) + (b) \approx ((a)+((b)))$

bin_op
a   +   b

5

---

# What Parsing doesn't do

- Doesn't check many things: type agreement, variables declared, variables initialized, etc.
  – int x = true;
  – int y; z = f(y);

- Deferred until semantic analysis

6

## Specifying Language Syntax

- First problem: how to describe language syntax precisely and conveniently
- Last time: can describe tokens using regular expressions
- Regular expressions easy to implement, efficient (by converting to DFA)
- Why not use regular expressions (on tokens) to specify programming language syntax?
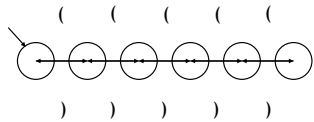
7

## Limits of REs

- Programming languages are not regular -- cannot be described by regular expressions

- Consider: language of all strings that contain balanced parentheses (easier than PLs)
  - () (()) ()() (())()((()()))
  - (( )( ()) ()()

- Problem: need to keep track of number of parentheses seen so far: unbounded counting

8

## Limits of REs

- RE = DFA
- DFA has only finite number of states; cannot perform unbounded counting

( ( ( ( (

*maximum depth: 5 parens*

9

Regexes

Compiler Writer

Scotty, we need more power!

## Context-Free Grammars

- A specification of the balanced-parenthesis language:
  - $S \rightarrow ( S ) S$
  - $S \rightarrow \varepsilon$
- The definition is recursive
- A **context-free grammar**
  - More expressive than regular expressions
  - $S = (S) \varepsilon = ((S) S) \varepsilon = ((\varepsilon) \varepsilon) \varepsilon = (())$
- If a grammar accepts a string, there is a *derivation* of that string using the productions of the grammar

11

## Definition of CFG

- Terminals
  - Token or $\varepsilon$
- Non-terminals $\qquad S \rightarrow ( S ) S$
  - Syntactic variables $\qquad S \rightarrow \varepsilon$
- Start symbol
  - A special nonterminal is designated: $S$
- Productions
  - Specify how non-terminals may be expanded to form strings
  - LHS: single non-terminal, RHS: string of terminals or non-terminals
- Vertical bar is shorthand for multiple prod'ns

12

## RE is subset of CFG

- Regular Expression for real numbers:
  
  $digit \rightarrow [\mathbf{0\text{-}9}]$
  $posint \rightarrow digit+$
  $int \rightarrow -?\ posint$
  $real \rightarrow int\ .\ (\varepsilon \mid posint)$

- RE symbolic names are only *shorthand:* no recursion, so all symbols can be fully expanded:
  
  $real \rightarrow -?\ [\mathbf{0\text{-}9}]+ . (\varepsilon \mid ([\mathbf{0\text{-}9}]+))$

## Sum grammar

$$S \rightarrow E + S \mid E$$
$$E \rightarrow \mathbf{number} \mid (S)$$

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

$S \rightarrow E + S$
$S \rightarrow E$
$E \rightarrow number$
$E \rightarrow (S)$

4 productions
2 non-terminals: S, E
4 terminals: (, ), +, number
start symbol S

14

## Derivation Example

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

Derive $(1+2+(3+4))+5$:

$S \rightarrow E + S \rightarrow (S) + S \rightarrow (E + S) + S$
$\rightarrow (1 + S)+S \rightarrow (1 + E + S)+S$
$\rightarrow (1 + 2 + S)+S \rightarrow (1 + 2 + E)+S$
$\rightarrow (1 + 2 + (S))+S \rightarrow (1 + 2 + (E + S))+S$
$\rightarrow (1 + 2 + (3 + S))+S$
$\rightarrow (1 + 2 + (3 + E))+S$
$\rightarrow (1 + 2+ (3+4))+S$
$\rightarrow (1 + 2+ (3+4))+E$
$\rightarrow (1 + 2+ (3+4))+5$

replacement string

non-terminal being expanded

15

## Constructing a derivation

- Start from start symbol: $S$
- Productions are used to derive a sequence of tokens from the start symbol
- For arbitrary strings α, β and γ and a production A → β
  A single step of derivation is
  
  $\alpha A \gamma \implies \alpha \beta \gamma$
  
  – *i.e.*, substitute β for an occurrence of A
  – $(S + E) + E \rightarrow (E + S + E)+E$
  
  $(A = S, \beta = E + S)$

16

## Derivation Example

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

Derive $(1+2+(3+4))+5$:

$S \rightarrow E + S \rightarrow (S) + S \rightarrow (E + S) + S$
$\rightarrow (1 + S)+S \rightarrow (1 + E + S)+S$
$\rightarrow (1 + 2 + S)+S \rightarrow (1 + 2 + E)+S$
$\rightarrow (1 + 2 + (S))+S \rightarrow (1 + 2 + (E + S))+S$
$\rightarrow (1 + 2 + (3 + S))+S$
$\rightarrow (1 + 2 + (3 + E))+S$
$\rightarrow (1 + 2+ (3+4))+S$
$\rightarrow (1 + 2+ (3+4))+E$
$\rightarrow (1 + 2+ (3+4))+5$

replacement string

non-terminal being expanded

17

## Derivation ⟹ Parse Tree

Parse Tree

Derivation

- Tree representation of the derivation
- Leaves of tree are terminals; in-order traversal yields string
- Internal nodes: non-terminals
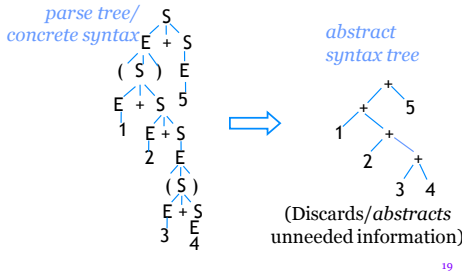- No information about order of derivation steps

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

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

$S \rightarrow E + S \rightarrow (S) + S \rightarrow (E + S) + S \rightarrow (1 + S)+S \rightarrow (1 + E + S)+S$
$\rightarrow (1 + 2 + S)+S \rightarrow ... \rightarrow (1 + 2 + (S))+S \rightarrow (1 + 2 + (E + S))+S$
$\rightarrow ... \rightarrow (1 + 2 + (3 + E))+S \rightarrow ... \rightarrow (1 + 2+ (3+4))+5$

18

3

## Parse Tree

- Also called "concrete syntax"

*parse tree/ concrete syntax*

```
        S
      E   S
    ( S ) E
  E + S   5
  1 E + S
    2 E
      ( S )
      E + S
      3 E
        4
```

*abstract syntax tree*

```
        +
      +   5
    1   +
      2   +
        3   4
```

(Discards/*abstracts* unneeded information)

19

## Derivation order

- Can choose to apply productions in any order; select any non-terminal

  $E + S \to 1 + S$ or $E + E + S$

- Two standard orders: left- and right-most -- useful for different kinds of automatic parsing

- **Leftmost derivation**: In the string, find the left-most non-terminal and apply a production to it. $E + S \to 1 + S$

- **Rightmost derivation**: find right-most non-terminal...etc. $E + S \to E + E + S$

20

## Example

$S \to E + S \,|\, E$
$E \to number \,|\, (\,S\,)$

- Left-most derivation
$S \to E+S \to (S) + S \to (E + S)+ S \to (1 + S)+S \to (1+E+S)+S \to (1+2+S)+S \to (1+2+E)+S \to (1+2+(S))+S \to (1+2+(E+S))+S \to (1+2+(3+S))+S \to (1+2+(3+E))+S \to (1+2+(3+4))+S \to (1+2+(3+4))+E \to (1+2+(3+4))+5$

- Right-most derivation
$S \to E+S \to E+E \to E+5 \to (S)+5 \to (E+S)+5 \to (E+E+S)+5 \to (E+E+E)+5 \to (E+E+(S))+5 \to (E+E+(E+S))+5 \to (E+E+(E+E))+5 \to (E+E+(E+4))+5 \to (E+E+(3+4))+5 \to (E+2+(3+4))+5 \to (1+2+(3+4))+5$

- Same parse tree: same productions chosen, diff. order

21

## Associativity

- + operator associates to right in parse tree regardless of derivation order

$(1+2+(3+4))+5$ ⇒

```
      +
    +   5
  1   +
    2   +
      3   4
```

- + associates to right because of ***right-recursive*** production $S \to E + S$

- In the example grammar, leftmost and rightmost derivations produce identical parse trees

22

## An Ambiguous Grammar

- Consider another grammar:

  $S \to S + S \,|\, S * S \,|\, number$
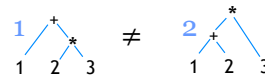
- Different derivations produce different parse trees: ambiguous grammar

23

## Differing Parse Trees

$S \to S + S \,|\, S * S \,|\, number$

- Consider expression $1 + 2 * 3$
- Derivation 1: $S \to S + S \to 1 + S \to 1 + S * S \to 1 + 2 * S \to 1 + 2 * 3$
- Derivation 2: $S \to S * S \to S * 3 \to S + S * 3 \to S + 2 * 3 \to 1 + 2 * 3$

```
1   +              2   *
  1   *      ≠       +   3
    2   3         1   2
```

1 2 3        1 2 3

24

## Impact of Ambiguity

- Different parse trees correspond to different evaluations!



$$\underset{1 \;\; 2 \;\; 3}{\overset{+}{\diagup}\overset{*}{\diagdown}} = 7 \qquad \underset{1 \;\; 2 \;\; 3}{\overset{*}{\diagup}\overset{+}{\diagdown}} = 9$$
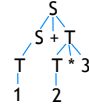
- Meaning of program not well defined

25

## Eliminating Ambiguity

- Often can eliminate ambiguity by adding non-terminals & allowing recursion only on right or left

$$S \rightarrow S + T \mid T$$
$$T \rightarrow T * \textbf{num} \mid \textbf{num}$$



- *S/T* separation enforces precedence
- Left-recursion : left-associativity

26

## if-then-else

- How to write a grammar for if stmts?

$$S \rightarrow \text{if } (E) \; S \text{ else } S$$
$$S \rightarrow \text{if } (E) \; S$$
$$S \rightarrow X = E \mid ...$$

- Is this grammar ok?

27

## No—Ambiguous!

- How to parse?
  if $(E_1)$ if $(E_2)$ $S_1$ else $S_2$

$$S \rightarrow \text{if } (E) \; S$$
$$S \rightarrow \text{if } (E) \; S \text{ else } S$$
$$S \rightarrow other$$
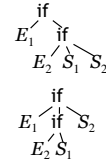
$S \;\; \rightarrow \text{if } (E) \; S$
$\rightarrow \text{if } (E) \; \underline{\text{if } (E) \; S \text{ else } S}$

$S \;\; \rightarrow \text{if } (E) \; S \text{ else } S$
$\rightarrow \text{if } (E) \; \underline{\text{if } (E) \; S} \text{ else } S$



- *Which "if" is the "else" attached to?*

28

## Grammar for Closest-if Rule

- Want to rule out if $(E)$ if $(E)$ $S$ else $S$
- Problem: unmatched if may not occur as the "then" (consequent) clause of a containing "if"

$$\begin{aligned} statement \quad &\rightarrow matched \mid unmatched \\ matched \quad &\rightarrow \text{if } (E) \; matched \text{ else } matched \\ &\qquad \mid \; other \\ unmatched \quad &\rightarrow \text{if } (E) \; statement \\ &\qquad \mid \; \text{if } (E) \; matched \text{ else } unmatched \end{aligned}$$

29

## Greedy ANTLR

- How to parse?
  if $(E_1)$ if $(E_2)$ $S_1$ else $S_2$

$$S \rightarrow \text{if } (E) \; S$$
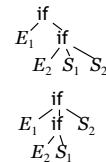$$S \rightarrow \text{if } (E) \; S \text{ else } S$$
$$S \rightarrow other$$

$S \;\; \rightarrow \text{if } (E) \; S$
$\rightarrow \text{if } (E) \; \underline{\text{if } (E) \; S \text{ else } S}$

$S \;\; \rightarrow \text{if } (E) \; S \text{ else } S$
$\rightarrow \text{if } (E) \; \underline{\text{if } (E) \; S} \text{ else } S$



- *Which "if" is the "else" attached to?*

30

## Greedy ANTLR

- ANTLR v4 grammar for if stmts:
    $S \rightarrow$ if $(E)$ $S$ (else $S$)?
    $S \rightarrow X = E \mid \ldots$

- Leftmost derivations
- Greedy derivations

31

## Limits of CFGs

- Syntactic analysis can't catch all "syntactic" errors
- Example: C++
  - HashTable<Key,Value> x;
- Need to know whether HashTable is the name of a type to understand syntax!
  Problem: "<", "," are overloaded
- Iota:
  - f(4)[1][2] = 0;
- Difficult to write grammar for LHS of assign
  - may be easier to allow all exprs, check later

32

## CFGs

- Context-free grammars allow concise specification of programming languages

- CFG specifies how to convert token stream to parse tree
  - If unambiguous
  - Or a derivation preference is designated

- Next time: implementing a top-down parser (leftmost derivation)

33