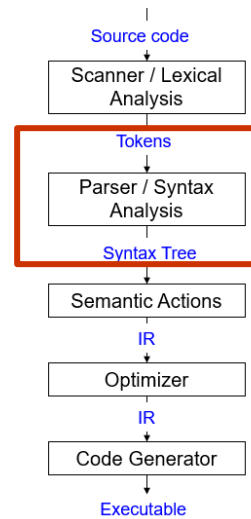# CS406: Compilers
## Spring 2021
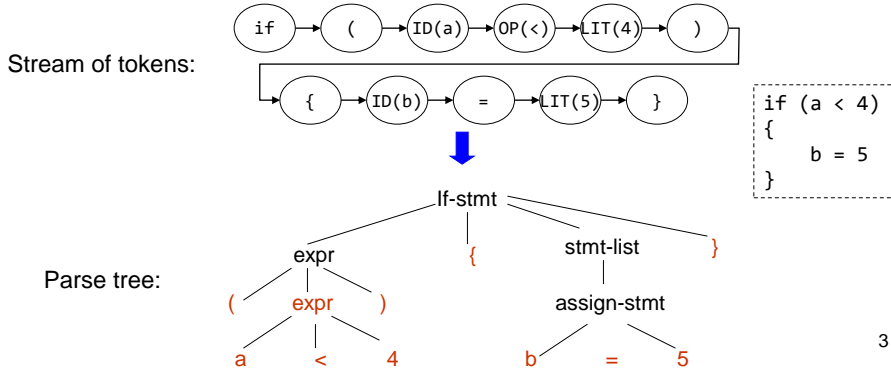
Week 3: Parsers

1

# Parsers - Overview

- Also called syntax analyzers
- Determine two things:
  1. Is a program syntactically valid?
     - Is an English sentence grammatically correct?
  2. What is the structure of programming language constructs? E.g. does the sequence `IF, ID(a), OP(<), ID(b), {, ID(a), ASSIGN, LIT(5), }, ;, }` refer to an `if-statement`?
     - Diagramming English sentences

Source code
↓
Scanner / Lexical Analysis
↓
Tokens
↓
Parser / Syntax Analysis
↓
Syntax Tree
↓
Semantic Actions
↓
IR
↓
Optimizer
↓
IR
↓
Code Generator
↓
Executable

2

# Parsers - Overview

- Input: stream of tokens
- Output: Parse tree
  - sometimes implicit

Stream of tokens:



```
if (a < 4)
{
      b = 5
}
```

Parse tree:



3

# Parsers – what do we need to know?

1. How do we define language constructs?

   – Context-free grammars

2. How do we determine: 1) valid strings in the language? 2) structure of program?

   – LL Parsers, LR Parsers

3. How do we write Parsers?

   – E.g. use a parser generator tool such as Bison

4

# Center Embeddings in English

```
The bird flew

The bird the boy saw flew

The bird the boy the dog chased saw flew

The bird the boy the dog the man owned chased saw flew

The bird the boy the dog the man the woman loved owned chased
     saw flew

...
```

*Exercise: write a regular expression that match the pattern. Note: the alphabets of your language are 'Noun', 'Verb' and 'the'*

5

You can construct arbitrarily long sentences like this in English.

# Languages

- A language is (possibly infinite) set of strings

- Regular expressions describe *regular languages*
  weakness: can't describe a string of the form:

$$\{ \ (^i \ )^i \ | \ i>=1\}$$   E.g. ((2+3)*5)

Parenthesized expressions:   ((( int x; )))
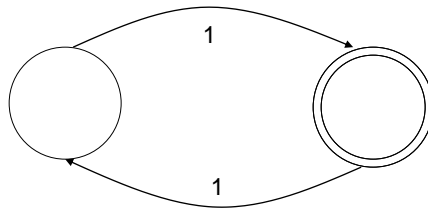
*Programming language syntax is i.e. recursive*

Nested structures:

```
        IF
          IF
            IF
            FI
          FI
        FI
```

is C regular?

6

# Trivia

- Regular expressions can describe strings:
  `{ mod k | k = # states in FA}`



"accept all strings having odd number of 1s"

What FAs and regular expressions can do is describe strings of the form "odd number of 1s", they can determine parity but cannot count.

# Context Free Grammar (CFG)

- Natural notation for describing <u>recursive structure</u> definitions. Hence, suitable for specifying language constructs.

- Consist of:

  - A set of *Terminals* (T)

  - A set of *Non-terminals* (N)

  - A *Start Symbol* (S$\in$N)
    ( aka. rules)
  - A *set of Productions* (X -> $Y_1 .. Y_N$)
    $P:X \longrightarrow Y_1 Y_2 Y_3 .. Y_N \mid X \in N, Y_i \in N \cup T \cup \epsilon/\lambda$

8

# Context Free Grammar (CFG)

- Grammar `G = (T, N, S, P)`

  `E.g. G = ({a,b}, {S, A, B}, S, {S→AB, A→Aa`
  `A→a, B→Bb, B→b})`

- Implicit meanings

  – <u>First rule</u> listed in the set of productions contains <u>start symbol</u> (on the left-hand side)

  – In the set of productions, <u>you can replace the symbol X</u> (appearing on the right-hand side only) with the <u>string of symbols</u> that are on the right-hand side of a rule, which has X (on the left-hand side)

9

# Context Free Grammar (CFG)

1. Begin with only S as the initial string

2. Replace S

   – S replaced with AB

3. Repeat 2 until the string contains only terminals

   – AB replaced with aB

   – aB replaced with bb

```
G = (T, N, S, P)
P:{ S->AB,
    A->Aa,
    A->a,
    B->Bb,
    B->b }
```

**Summary:** we move from S to a string of terminals through a series of <u>transformations:</u>

$$\alpha_0 \text{-> } ... \text{-> } \alpha_n \text{ where } \alpha_1 \ldots \alpha_n \text{ are strings}$$

Shorthand notation: $\alpha_0 \overset{*}{\text{->}} \alpha_n$

# Language of the Grammar

- Language L(G) of the context-free grammar G

  - Set of strings that can be derived from S

  - $\{a_1 a_2 a_3 .. a_N \mid a_i \in T \; \forall i \; \text{and} \; S \overset{*}{\text{->}} a_1 a_2 a_3 .. a_N \}$

  - Is called context-free language

    - All regular languages are context-free but not vice-versa.

    - Can have many grammars generating same language.

11

# Context-Sensitive Grammar

- Can have context-sensitive grammar and languages (think: aB->ab)

    - Cannot replace right-hand side with left-hand side irrespective of the context.

    - E.g. aB->ab lays down a context: 'a' must be a prefix in order to transform the string "aB" to a string of terminals "ab"

        - ccaBb can be replaced by ccabb

```
                                 G = (T, N, S, P)
                                 P:{ S->AB,
        Is grammar G context-free?    A->Aa,
                                      A->a,
                                      B->Bb,
                                      B->b }
```

12

# Does a string belong to the Language?

- How do we apply the grammar rules to determine the validity of a string? (i.e. string belongs to the language specified by the context-free grammar)

  - Begin with S

  - Replace S

  - Repeat till string contains terminals only
    *L(G) must contain strings of terminals only*

- Notation:

  - We will use Greek letters to denote strings containing non-terminals and terminals

13

# Simple grammar

$$A \rightarrow Aa$$
$$\mid a$$

Start symbol ——→ S → A B

A → A a ←—— Terminals

Non-terminals A → a

B → B b

B → b ←—— Production

*Backus Naur Form (BNF)*

14

# Generating strings

S → A B

A → A a

A → a

B → B b

B → b

- Given a start rule, productions tell us how to rewrite a non-terminal into a different set of symbols

- Some productions may rewrite to $\lambda$. That just removes the non-terminal

To derive the string "a a b b b" we can do the following rewrites:

S ⇒ A B ⇒ A a B ⇒ a a B ⇒ a a B b ⇒

a a B b b ⇒ a a b b b

15

# Exercise

Which of the below strings are accepted by the grammar:

```
1: A -> aAa
2: A -> bBb
3: A -> λ
4: B -> cA
5: B -> λ
```

1. abcba      1->2->4->3
2. abcbca
3. abba      1->2->5
4. abca

16

# Programming language syntax

- Programming language syntax is defined with CFGs

- Constructs in language become non-terminals

  - May use auxiliary non-terminals to make it easier to define constructs

    if_stmt  → if ( cond_expr ) then statement else_part

    else_part  → else statement

    else_part  → λ

- Tokens in language become terminals

17

# CFG Contd..

- Is it enough if parsers answer "yes" or "no" to check if a string belongs to context-free language?

  - Also need a parse tree

- What if the answer is a "no"?

  - Handle errors

- How do we implement CFGs?

  - E.g. Bison

18

# Parse trees

- Tree which shows how a string was produced by a language

  - Interior nodes of tree: non-terminals

    - Children: the terminals and non-terminals generated by applying a production rule

  - Leaf nodes: terminals



19

# Parse Trees and String Derivations

- Recall: sequence of rules applied to produce a string is a derivation

- A derivation defines a parse tree

  - A parse tree may have many derivations

20

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id\*id+id**

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id*id+id**

**Apply 1: Start with E, the start symbol      Parse Tree**

E

(E)

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id\*id+id**

**Apply 1: Replace E with E + E**         **Parse Tree**

```
E
E+E
```



23

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id*id+id**

**Apply 2: Replace E with E \* E**                    **Parse Tree**

```
E
E+E
E*E+E
```

24

# Derivations and Parse Trees

- Consider the grammar with the following rules:
  ```
  1: E -> E + E
  2:    | E * E
  3:    | id
  ```
- Produce derivations for the string: **id*id+id**

**Apply 3: Replace E with id**                **Parse Tree**

```
E
E+E
E*E+E
id*E+E
```



25

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id\*id+id**

**Apply 3: Replace E with id**                    **Parse Tree**

```
E
E+E
E*E+E
id*E+E
id*id+E
```



26

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id*id+id**

**Apply 3: Replace E with id**          **Parse Tree**

```
E
E+E
E*E+E
id*E+E
id*id+E
id*id+id
```

27

# Derivations and Parse Trees

- Note in previous slides:
  - Replacement done on left-most non-terminal in the string - **called left-most derivation**
  - Terminals at leaves and non-terminal as interior nodes
  - Inorder traversal produces input string id*id+id

28

# Derivations and Parse Trees

- Note in previous slides:
  - Replacement done on left-most non-terminal in the string - **called left-most derivation**
  - Terminals at leaves and non-terminal as interior nodes
  - Inorder traversal produces input string id*id+id
  - Parse tree shows association of operations. Input string doesn't
    - * associated with identifiers in the subtree

$$(id * id)+id$$

29

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id*id+id**
  - Using right-most derivations

    i.e. replace the right-most non-terminal

30

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id\*id+id**

**Start with E, the start symbol**

E

( E )  **Parse Tree**

# Derivations and Parse Trees

- Consider the grammar with the following rules:

  ```
  1: E -> E + E
  2:    | E * E
  3:    | id
  ```

- Produce derivations for the string: **id*id+id**

**Apply 2: Replace E with E+E**

```
E
E+E
```



**Parse Tree**

32

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id\*id+id**

**Apply 1: Replace E with id**

```
E
E+E
E+id
```



Parse Tree

33

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id*id+id**

**Apply 3: Replace E with E * E**

```
E
E+E
E+id
E*E+id
```

Parse Tree



34

# Derivations and Parse Trees

- Consider the grammar with the following rules:
  ```
  1: E -> E + E
  2:    | E * E
  3:    | id
  ```
- Produce derivations for the string: **id*id+id**

**Apply 3: Replace E with id**

```
E
E+E
E+id
E*E+id
E*id+id
```

Parse Tree

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id\*id+id**

**Apply 3: Replace E with id**

```
E
E+E
E+id
E*E+id
E*id+id
id*id+id
```
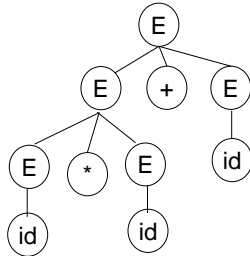


Parse Tree

36

# Derivations and Parse Trees

- We get the same parse tree using left-most and right-most derivations.
  - Every parse tree has left-most and right-most (and any random order) derivations.



37

# Derivations and Parse Trees

- We get the same parse tree using left-most and right-most derivations.
  - Every parse tree has left-most and right-most (and any random order) derivations.



- But there could be a string (or more than one strings) for which there exists derivations that would get different parse trees

38

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id*id+id**

**Start with E, the start symbol**

E                                    ( E )   **Parse Tree**

39

# Derivations and Parse Trees

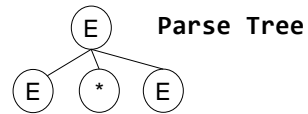- Consider the grammar with the following rules:
  ```
  1: E -> E + E
  2:    | E * E
  3:    | id
  ```
- Produce derivations for the string: **id*id+id**

**Apply 2: Replace E with E*E** *Earlier it was replace E with E+E*

```
E
E*E
```

**Parse Tree**

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id*id+id**

**Apply 1: Replace E with E+E**

```
E
E*E
E*E+E
```

**Parse Tree**



41

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

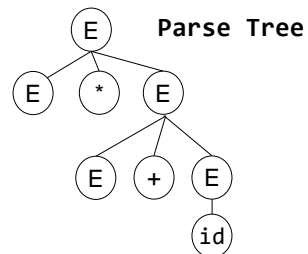- Produce derivations for the string: **id*id+id**

**Apply 3: Replace E with id**

```
E
E*E
E*E+E
E*E+id
```



Parse Tree

42

# Derivations and Parse Trees

- Consider the grammar with the following rules:
  ```
  1: E -> E + E
  2:    | E * E
  3:    | id
  ```
- Produce derivations for the string: **id*id+id**

**Apply 3: Replace E with id**

```
E
E*E
E*E+E
E*E+id
E*id+id
```

**Parse Tree**



43

# Derivations and Parse Trees

- Consider the grammar with the following rules:

```
1: E -> E + E
2:    | E * E
3:    | id
```

- Produce derivations for the string: **id\*id+id**

**Apply 3: Replace E with id**

```
E
E*E
E*E+E
E*E+id
E*id+id
id*id+id
```
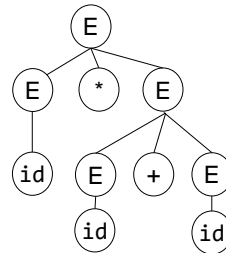
**Parse Tree**

44

# Derivations and Parse Trees

- Input string: `id*id+id`



**earlier**

**now**

- `Inorder` traversal of both trees produces the same string

45

# Ambiguous Grammar

- Grammar that produces more than one parse tree for <u>some</u> string

```
1: E -> E + E
2:    | E * E
3:    | id
```
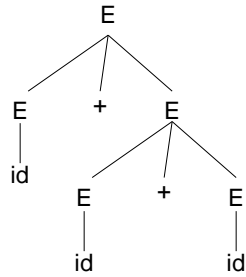
46

# Ambiguity – what to do?

- Ignore it (let it be ambiguous)
  - Give hints to other components of the compiler on how to resolve it
- Fix it
  - Manually
  - May make the grammar complicated and difficult to maintain
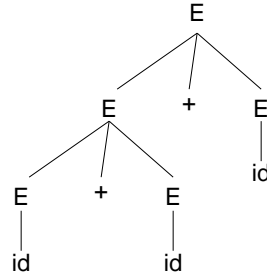
47

# Ambiguity – ignore

- E -> E + E | id

```
E->E+E
E->id+E
E->id+E+E
E->id+id+E
E->id+id+id

Produces:
id+(id+id)
```

```
E->E+E
E->E+E+E
E->id+E+E
E->id+id+E
E->id+id+id

Produces:
(id+id)+id
```
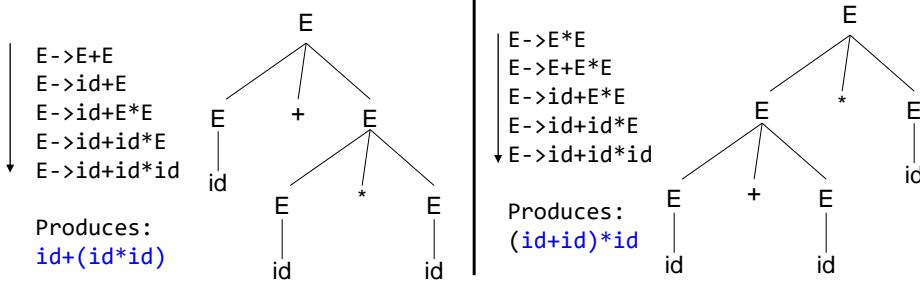
- Associativity declaration in Bison:
  %left +          Picks the parse tree on the right

# Ambiguity - ignore

- E -> E + E | E * E | **id**



```
E->E+E
E->id+E
E->id+E*E
E->id+id*E
E->id+id*id
```

Produces:
id+(id*id)

```
E->E*E
E->E+E*E
E->id+E*E
E->id+id*E
E->id+id*id
```
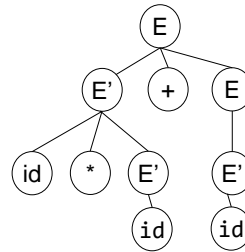
Produces:
(id+id)*id

**%left +**
**%left ***

*Tells that \* has higher precedence over + and both are left associative. So, we get the tree on left.*

49

# Ambiguity – fixing

- Rewrite `E -> E + E`   as:   `E -> E' + E | E'`
                 `| E * E`          `E' -> id * E' | id`
                 `| id`             `| (E) * E' | (E)`

```
E->E'+E
E'->id*E'
E'->id
E->E'
E'->id
```

*Is the above sequence left-most or right-most derivation?*

E controls generation of +

E' controls generation of *. *'s are nested deeper in the parse tree.

50

# Ambiguity Fixing - Exercise

**Exercise:** *Is this grammar ambiguous? Draw parse trees for the following* **String: if E1 then if E2 then S1 else S2**

```
1: STMT -> if EXPR then STMT
2:       |  if EXPR then STMT else STMT
3:       |  s1
4:       |  s2
5: EXPR -> e1 | e2
```

# Ambiguity Fixing - Exercise

**Exercise:** *Is this grammar ambiguous? Draw parse trees for the following* **String: if E1 then if E2 then S1 else S2**

```
1: STMT -> if EXPR then STMT
2:       | if EXPR then STMT else STMT
3:       | s1
4:       | s2
5: EXPR -> e1 | e2
```

# Ambiguity Fixing - Exercise

**Exercise:** *Is this grammar ambiguous? Draw parse trees for the following* String: `if E1 then if E2 then S1 else S2`

```
1: STMT -> if EXPR then STMT
2:       |  if EXPR then STMT else STMT
3:       |  s1
4:       |  s2
5: EXPR -> e1 | e2
```
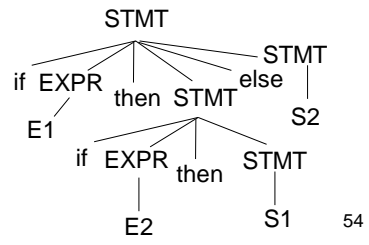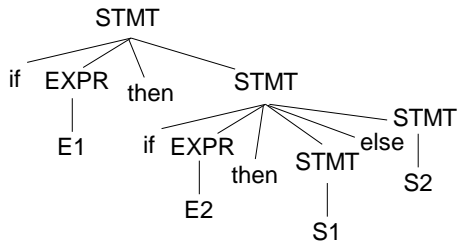
# Ambiguity Fixing - Exercise

**Exercise:** *Is this grammar ambiguous? Draw parse trees for the following*   **String: if E1 then if E2 then S1 else S2**

```
1: STMT -> if EXPR then STMT
2:       | if EXPR then STMT else STMT
3:       | s1
4:       | s2
5: EXPR -> e1 | e2
```



54

# Ambiguity Fixing - Exercise

**Exercise:** Which `if` is the `else` associated with?

**String:** `if E1 then if E2 then S1 else S2`

# Ambiguity Fixing - Exercise

**Exercise:** Which `if` is the `else` associated with?

**String: if E1 then if E2 then S1 else S2**

# Ambiguity Fixing - Exercise

**Exercise:**  Rewrite the grammar to make it unambiguous.

```
1: STMT -> if EXPR then STMT
2:       |  if EXPR then STMT else STMT
3:       |  s1
4:       |  s2
5: EXPR -> e1 | e2
```
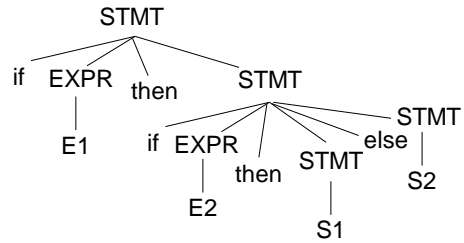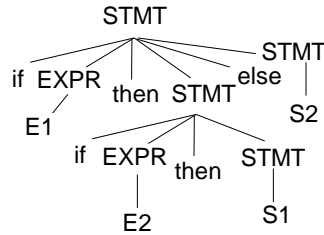
# Ambiguity Fixing - Exercise

**Exercise:** Rewrite the grammar to make it unambiguous.

```
1: STMT -> if EXPR then STMT
2:       |  if EXPR then STMT else STMT
3:       |  s1
4:       |  s2
5: EXPR -> e1 | e2
```

```
STMT -> MATCHED | OPEN
MATCHED -> if EXPR then MATCHED else MATCHED | s1 | s2
OPEN -> if EXPR then STMT | if EXPR then MATCHED else OPEN
EXPR -> e1 | e2
```

58

# Error Handling

- Objective: detect invalid programs and provide meaningful feedback to programmer
    - Report errors accurately
    - Recover from errors quickly
    - Don't slow down compilation

59

# Error Types

- Many types of errors:
  - Lexical – use  int instead of INT
  - Syntactic – extra brace inserted {
  - Semantic – `float sqr; sqr(2);`
  - Logical – use =  instead of ==

60

# Error Handling - Types

1. Panic mode
2. Error production
3. Automatic local or global correction

61

# Panic Mode Error Handling

- Simplest, most popular
- Discards tokens until one from a set of *synchronizing tokens* is found
- Synchronizing tokens have a clear role
  e.g. semicolons, braces
- E.g. i= i++j

  *policy:* while parsing an expression, discard all tokens until an identifier is found. *This policy skips the additional +*
- Specifying policy in bison: error keyword

  ```
  E -> E + E | (E) | id | error id | error
  ```

62

# Error Productions

- Anticipate common errors
  - 2 x instead of 2 *
- Augment the grammar
  - E -> EE | …
- Disadvantages:
  - Complicates the grammar

63

# Error Corrections

- Rewrite the program – find a "nearby" correct program
  - Local corrections – insert a semicolon, replace a comma with semicolon etc.
  - Global corrections – modify the parse tree with "edit distance" metric in mind
- Disadvantages?
  - Implementation difficulty
  - Slows down compilation
  - Not sure if "nearby" program is intended

64