Compilers INF-400

Burak Arslan ext-inf400@burakarslan.com

Galatasaray Üniversitesi

Lecture V 2023-11-02

Course website

burakarslan.com/inf400

Homework II

Homework II will be out today!

Deadline: 2023-11-08 23:59

Homework I

Any questions?

Can we keep on using regular expressions?

Are regular expressions enough to specify a complete programming language?

Can we keep on using regular expressions?

What is the regular expression for the language of balanced parentheses?

$L = \{(i^{i})^{i} \mid i \geq 0\}$

Can we keep on using regular expressions?

Finite Automata can't:

- Count the number of times a state was visited.
- As a consequence, they can't represent nested structures.



Purpose of parsing





Purpose of parsing

We generate an **Abstract Syntax Tree** – with as little information as possible.

- Parse Tree / Concrete syntax tree: contains every detail. Language-dependent structure
- Abstract Syntax Tree: Just enough to do semantic analysis. Mostly language-independent



Purpose of parsing

Let's see the abstract syntax tree for: 1 + 2;



Purpose of parsing

Let's see the parse tree for: 1 + 2;





Parsers also apply further validation:

Not all token arrays are valid programs!



Statements in programming languages are **nested structures**.



Context Free Grammars. are a perfect match.

Context Free Grammars

Context Free Grammars: Definition:

- \blacktriangleright A set of terminals T
- ► A set of nonterminals N
- A start symbol s
- A set of productions

Context Free Grammars

We will use a tool named **Bison**: Conventions:

- ► Terminal names are UPPERCASE.
- Non-terminals are all lowercase.
- Starting symbol is the first symbol.



Bison algorithm:

- Start from the starting symbol given input token array.
- Replace non terminals by one of the productions on the right ¹
- Repeat until only terminals remain

 $^{^1 {\}rm This}$ point is doint a lot of work here \odot



Bison algorithm:

- Start from the starting symbol given input token array.
- Replace non terminals by one of the productions on the right ¹
- Repeat until only terminals remain

In other words:

$$\blacktriangleright$$
 $s \rightarrow T_0 \dots T_n$

 $^{^1 {\}rm This}$ point is doint a lot of work here \odot

Context Free Grammars

Terminals;

- are supposed to be the tokens recognized by the lexer
- can't be replaced once generated, hence the name "terminal"

Context Free Grammars

Let's see the grammar for the following language: $L = \{(i)^i, i \ge 0\}$

 $ext{expr} o (ext{expr}) ert arepsilon$

Or, in bison notation:

Context Free Grammars

Simple arithmetic example:

$$e
ightarrow e + e | e imes e | (e) |$$
ID 2

Some strings from the above language:

$$x + (y)$$

$$x + y * z$$

$$(x + y) * z$$

 $^{^2\}mathsf{Remember}$ the IDENTIFIER from lexer project? It's the same thing



Derivation: A sequence of productions leading to a string of terminals.

Given:



 $e
ightarrow e + e | e imes e | (e) | ext{ID}$

Let's see the derivation of:

a * b + c



Derivation: A sequence of productions leading to a string of terminals.

Given:

 $e
ightarrow e + e | e imes e | (e) | ext{ID}$

Let's see the derivation of:

a * b + c





Derivation: A sequence of productions leading to a string of terminals.

Given:

 $e
ightarrow e + e | e imes e | (e) | ext{ID}$

Let's see the derivation of:

a * b + c





Derivation: A sequence of productions leading to a string of terminals.

Given:

 $e
ightarrow e + e | e imes e | (e) | ext{ID}$

Let's see the derivation of:

a * b + c

e
e+e
e×e+e
ID×e+e



Derivation: A sequence of productions leading to a string of terminals.

Given:

 $e
ightarrow e + e | e imes e | (e) | ext{ID}$

Let's see the derivation of:

a * b + c

e
e + e
e × e + e
ID × e + e
ID × ID + e



Derivation: A sequence of productions leading to a string of terminals.

Given:

 $e
ightarrow e + e | e imes e | (e) | ext{ID}$

Let's see the derivation of:

a * b + c

e
e + e
e × e + e
ID × e + e
ID × ID + e
ID × ID + ID

Context Free Grammars





Context Free Grammars

Derivation of a * b + c: е е e ID е Х е ID ID





Ambiguity is **BAD**

Context Free Grammars

Ambiguity means your language contains **ill-defined** code fragments³.

³ie. Code fragments with more than one meaning

Context Free Grammars

Dealing with ambiguity:

- Add visible precedence markers (tokens)
- Add implicit precedence markers (%left and %right)
- Write a non-ambiguous grammar

Context Free Grammars

Dealing with ambiguity:

- Add visible precedence markers (tokens, eg. parentheses)
- Add implicit precedence markers (%left and %right)
- Write a non-ambiguous grammar



A non-ambiguous version of the below grammar

$$e
ightarrow e + e | e imes e |$$
ID

could be as follows:

$$e
ightarrow e + e|f + e|f$$

 $f
ightarrow f imes e|f imes f|$ ID

... where the multiplication has precedence over addition

Next Up

- Shift-reduce or Bottom-up parsers (what bison does)
- (Maybe) other parsing algorithms
- The rest of the kiraz grammar