

Automata and Formal Languages (6)

Email: christian.urban at kcl.ac.uk
Office: SI.27 (1st floor Strand Building)
Slides: KEATS (also home work is there)

Regular Languages

While regular expressions are very useful for lexing, there is no regular expression that can recognise the language $a^n b^n$.

$((((()()))))$ vs. $((((()()))()))$

Grammars

A (context-free) grammar G consists of

- a finite set of nonterminal symbols (upper case)
- a finite terminal symbols or tokens (lower case)
- a start symbol (which must be a nonterminal)
- a set of rules

$$A \rightarrow \text{rhs}$$

where rhs are sequences involving terminals and nonterminals, including the empty sequence ϵ .

Grammars

A (context-free) grammar G consists of

- a finite set of nonterminal symbols (upper case)
- a finite terminal symbols or tokens (lower case)
- a start symbol (which must be a nonterminal)
- a set of rules

$$A \rightarrow \text{rhs}$$

where rhs are sequences involving terminals and nonterminals, including the empty sequence ϵ .

We also allow rules

$$A \rightarrow \text{rhs}_1 | \text{rhs}_2 | \dots$$

Palindromes

$$S \rightarrow \epsilon$$

$$S \rightarrow a \cdot S \cdot a$$

$$S \rightarrow b \cdot S \cdot b$$

Palindromes

$$\begin{aligned} S &\rightarrow \epsilon \\ S &\rightarrow a \cdot S \cdot a \\ S &\rightarrow b \cdot S \cdot b \end{aligned}$$

or

$$S \rightarrow \epsilon \mid a \cdot S \cdot a \mid b \cdot S \cdot b$$

Arithmetic Expressions

$E \rightarrow num_token$

$E \rightarrow E \cdot + \cdot E$

$E \rightarrow E \cdot - \cdot E$

$E \rightarrow E \cdot * \cdot E$

$E \rightarrow (\cdot E \cdot)$

Arithmetic Expressions

$E \rightarrow num_token$

$E \rightarrow E \cdot + \cdot E$

$E \rightarrow E \cdot - \cdot E$

$E \rightarrow E \cdot * \cdot E$

$E \rightarrow (\cdot E \cdot)$

1 + 2 * 3 + 4

A CFG Derivation

- 1 Begin with a string containing only the start symbol, say S
- 2 Replace any nonterminal X in the string by the right-hand side of some production $X \rightarrow \text{rhs}$
- 3 Repeat 2 until there are no nonterminals

$S \rightarrow \dots \rightarrow \dots \rightarrow \dots \rightarrow \dots$

Example Derivation

$$S \rightarrow \epsilon \mid a \cdot S \cdot a \mid b \cdot S \cdot b$$

$$\begin{aligned} S &\rightarrow aSa \\ &\rightarrow abSba \\ &\rightarrow abaSaba \\ &\rightarrow abaaba \end{aligned}$$

Example Derivation

$$E \rightarrow num_token$$

$$E \rightarrow E \cdot + \cdot E$$

$$E \rightarrow E \cdot - \cdot E$$

$$E \rightarrow E \cdot * \cdot E$$

$$E \rightarrow (\cdot E \cdot)$$

$$E \rightarrow E * E$$

$$\rightarrow E + E * E$$

$$\rightarrow E + E * E + E$$

$$\rightarrow^+ 1 + 2 * 3 + 4$$

Example Derivation

$$E \rightarrow \textit{num_token}$$

$$E \rightarrow E \cdot + \cdot E$$

$$E \rightarrow E \cdot - \cdot E$$

$$E \rightarrow E \cdot * \cdot E$$

$$E \rightarrow (\cdot E \cdot)$$

$$E \rightarrow E * E$$

$$\rightarrow E + E * E$$

$$\rightarrow E + E * E + E$$

$$\rightarrow^+ 1 + 2 * 3 + 4$$

$$E \rightarrow E + E$$

$$\rightarrow E + E + E$$

$$\rightarrow E + E * E + E$$

$$\rightarrow^+ 1 + 2 * 3 + 4$$

Language of a CFG

Let G be a context-free grammar with start symbol S . Then the language $L(G)$ is:

$$\{c_1 \dots c_n \mid \forall i. c_i \in T \wedge S \rightarrow^* c_1 \dots c_n\}$$

Language of a CFG

Let G be a context-free grammar with start symbol S . Then the language $L(G)$ is:

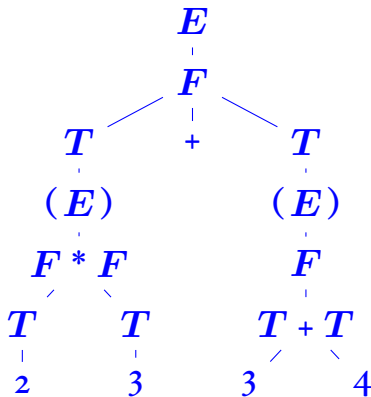
$$\{c_1 \dots c_n \mid \forall i. c_i \in T \wedge S \rightarrow^* c_1 \dots c_n\}$$

- Terminals, because there are no rules for replacing them.
- Once generated, terminals are “permanent”.
- Terminals ought to be tokens of the language (but can also be strings).

Parse Trees

$$E \rightarrow F \mid F \cdot * \cdot F$$
$$F \rightarrow T \mid T \cdot + \cdot T \mid T \cdot - \cdot T$$
$$T \rightarrow \text{num_token} \mid (\cdot E \cdot)$$

$(2 * 3) + (3 + 4)$



Arithmetic Expressions

$E \rightarrow num_token$

$E \rightarrow E \cdot + \cdot E$

$E \rightarrow E \cdot - \cdot E$

$E \rightarrow E \cdot * \cdot E$

$E \rightarrow (\cdot E \cdot)$

Arithmetic Expressions

$$E \rightarrow \textit{num_token}$$

$$E \rightarrow E \cdot + \cdot E$$

$$E \rightarrow E \cdot - \cdot E$$

$$E \rightarrow E \cdot * \cdot E$$

$$E \rightarrow (\cdot E \cdot)$$

A CFG is **left-recursive** if it has a nonterminal E such that $E \rightarrow^+ E \cdot \dots$

Ambiguous Grammars

A grammar is **ambiguous** if there is a string that has at least two different parse trees.

$$E \rightarrow \textit{num_token}$$

$$E \rightarrow E \cdot + \cdot E$$

$$E \rightarrow E \cdot - \cdot E$$

$$E \rightarrow E \cdot * \cdot E$$

$$E \rightarrow (\cdot E \cdot)$$

1 + 2 * 3 + 4

Dangling Else

Another ambiguous grammar:

$$\begin{array}{l} E \rightarrow \text{if } E \text{ then } E \\ \quad | \text{if } E \text{ then } E \text{ else } E \\ \quad | \dots \end{array}$$

if a then if x then y else c

Parser Combinators

Parser combinators:

$\underbrace{\text{list of tokens}}_{\text{input}} \Rightarrow \underbrace{\text{set of (parsed input, unparsed input)}}_{\text{output}}$

- sequencing
- alternative
- semantic action

Alternative parser (code $p \parallel q$)

- apply p and also q ; then combine the outputs

$$p(\text{input}) \cup q(\text{input})$$

Sequence parser (code $p \sim q$)

- apply first p producing a set of pairs
- then apply q to the unparsed parts
- then combine the results:
((output₁, output₂), unparsed part)

$$\{((o_1, o_2), u_2) \mid (o_1, u_1) \in p(\text{input}) \wedge (o_2, u_2) \in q(u_1)\}$$

Function parser (code $p \Rightarrow f$)

- apply p producing a set of pairs
- then apply the function f to each first component

$$\{(f(o_1), u_1) \mid (o_1, u_1) \in p(\text{input})\}$$

Function parser (code $p \Rightarrow f$)

- apply p producing a set of pairs
- then apply the function f to each first component

$$\{(f(o_1), u_1) \mid (o_1, u_1) \in p(\text{input})\}$$

f is the semantic action (“what to do with the parsed input”)

Semantic Actions

Addition

$$T \sim + \sim E \Rightarrow \underbrace{f((x, y), z) \Rightarrow x + z}_{\text{semantic action}}$$

Semantic Actions

Addition

$$T \sim + \sim E \Rightarrow \underbrace{f((x, y), z) \Rightarrow x + z}_{\text{semantic action}}$$

Multiplication

$$F \sim * \sim T \Rightarrow f((x, y), z) \Rightarrow x * z$$

Semantic Actions

Addition

$$T \sim + \sim E \Rightarrow \underbrace{f((x, y), z) \Rightarrow x + z}_{\text{semantic action}}$$

Multiplication

$$F \sim * \sim T \Rightarrow f((x, y), z) \Rightarrow x * z$$

Parenthesis

$$(\sim E \sim) \Rightarrow f((x, y), z) \Rightarrow y$$

Types of Parsers

- **Sequencing:** if p returns results of type T , and q results of type S , then $p \sim q$ returns results of type

$$T \times S$$

Types of Parsers

- **Sequencing:** if p returns results of type T , and q results of type S , then $p \sim q$ returns results of type

$$T \times S$$

- **Alternative:** if p returns results of type T then q **must** also have results of type T , and $p \parallel q$ returns results of type

$$T$$

Types of Parsers

- **Sequencing:** if p returns results of type T , and q results of type S , then $p \sim q$ returns results of type

$$T \times S$$

- **Alternative:** if p returns results of type T then q **must** also have results of type T , and $p \parallel q$ returns results of type

$$T$$

- **Semantic Action:** if p returns results of type T and f is a function from T to S , then $p \Rightarrow f$ returns results of type

$$S$$

Input Types of Parsers

- input: **string**
- output: set of (output_type, **string**)

Input Types of Parsers

- input: **string**
- output: set of (output_type, **string**)

actually it can be any input type as long as it is a kind of sequence (for example a string)

Scannerless Parsers

- input: **string**
- output: set of (output_type, **string**)

but lexers are better when whitespaces or comments need to be filtered out; then input is a sequence of tokens

Successful Parses

- input: string
- output: **set of** (output_type, string)

a parse is successful whenever the input has been fully “consumed” (that is the second component is empty)

Abstract Parser Class

```
1 abstract class Parser[I, T] {  
2   def parse(ts: I): Set[(T, I)]  
3  
4   def parse_all(ts: I) : Set[T] =  
5     for ((head, tail) <- parse(ts); if (tail.isEmpty))  
6       yield head  
7 }
```

```

1 class AltParser[I, T](p: => Parser[I, T],
2                       q: => Parser[I, T])
3     extends Parser[I, T] {
4   def parse(sb: I) = p.parse(sb) ++ q.parse(sb)
5 }
6
7 class SeqParser[I, T, S](p: => Parser[I, T],
8                          q: => Parser[I, S])
9     extends Parser[I, (T, S)] {
10  def parse(sb: I) =
11    for ((head1, tail1) <- p.parse(sb);
12         (head2, tail2) <- q.parse(tail1))
13      yield ((head1, head2), tail2)
14 }
15
16 class FunParser[I, T, S](p: => Parser[I, T], f: T => S)
17     extends Parser[I, S] {
18  def parse(sb: I) =
19    for ((head, tail) <- p.parse(sb))
20      yield (f(head), tail)
21 }

```

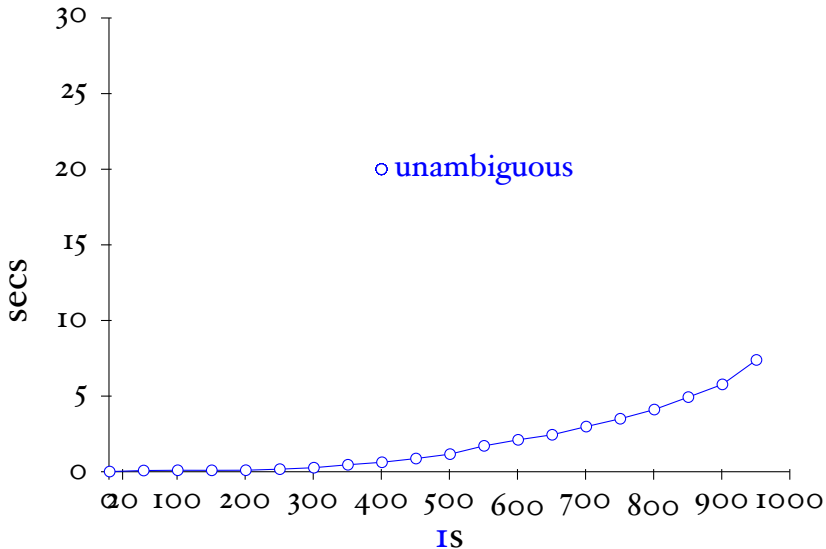
Two Grammars

Which languages are recognised by the following two grammars?

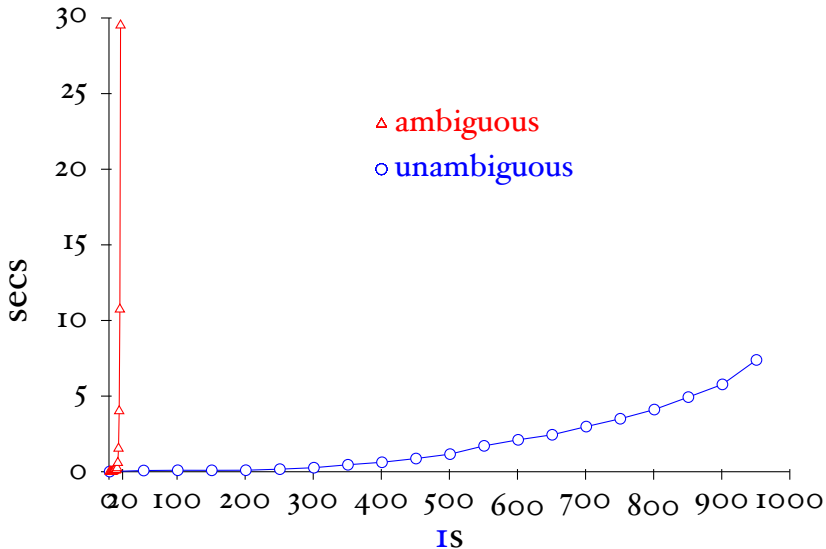
$$\begin{array}{c} S \rightarrow 1 \cdot S \cdot S \\ | \quad \epsilon \end{array}$$

$$\begin{array}{c} U \rightarrow 1 \cdot U \\ | \quad \epsilon \end{array}$$

Ambiguous Grammars



Ambiguous Grammars



While-Language

Stmt → skip
| *Id := AExp*
| if *BExp* then *Block* else *Block*
| while *BExp* do *Block*

Stmts → *Stmt ; Stmts*
| *Stmt*

Block → {*Stmts*}
| *Stmt*

AExp → ...

BExp → ...

An Interpreter

```
{  
   $x := 5;$   
   $y := x * 3;$   
   $y := x * 4;$   
   $x := u * 3$   
}
```

- the interpreter has to record the value of x before assigning a value to y

An Interpreter

```
{  
   $x := 5;$   
   $y := x * 3;$   
   $y := x * 4;$   
   $x := u * 3$   
}
```

- the interpreter has to record the value of x before assigning a value to y
- `eval(stmt, env)`