# Compilers and Formal Languages (6)
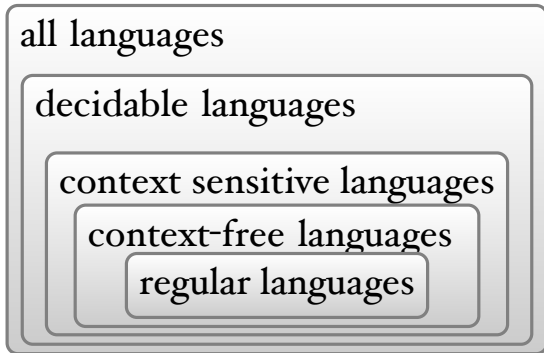
Email:   christian.urban at kcl.ac.uk
Office:  N7.07 (North Wing, Bush House)
Slides:  KEATS (also home work is there)

# Hierarchy of Languages

Recall that languages are sets of strings.



all languages
  decidable languages
    context sensitive languages
      context-free languages
        regular languages

Atomic parsers, for example

$$1 :: \textit{rest} \;\Rightarrow\; \big\{ (1, \textit{rest}) \big\}$$

- you consume one or more tokens from the input (stream)
- also works for characters and strings

Alternative parser (code $p \;||\; q$)

- 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:

$$((\text{output}_1, \text{output}_2), \text{unparsed part})$$

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

Function parser (code $p \Rightarrow f$)

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

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

# 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 $p \mid\mid q$ 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 $p \mathbin{||} q$ 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
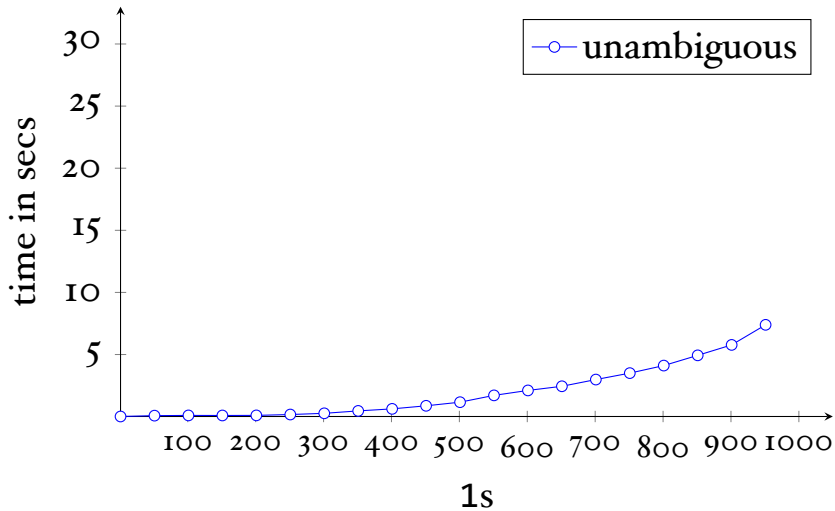
$$S$$

# Two Grammars

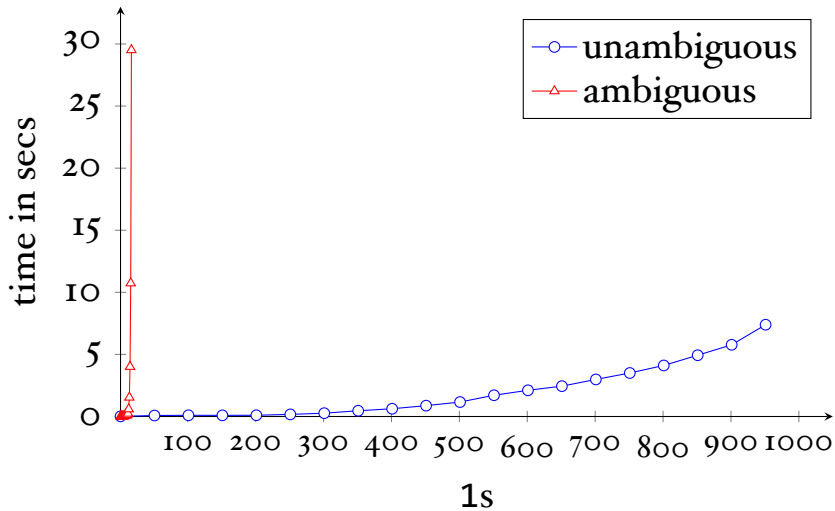Which languages are recognised by the following two grammars?

$$S ::= \mathbf{1} \cdot S \cdot S \mid \epsilon$$

$$U ::= \mathbf{1} \cdot U \mid \epsilon$$

# Ambiguous Grammars

# Ambiguous Grammars

# Arithmetic Expressions

A grammar for arithmetic expressions and numbers:

$$E ::= E \cdot + \cdot E \mid E \cdot * \cdot E \mid (\cdot E \cdot) \mid N$$
$$N ::= N \cdot N \mid 0 \mid 1 \mid \ldots \mid 9$$

Unfortunately it is left-recursive (and ambiguous).

A problem for <span style="color:red">recursive descent parsers</span> (e.g. parser combinators).

# Arithmetic Expressions

A grammar for arithmetic expressions and numbers:

$$E ::= E \cdot + \cdot E \mid E \cdot * \cdot E \mid (\cdot E \cdot) \mid N$$
$$N ::= N \cdot N \mid 0 \mid 1 \mid \ldots \mid 9$$

Unfortunately it is left-recursive (and ambiguous).

A problem for recursive descent parsers (e.g. parser combinators).

# Numbers

$$N ::= N \cdot N \mid 0 \mid 1 \mid \dots \mid 9$$

A non-left-recursive, non-ambiguous grammar for numbers:

$$N ::= 0 \cdot N \mid 1 \cdot N \mid \dots \mid 0 \mid 1 \mid \dots \mid 9$$

# Removing Left-Recursion

The rule for numbers is directly left-recursive:

$$N \ ::= \ N \cdot N \mid 0 \mid 1 \quad (\ldots)$$

Translate

$$
\begin{aligned}
N \ &::= \ N \cdot \alpha \\
&\mid \ \beta
\end{aligned}
\quad \Longrightarrow \quad
\begin{aligned}
N \ &::= \ \beta \cdot N' \\
N' \ &::= \ \alpha \cdot N' \\
&\mid \ \epsilon
\end{aligned}
$$

# Removing Left-Recursion

The rule for numbers is directly left-recursive:

$$N \; ::= \; N \cdot N \mid 0 \mid 1 \quad (\ldots)$$

Translate

$$
\begin{array}{llll}
N & ::= & N \cdot \alpha & \\
& \mid & \beta & \Longrightarrow
\end{array}
\qquad
\begin{array}{lll}
N & ::= & \beta \cdot N' \\
N' & ::= & \alpha \cdot N' \\
& \mid & \epsilon
\end{array}
$$

Which means in this case:

$$
\begin{array}{lll}
N & \rightarrow & 0 \cdot N' \mid 1 \cdot N' \\
N' & \rightarrow & N \cdot N' \mid \epsilon
\end{array}
$$

# Operator Precedences

To disambiguate

$$E \quad ::= \quad E \cdot + \cdot E \mid E \cdot * \cdot E \mid (\cdot E \cdot) \mid N$$

Decide on how many precedence levels, say

highest for $()$, medium for $*$, lowest for $+$

$$
\begin{aligned}
E_{low} \quad &::= \quad E_{med} \cdot + \cdot E_{low} \mid E_{med} \\
E_{med} \quad &::= \quad E_{hi} \cdot * \cdot E_{med} \mid E_{hi} \\
E_{hi} \quad &::= \quad (\cdot E_{low} \cdot) \mid N
\end{aligned}
$$

# Operator Precedences

To disambiguate

$$E \quad ::= \quad E \cdot + \cdot E \mid E \cdot * \cdot E \mid (\cdot E \cdot) \mid N$$

Decide on how many precedence levels, say

highest for $()$, medium for $*$, lowest for $+$

$$
\begin{aligned}
E_{low} \quad &::= \quad E_{med} \cdot + \cdot E_{low} \mid E_{med} \\
E_{med} \quad &::= \quad E_{hi} \cdot * \cdot E_{med} \mid E_{hi} \\
E_{hi} \quad &::= \quad (\cdot E_{low} \cdot) \mid N
\end{aligned}
$$

What happens with $1 + 3 + 4$?

# Chomsky Normal Form

All rules must be of the form

$$A ::= a$$

or

$$A ::= B \cdot C$$

No rule can contain $\epsilon$.

# $\epsilon$-**Removal**

1. If $A ::= \alpha \cdot B \cdot \beta$ and $B ::= \epsilon$ are in the grammar, then add $A ::= \alpha \cdot \beta$ (iterate if necessary).
2. Throw out all $B ::= \epsilon$.

$$N ::= \text{o} \cdot N' \mid \text{1} \cdot N'$$
$$N' ::= N \cdot N' \mid \epsilon$$

$$N ::= \text{o} \cdot N' \mid \text{1} \cdot N' \mid \text{o} \mid \text{1}$$
$$N' ::= N \cdot N' \mid N \mid \epsilon$$

$$N ::= \text{o} \cdot N' \mid \text{1} \cdot N' \mid \text{o} \mid \text{1}$$
$$N' ::= N \cdot N' \mid N$$

# $\epsilon$-**Removal**

① If $A ::= \alpha \cdot B \cdot \beta$ and $B ::= \epsilon$ are in the grammar, then add $A ::= \alpha \cdot \beta$ (iterate if necessary).

② Throw out all $B ::= \epsilon$.

$$N ::= \text{o} \cdot N' \mid \text{1} \cdot N'$$
$$N' ::= N \cdot N' \mid \epsilon$$

$$N ::= \text{o} \cdot N' \mid \text{1} \cdot N' \mid \text{o} \mid \text{1}$$
$$N' ::= N \cdot N' \mid N \mid \epsilon$$

$$N ::= \text{o} \cdot N' \mid \text{1} \cdot N' \mid \text{o} \mid \text{1}$$
$$N' ::= N \cdot N' \mid N$$

$$N ::= \text{o} \cdot N \mid \text{1} \cdot N \mid \text{o} \mid \text{1}$$

# CYK Algorithm

If grammar is in Chomsky normalform ...

$S$  ::=  $N \cdot P$
$P$  ::=  $V \cdot N$
$N$  ::=  $N \cdot N$
$N$  ::=  students | Jeff | geometry | trains
$V$  ::=  trains

Jeff trains geometry students

# CYK Algorithm

- fastest possible algorithm for recognition problem
- runtime is $O(n^3)$

- grammars need to be transformed into CNF

# The Goal of this Course

## Write a Compiler



We have lexer and parser.

$$
\begin{array}{lll}
Stmt & ::= & \texttt{skip} \\
& | & Id := AExp \\
& | & \texttt{if } BExp \texttt{ then } Block \texttt{ else } Block \\
& | & \texttt{while } BExp \texttt{ do } Block \\
& | & \texttt{read } Id \\
& | & \texttt{write } Id \\
& | & \texttt{write } String \\
\\
Stmts & ::= & Stmt \; ; \; Stmts \\
& | & Stmt \\
\\
Block & ::= & \{ \; Stmts \; \} \\
& | & Stmt \\
\\
AExp & ::= & ... \\
BExp & ::= & ...
\end{array}
$$

```
write "Fib";
read n;
minus1 := 0;
minus2 := 1;
while n > 0 do {
        temp := minus2;
        minus2 := minus1 + minus2;
        minus1 := temp;
        n := n - 1
};
write "Result";
write minus2
```

# An Interpreter

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

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

# An Interpreter

$$
\begin{aligned}
&\{ \\
&\quad x := 5; \\
&\quad y := x * 3; \\
&\quad y := x * 4; \\
&\quad x := u * 3 \\
&\}
\end{aligned}
$$

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

# Interpreter

$$\text{eval}(n, E) \overset{\text{def}}{=} n$$

$$\text{eval}(x, E) \overset{\text{def}}{=} E(x) \quad \text{lookup } x \text{ in } E$$

$$\text{eval}(a_1 + a_2, E) \overset{\text{def}}{=} \text{eval}(a_1, E) + \text{eval}(a_2, E)$$

$$\text{eval}(a_1 - a_2, E) \overset{\text{def}}{=} \text{eval}(a_1, E) - \text{eval}(a_2, E)$$

$$\text{eval}(a_1 * a_2, E) \overset{\text{def}}{=} \text{eval}(a_1, E) * \text{eval}(a_2, E)$$

$$\text{eval}(a_1 = a_2, E) \overset{\text{def}}{=} \text{eval}(a_1, E) = \text{eval}(a_2, E)$$

$$\text{eval}(a_1 \,!= a_2, E) \overset{\text{def}}{=} \neg(\text{eval}(a_1, E) = \text{eval}(a_2, E))$$

$$\text{eval}(a_1 < a_2, E) \overset{\text{def}}{=} \text{eval}(a_1, E) < \text{eval}(a_2, E)$$

# Interpreter (2)

$$\text{eval}(\text{skip}, E) \quad \overset{\text{def}}{=} \quad E$$

$$\text{eval}(x := a, E) \quad \overset{\text{def}}{=} \quad E(x \mapsto \text{eval}(a, E))$$

$$\text{eval}(\text{if } b \text{ then } cs_1 \text{ else } cs_2, E) \overset{\text{def}}{=}$$
$$\text{if eval}(b, E) \text{ then eval}(cs_1, E)$$
$$\text{else eval}(cs_2, E)$$

$$\text{eval}(\text{while } b \text{ do } cs, E) \overset{\text{def}}{=}$$
$$\text{if eval}(b, E)$$
$$\text{then eval}(\text{while } b \text{ do } cs, \text{eval}(cs, E))$$
$$\text{else } E$$

$$\text{eval}(\text{write } x, E) \quad \overset{\text{def}}{=} \quad \{ \text{ println}(E(x)) \; ; \; E \}$$

# Test Program

```
start := 1000;
x := start;
y := start;
z := start;
while 0 < x do {
 while 0 < y do {
  while 0 < z do { z := z - 1 };
  z := start;
  y := y - 1
 };
 y := start;
 x := x - 1
}
```

# Interpreted Code

# Java Virtual Machine

- introduced in 1995
- is a stack-based VM (like Postscript, CLR of .Net)
- contains a JIT compiler
- many languages take advantage of JVM's infrastructure (JRE)
- is garbage collected ⇒ no buffer overflows
- some languages compile to the JVM: Scala, Clojure...