

Compilers and Formal Languages

Email: christian.urban at kcl.ac.uk

Slides & Progs: KEATS (also homework is there)

1 Introduction, Languages	6 While-Language
2 Regular Expressions, Derivatives	7 Compilation, JVM
3 Automata, Regular Languages	8 Compiling Functional Languages
4 Lexing, Tokenising	9 Optimisations
5 Grammars, Parsing	10 LLVM

Parser Combinators

One of the simplest ways to implement a parser, see <https://vimeo.com/142341803> (by Haoyi Li)

- build-in library in Scala
- fastparse (2) library by Haoyi Li; is part of Ammonite
- possible exponential runtime behaviour

Parser Combinators

Parser combinators:

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

atomic parsers

sequencing

alternative

semantic action (map-parser)

Atomic parsers, for example, number tokens

$$\text{Num}(123) :: \text{rest} \Rightarrow \{(\text{Num}(123), \text{rest})\}$$

you consume one or more token from the
input (stream)

also works for characters and strings

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 part
then combine the results:

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

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

Map-parser (code $p.map(f)$)

apply p producing a set of pairs

then apply the function f to each first component

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

Map-parser (code $p.map(f)$)

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 returns 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 returns 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 returns 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: token list

output: set of (output_type, token list)

Input Types of Parsers

input: **token list**

output: set of (output_type, **token list**)

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 using lexers is better because whitespaces or comments can 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

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

```

class AltParser[I, T](p: => Parser[I, T],
                    q: => Parser[I, T])
    extends Parser[I, T] {
  def parse(sb: I) = p.parse(sb) ++ q.parse(sb)
}

class SeqParser[I, T, S](p: => Parser[I, T],
                       q: => Parser[I, S])
    extends Parser[I, (T, S)] {
  def parse(sb: I) =
    for ((head1, tail1) <- p.parse(sb);
         (head2, tail2) <- q.parse(tail1))
      yield ((head1, head2), tail2)
}

class FunParser[I, T, S](p: => Parser[I, T], f: T => S)
    extends Parser[I, S] {
  def parse(sb: I) =
    for ((head, tail) <- p.parse(sb))
      yield (f(head), tail)
}

```

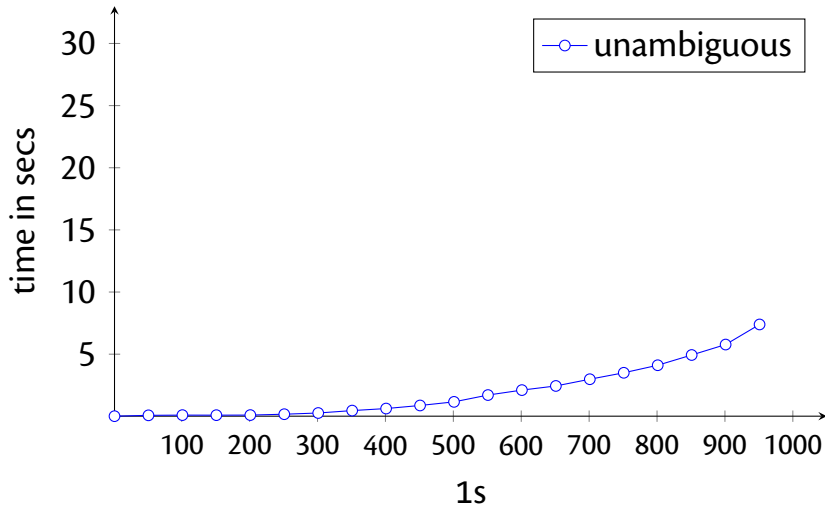
Two Grammars

Which languages are recognised by the following two grammars?

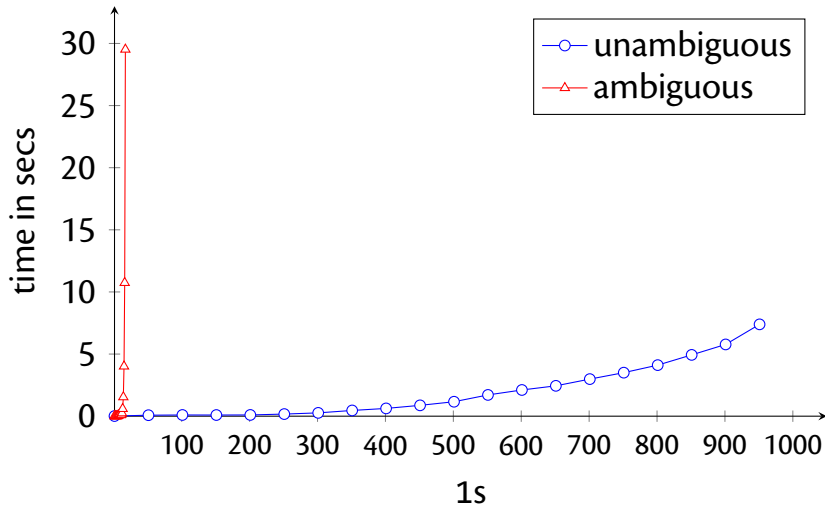
$$\begin{array}{l} S \rightarrow 1 \cdot S \cdot S \\ | \quad \epsilon \end{array}$$

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

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)
```

Interpreter

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

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

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

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

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

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

$$\text{eval}(a_1 \neq a_2, E) \stackrel{\text{def}}{=} \neg(\text{eval}(a_1, E) = \text{eval}(a_2, E))$$

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

Interpreter (2)

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

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

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

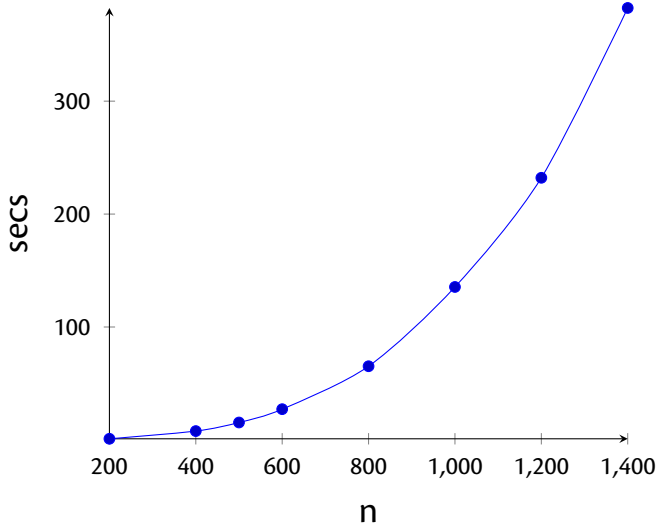
$$\text{eval}(\text{while } b \text{ do } cs, E) \stackrel{\text{def}}{=} \\ \text{if } \text{eval}(b, E) \\ \text{then } \text{eval}(\text{while } b \text{ do } cs, \text{eval}(cs, E)) \\ \text{else } E$$

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

Test Program

??

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 \Rightarrow no buffer overflows

some languages compile to the JVM: Scala, Clojure...

Compilers and Formal Languages

Email: christian.urban at kcl.ac.uk

Slides & Progs: KEATS (also homework is there)

1 Introduction, Languages	6 While-Language
2 Regular Expressions, Derivatives	7 Compilation, JVM
3 Automata, Regular Languages	8 Compiling Functional Languages
4 Lexing, Tokenising	9 Optimisations
5 Grammars, Parsing	10 LLVM

Starting Symbol

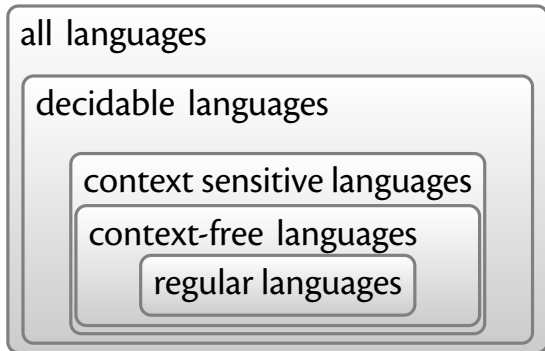
$$S ::= A \cdot S \cdot B \mid B \cdot S \cdot A \mid \epsilon$$
$$A ::= a \mid \epsilon$$
$$B ::= b$$

TODO: Testcases for math expressions

<https://github.com/ArashPartow/math-parser-benchmark-project>

Hierarchy of Languages

Recall that languages are sets of strings.



Parser Combinators

Atomic parsers, for example

$$1 \text{ :: } rest \Rightarrow \{(1, rest)\}$$

you consume one or more tokens from the
input (stream)

also works for characters and strings

Alternative parser (code $p \mid q$)

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 returns 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 returns 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 $p \mid q$ returns results of type

$$T$$

Types of Parsers

Sequencing: if p returns results of type T , and q returns 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 $p \mid q$ 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

$$S$$

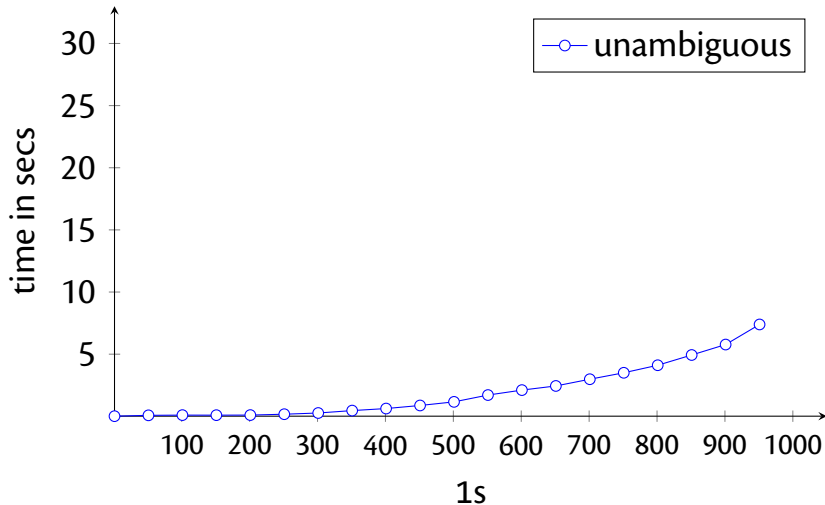
Two Grammars

Which languages are recognised by the following two grammars?

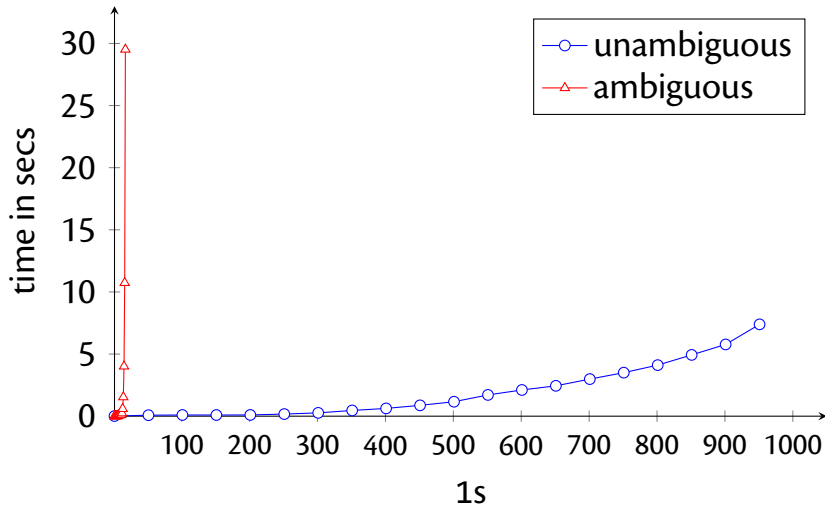
$$S ::= 1 \cdot S \cdot S \mid \epsilon$$

$$U ::= 1 \cdot U \mid \epsilon$$

Ambiguous Grammars



Ambiguous Grammars



Arithmetic Expressions

A grammar for arithmetic expressions and numbers:

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

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

Unfortunately it is left-recursive (and ambiguous).

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

Arithmetic Expressions

A grammar for arithmetic expressions and numbers:

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

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

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:

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

Translate

$$\begin{array}{l} N ::= N \cdot \alpha \\ \quad \mid \beta \end{array} \Rightarrow \begin{array}{l} N ::= \beta \cdot N' \\ N' ::= \alpha \cdot N' \\ \quad \mid \epsilon \end{array}$$

Removing Left-Recursion

The rule for numbers is directly left-recursive:

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

Translate

$$\begin{array}{l} N ::= N \cdot \alpha \\ \quad \mid \beta \end{array} \Rightarrow \begin{array}{l} N ::= \beta \cdot N' \\ N' ::= \alpha \cdot N' \\ \quad \mid \epsilon \end{array}$$

Which means in this case:

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

Chomsky Normal Form

All rules must be of the form

$$A ::= a$$

or

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

No rule can contain ϵ .

ϵ -Removal

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

Throw out all $B ::= \epsilon$.

$$N ::= 0 \cdot N' \mid 1 \cdot N'$$

$$N' ::= N \cdot N' \mid \epsilon$$

$$N ::= 0 \cdot N' \mid 1 \cdot N' \mid 0 \mid 1$$

$$N' ::= N \cdot N' \mid N \mid \epsilon$$

$$N ::= 0 \cdot N' \mid 1 \cdot N' \mid 0 \mid 1$$

$$N' ::= N \cdot N' \mid N$$

ϵ -Removal

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

Throw out all $B ::= \epsilon$.

$$N ::= 0 \cdot N' \mid 1 \cdot N'$$

$$N' ::= N \cdot N' \mid \epsilon$$

$$N ::= 0 \cdot N' \mid 1 \cdot N' \mid 0 \mid 1$$

$$N' ::= N \cdot N' \mid N \mid \epsilon$$

$$N ::= 0 \cdot N' \mid 1 \cdot N' \mid 0 \mid 1$$

$$N' ::= N \cdot N' \mid N$$

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

CYK Algorithm

If grammar is in Chomsky normalform ...

$S ::= N \cdot P$

$P ::= V \cdot N$

$N ::= N \cdot N$

$N ::= \text{students} \mid \text{Jeff} \mid \text{geometry} \mid \text{trains}$

$V ::= \text{trains}$

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 a lexer and a parser...

Stmt ::= skip
| *Id* := *AExp*
| if *BExp* then *Block* else *Block*
| while *BExp* do *Block*
| read *Id*
| write *Id*
| write *String*

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)`

An Interpreter

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

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

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

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

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

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

$$\text{eval}(a_1 \neq a_2, E) \stackrel{\text{def}}{=} \neg(\text{eval}(a_1, E) = \text{eval}(a_2, E))$$

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

An Interpreter (2)

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

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

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

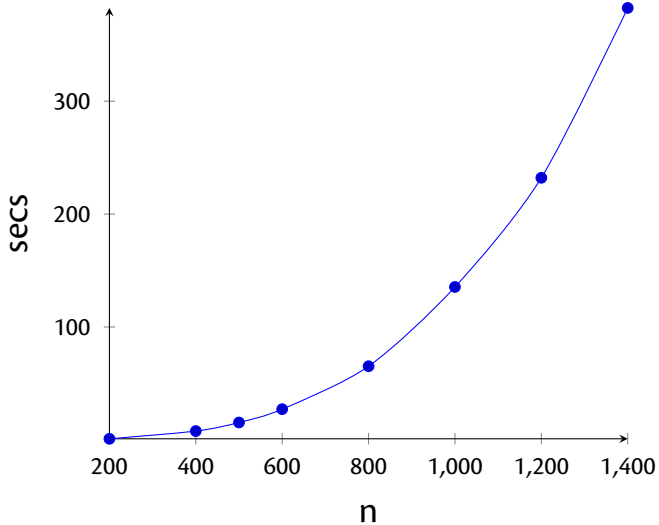
$\text{eval}(\text{while } b \text{ do } cs, E) \stackrel{\text{def}}{=} \\ \text{if } \text{eval}(b, E) \\ \text{then } \text{eval}(\text{while } b \text{ do } cs, \text{eval}(cs, E)) \\ \text{else } E$

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

Test Program

??

Interpreted Code



Java Virtual Machine

introduced in 1995

is a stack-based VM (like Postscript, CLR of .Net)

contains a JIT compiler

From the Cradle to the Holy Graal - the JDK Story

<https://www.youtube.com/watch?v=h419kfbLhUI>

is garbage collected \Rightarrow no buffer overflows

some languages compile to the JVM: Scala,
Clojure...

LLVM

LLVM started by academics in 2000 (University of Illinois in Urbana-Champaign)

suite of compiler tools

SSA-based intermediate language

no need to allocate registers

source languages: C, C++, Rust, Go, Swift

target CPUs: x86, ARM, PowerPC, ...