



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]

Haskell, Do You Read Me?

Constructing and Composing Efficient Top-down Parsers at Runtime

Marcos Viera Doaitse Swierstra Eelco Lempink

Instituto de Computación, Universidad de la República, Uruguay
Dept. of Information and Computing Sciences, Utrecht University

September 25, 2008

Symptoms

```
data T1 = T1 <: T1
        | T1 >: T1
        | C
    deriving (Read, Show)
```

```
infixl 5 <:
```

```
infixr 6 >:
```

```
x      :: T1
```

```
x = C <: C <: C
```

```
*Main> x
```

```
(C <: C) <: C
```



Symptoms

```
data T1 = T1 <: T1
        | T1 >: T1
        | C
    deriving (Read, Show)
```

```
infixl 5 <:
```

```
infixr 6 >:
```

```
x, x'    :: T1
```

```
x = C <: C <: C
```

```
x' = (read ∘ show) $ C <: C <: C
```

```
*Main> x'
```

```
(C <: C) <: C
```



Symptoms

```
data T1 = T1 <: T1
        | T1 >: T1
        | C
    deriving (Read, Show)
```

```
infixl 5 <:
```

```
infixr 6 >:
```

```
 $x, x', x'' :: T1$ 
```

```
 $x = C <: C <: C$ 
```

```
 $x' = (read \circ show) \$ C <: C <: C$ 
```

```
 $x'' = read "C <: C <: C"$ 
```

```
*Main> x''
```

```
*** Exception: Prelude.read: no parse
```



Symptoms

```
data T1 = T1 <: T1
        | T1 >: T1
        | C
    deriving (Read, Show)
```

infixl 5 <:

infixr 6 >:

$x, x', x'' :: T1$

$x = C <: C <: C$

$x' = (read \circ show) \$ C <: C <: C$

$x'' = read "C <: C <: C"$

Ideally, you should be able to *read* every valid constant expression.



Parentheses

```
*Main> time (read "C" :: T1)
C
CPU Time: 0 ms
```



Parentheses

```
*Main> time (read "C" :: T1)
```

```
C
```

```
CPU Time: 0 ms
```

```
*Main> time (read "((((((C))))))" :: T1)
```

```
C
```

```
CPU Time: 74 ms
```



Parentheses

```
*Main> time (read "C" :: T1)
```

```
C
```

```
CPU Time: 0 ms
```

```
*Main> time (read "((((C)))") :: T1)
```

```
C
```

```
CPU Time: 74 ms
```

```
*Main> time (read "((((((C))))))" :: T1)
```

```
C
```

```
CPU Time: 389 ms
```



Parentheses

```
*Main> time (read "C" :: T1)
```

```
C
```

```
CPU Time: 0 ms
```

```
*Main> time (read "((((C))))" :: T1)
```

```
C
```

```
CPU Time: 74 ms
```

```
*Main> time (read "((((((C))))))" :: T1)
```

```
C
```

```
CPU Time: 389 ms
```

```
*Main> time (read "((((((((C)))))))" :: T1)
```

```
C
```

```
CPU Time: 1753 ms
```



Breadth-first Parsing

The language which is actually recognised by the generated *read* function is described by the non left-recursive grammar:

$T1(n) \rightarrow T1(6) \text{ ":"} < \text{" } T1(6)$	$(n \leq 5)$
$T1(7) \text{ ":"} > \text{" } T1(7)$	$(n \leq 6)$
"C"	$(n \leq 10)$
$\text{"(" } T1(0) \text{ "}"$	$(n \leq 10)$



Breadth-first Parsing

The language which is actually recognised by the generated *read* function is described by the non left-recursive grammar:

$T1(n) \rightarrow T1(6) \text{ ":<:" } T1(6)$	$(n \leq 5)$
$T1(7) \text{ ":>:" } T1(7)$	$(n \leq 6)$
"C"	$(n \leq 10)$
$\text{"(" } T1(0) \text{ ")"}$	$(n \leq 10)$

Three parallel parsers are started up for the first '(', and so on recursively.



Common Left-factors

Unfortunately the problem also shows up for more reasonable expressions such as $C :>: (C :>: (C :>: \dots))$.

We remove the conditions, and encode them in the non-terminals.

$$\begin{aligned} T1 (0..5) &\rightarrow T1 (6) " :<:" T1 (6) \mid T1 (6) \\ T1 (6) &\rightarrow T1 (7) " :>:" T1 (7) \mid T1 (7) \\ T1 (7..10) &\rightarrow "C" \\ &\mid "(" T1 (0) ")" \end{aligned}$$


Common Left-factors

Unfortunately the problem also shows up for more reasonable expressions such as $C :>: (C :>: (C :>: \dots))$.

We remove the conditions, and encode them in the non-terminals.

$$\begin{aligned} T1 (0..5) &\rightarrow T1 (6) \text{ ":"} T1 (6) \mid T1 (6) \\ T1 (6) &\rightarrow T1 (7) \text{ ":"} T1 (7) \mid T1 (7) \\ T1 (7..10) &\rightarrow \text{"C"} \\ &\mid \text{"(" } T1 (0) \text{ ")" } \end{aligned}$$

We see that some alternatives start with the same non-terminal symbol.



The Problem

- ▶ Derived *read* functions treat all operators as being *non-associative*, despite their declared associativities and precedences.



The Problem

- ▶ Derived *read* functions treat all operators as being *non-associative*, despite their declared associativities and precedences.
- ▶ Derived *show* functions generate the needed extra parentheses, in order to make *read* \circ *show* work.



The Problem

- ▶ Derived *read* functions treat all operators as being *non-associative*, despite their declared associativities and precedences.
- ▶ Derived *show* functions generate the needed extra parentheses, in order to make *read* \circ *show* work.
- ▶ These extra parentheses make parsing take exponential time.



The Problem

- ▶ Derived *read* functions treat all operators as being *non-associative*, despite their declared associativities and precedences.
- ▶ Derived *show* functions generate the needed extra parentheses, in order to make *read* \circ *show* work.
- ▶ These extra parentheses make parsing take exponential time.
- ▶ Common left-factors have a similar effect.



How Does the Problem Arise?

```
infix 5 :+:  
infix 6 :*:  
data T2 a = T2 a :+ : T2 a  
           | a    :* : T2 a  
           | C2  
           deriving Read  
t2 :: T2 (T2 Int)  
t2 = read "(3 :* C2) :* C2"
```

The function *read* is a member of the class *Read*:



How Does the Problem Arise?

```
infix 5 :+:  
infix 6 :*:  
data T2 a = T2 a :+ : T2 a  
           | a    :* : T2 a  
           | C2  
           deriving Read  
t2 :: T2 (T2 Int)  
t2 = read "(3 :* C2) :* C2"
```

The function *read* is a member of the class *Read*:

- ▶ *read* functions are elements in dictionaries



How Does the Problem Arise?

```
infix 5 :+:  
infix 6 :*:  
data T2 a = T2 a :+ : T2 a  
           | a    :* : T2 a  
           | C2  
           deriving Read  
t2 :: T2 (T2 Int)  
t2 = read "(3 :* C2) :* C2"
```

The function *read* is a member of the class *Read*:

- ▶ *read* functions are elements in dictionaries
- ▶ **instance**-declarations compose new dictionaries out of existing dictionaries at run-time



How Does the Problem Arise?

```
infix 5 :+:  
infix 6 :*:  
data T2 a = T2 a :+ : T2 a  
           | a    :* : T2 a  
           | C2  
           deriving Read  
t2 :: T2 (T2 Int)  
t2 = read "(3 :* C2) :* C2"
```

The function *read* is a member of the class *Read*:

- ▶ *read* functions are elements in dictionaries
- ▶ **instance**-declarations compose new dictionaries out of existing dictionaries at run-time
- ▶ hence *read* functions are to be composed at run-time



The Bad News

- ▶ Bottom-up parsers do not compose at all, and all perform an analysis of the complete grammar (YACC, Happy)
- ▶ Top-down parsers do not compose efficiently for arbitrary grammars, and may lead to left-recursive parsers if no care is taken:

```
data T1 a = a      : * : Int deriving Read
data T2    = T1 T2 : + : Int deriving Read
```



The Bad and the Good News

- ▶ Bottom-up parsers do not compose at all, and all perform an analysis of the complete grammar (YACC, Happy)
- ▶ Top-down parsers do not compose efficiently for arbitrary grammars, and may lead to left-recursive parsers if no care is taken:

```
data T1 a = a      : * : Int deriving Read
data T2    = T1 T2 : + : Int deriving Read
```

- ▶ Grammars can be composed!



Using Grammars instead of Parsers

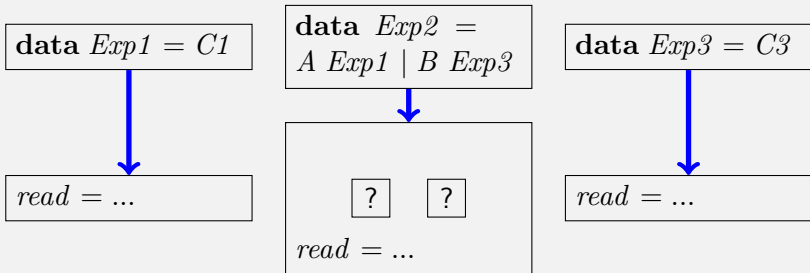
data $Exp1 = C1$

data $Exp2 =$
 $A Exp1 \mid B Exp3$

data $Exp3 = C3$



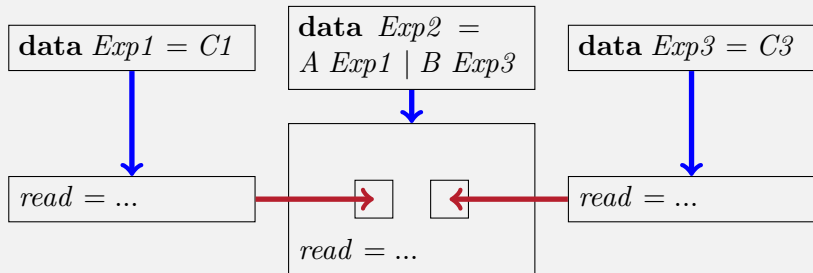
Using Grammars instead of Parsers



derive



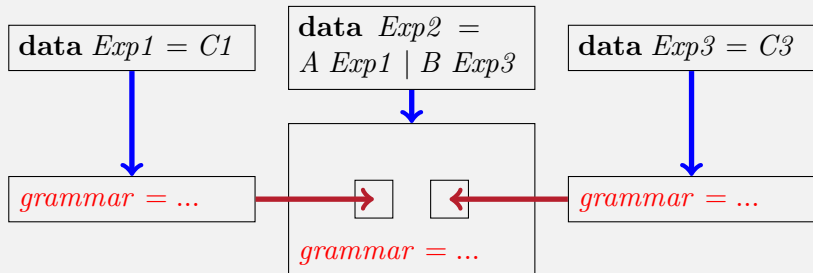
Using Grammars instead of Parsers



derive parameterise



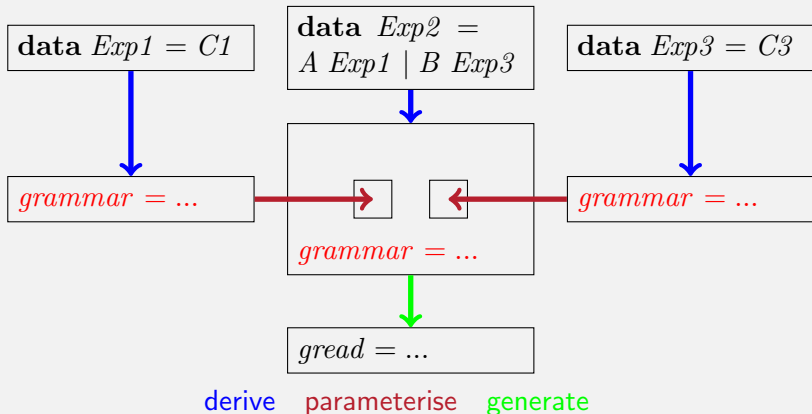
Using Grammars instead of Parsers



derive parameterise



Using Grammars instead of Parsers



The Class *Gram*

Instead of the class *Read* we introduce:

```
class Gram a where  
  grammar :: DGrammar a
```

where *DDGrammar* is a data type describing grammatical structures, including information about precedences.



The Class *Gram*

Instead of the class *Read* we introduce:

```
class Gram a where  
  grammar :: DGrammar a
```

where *DGrammar* is a data type describing grammatical structures, including information about precedences.

Note that it is labelled with a type *a*, which is the data type described by a value of type *DGrammar a*.



The Class *Gram*

Instead of the class *Read* we introduce:

```
class Gram a where  
  grammar :: DGrammar a
```

where *DGrammar* is a data type describing grammatical structures, including information about precedences.

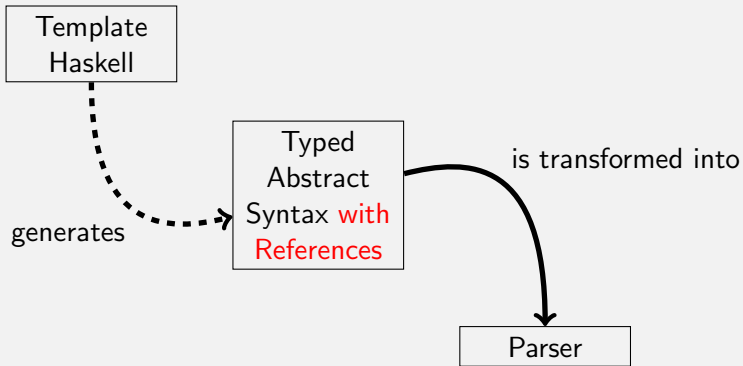
Note that it is labelled with a type *a*, which is the data type described by a value of type *DGrammar a*.

Now we can, just as for *read* define:

```
read  :: Read a ⇒ String → a  
gread :: Gram a ⇒ String → a
```



Generating parsers from Data Types



The Steps to be Taken

group Combine pieces of grammar together, introduce extra non-terminals to represent the precedences.



The Steps to be Taken

group Combine pieces of grammar together, introduce extra non-terminals to represent the precedences.

leftcorner Remove possible left-recursion by applying the *Left-Corner Transform*



The Steps to be Taken

group Combine pieces of grammar together, introduce extra non-terminals to represent the precedences.

leftcorner Remove possible left-recursion by applying the *Left-Corner Transform*

leftfactoring Combine common prefixes of alternatives



The Steps to be Taken

group Combine pieces of grammar together, introduce extra non-terminals to represent the precedences.

leftcorner Remove possible left-recursion by applying the *Left-Corner Transform*

leftfactoring Combine common prefixes of alternatives

compile Map the resulting *Grammar* onto a parser



The Steps to be Taken

group Combine pieces of grammar together, introduce extra non-terminals to represent the precedences.

leftcorner Remove possible left-recursion by applying the *Left-Corner Transform*

leftfactoring Combine common prefixes of alternatives

compile Map the resulting *Grammar* onto a parser

parse Use the parser to read a value



The Steps to be Taken

group Combine pieces of grammar together, introduce extra non-terminals to represent the precedences.

leftcorner Remove possible left-recursion by applying the *Left-Corner Transform*

leftfactoring Combine common prefixes of alternatives

compile Map the resulting *Grammar* onto a parser

parse Use the parser to read a value



The Steps to be Taken

group Combine pieces of grammar together, introduce extra non-terminals to represent the precedences.

leftcorner Remove possible left-recursion by applying the *Left-Corner Transform*

leftfactoring Combine common prefixes of alternatives

compile Map the resulting *Grammar* onto a parser

parse Use the parser to read a value

$$\begin{aligned} \text{gread} &:: \text{Gram } a \Rightarrow \text{String} \rightarrow a \\ \text{gread} &= (\text{parse} \circ \text{compile} \quad \circ \text{leftfactoring} \\ &\quad \circ \text{leftcorner} \circ \text{group}) \text{ grammar} \end{aligned}$$


Types Abstract Syntax with Explicit References

- ▶ Right hand sides of productions contain references to non-terminals



Types Abstract Syntax with Explicit References

- ▶ Right hand sides of productions contain references to non-terminals
- ▶ We want to be able to inspect and transform the grammar



Types Abstract Syntax with Explicit References

- ▶ Right hand sides of productions contain references to non-terminals
- ▶ We want to be able to inspect and transform the grammar
- ▶ We want to inspect the underlying graph structure



Types Abstract Syntax with Explicit References

- ▶ Right hand sides of productions contain references to non-terminals
- ▶ We want to be able to inspect and transform the grammar
- ▶ We want to inspect the underlying graph structure
- ▶ Of which the nodes are labelled with different types



Types Abstract Syntax with Explicit References

- ▶ Right hand sides of productions contain references to non-terminals
- ▶ We want to be able to inspect and transform the grammar
- ▶ We want to inspect the underlying graph structure
- ▶ Of which the nodes are labelled with different types
- ▶ So we use heterogeneous collections, i.e. we use nested cartesian products, henceforth called *Env*-ironments



References and Environments I

We introduce natural numbers, labelled with a type a describing what is referred to, and a list of types env describing the structure in which this a labelled object lives:

data $Ref\ a\ env$ **where**

$Zero :: Ref\ a\ (a, env)$

$Suc :: Ref\ a\ env' \rightarrow Ref\ a\ (x, env')$

data $Equal\ a\ b$ **where**

$Eq :: Equal\ a\ a$

$match :: Ref\ a\ env \rightarrow Ref\ b\ env \rightarrow Maybe\ (Equal\ a\ b)$

$match\ Zero\ Zero = Just\ Eq$

$match\ (Suc\ x)\ (Suc\ y) = match\ x\ y$

$match\ _ _ = Nothing$



References and Environments II

data $Env\ t\ use\ def$ **where**

$Empty :: Env\ t\ use\ ()$

$Cons :: t\ a\ use \rightarrow Env\ t\ use\ def'$
 $\rightarrow Env\ t\ use\ (a, def')$

$t\ a\ use ::$ a term of type t , describing a value of type a contains references pointing into an environment labelled by use . The parameter def describes the values actually existing in the Env . If use equals def the environment is closed.



data *DGrammar a*

$= \forall env. DGrammar (Ref\ a\ env)$
 $(Env\ DGram\ env\ env)$

data *DGram a env = DGD (DLNontDefs a env)*
 $| DGG (DGrammar\ a)$

newtype *DRef a env = DRef (Ref a env, Int)*

newtype *DLNontDefs a env*

$= DLNontDefs [(DRef\ a\ env, DProductions\ a\ env)]$



Continued ..

```
newtype DProductions a env  
  = DPS { unDPS :: [DProd a env] }
```

```
data DProd a env where  
  DSeq :: DSymbol b env → DProd (b → a) env  
                                               → DProd a env  
  DEnd :: a                               → DProd a env
```

```
data DSymbol a env where  
  DNont :: DRef a env → DSymbol a env  
  DTerm :: Token     → DSymbol Token env
```



Typed Abstract Syntax

```
data Exp = Exp :+: Exp
         | C
infixl 6 :+:
```

```
[
    _Exp      :+:      _Exp
,
    "C"
]
```



Typed Abstract Syntax

```
data Exp = Exp :+: Exp
         | C
infixl 6 :+:
```

```
[
    _Exp      :+:      _Exp
,
    "C"
,
    "("      _Exp
    ")"
]
```



Typed Abstract Syntax

```
data Exp = Exp :+ : Exp  
          | C  
infixl 6 :+ :
```

```
[  
    dNont (_Exp ) .#. dTerm " :+ : " .#.  
    dNont (_Exp )  
  
    ,  
    dTerm "C"  
    , dTerm "(" .#. dNont (_Exp ) .#.  
    dTerm ")"  
  
]
```



Typed Abstract Syntax

```
data Exp = Exp :+ : Exp
           | C
infixl 6 :+ :
```

```
[
  DPS [ dNont (_Exp, 6) .#. dTerm " :+ : " .#.
        dNont (_Exp, 7)
      ]
,
  DPS [ dTerm "C"
        , dTerm "(" .#. dNont (_Exp, 0) .#.
          dTerm ")"
      ]
]
```



Typed Abstract Syntax

```
data Exp = Exp :+ : Exp
          | C
infixl 6 :+ :
```

```
[ (DRef (_Exp, 6)
  , DPS [dNont (_Exp, 6) .#. dTerm ":+:" .#.
        dNont (_Exp, 7)
        ]
  )
, (DRef (_Exp, 10)
  , DPS [dTerm "C"
        , dTerm "(" .#. dNont (_Exp, 0) .#.
        dTerm ")"
        ]
  )
]
```



Typed Abstract Syntax

```
data Exp = Exp :+: Exp
          | C
infixl 6 :+:
```

```
[ (DRef (_Exp, 6)
  , DPS [dNont (_Exp, 6) .#. dTerm ":+:" .#.
         dNont (_Exp, 7) .#. dEnd plus]
  )
, (DRef (_Exp, 10)
  , DPS [dTerm "C" .#. dEnd (const C)
         , dTerm "(" .#. dNont (_Exp, 0) .#.
         dTerm ")" .#. dEnd parenT]
  )
]
plus e1 _ e2 = e2 :+: e1
```



Typed Abstract Syntax

instance *Gram Exp* **where**

grammar = *DGrammar* *_0 envExp*

envExp :: *Env DGram (Exp, ()) (Exp, ())*

envExp = *consD (nonts _0) Empty*

where

nonts *_Exp* = *DLNontDefs*

```
[ (DRef (_Exp, 6)
  , DPS [dNont (_Exp, 6) .#. dTerm ":+:" .#.
        dNont (_Exp, 7) .#. dEnd plus]
  )
  , (DRef (_Exp, 10)
    , DPS [dTerm "C" .#. dEnd (const C)
          , dTerm "(" .#. dNont (_Exp, 0) .#.
            dTerm ")" .#. dEnd parenT]
    )
  ]
```

plus *e1* *_ e2* = *e2* *:+:* *e1*



An Intermediate result

$$\begin{aligned}A &\rightarrow \text{"C1"} A_C1 \mid \text{"(" } A_(\ \\A_A &\rightarrow \text{":<:" } B A_A \mid \text{":<:" } B \\A_B &\rightarrow A_A \mid \epsilon \\A_C &\rightarrow \text{":>:" } B A_B \mid A_B \\A_C1 &\rightarrow A_C \\A_(\ &\rightarrow A \text{"})" } A_C \\B &\rightarrow \text{"C1"} B_C1 \mid \text{"(" } B_(\ \\B_C &\rightarrow \text{":>:" } B \mid \epsilon \\B_C1 &\rightarrow B_C \\B_(\ &\rightarrow A \text{"})" } B_C \\C &\rightarrow \text{"C1"} C_C1 \mid \text{"(" } C_(\ \\C_C1 &\rightarrow \epsilon \\C_(\ &\rightarrow A \text{"})" }\end{aligned}$$

1. We have introduced new non-terminals
2. Old non-terminals have new productions



The Transformations

All the transformations can be expressed in terms of an arrow-like type:

```
data Trafo m t a b =  
  Trafo ( $\forall env1.m env1 \rightarrow$   
     $\exists env2.$   
      ( $m env2$   
        ,  $\forall s. a s \rightarrow T env2 s \rightarrow Env t s env1 \rightarrow$   
          ( $b s, T env1 s, Env t s env2$ )  
        )  
      )
```



Results I

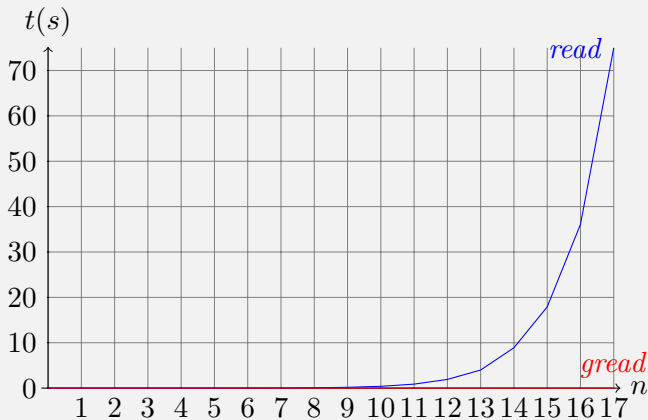
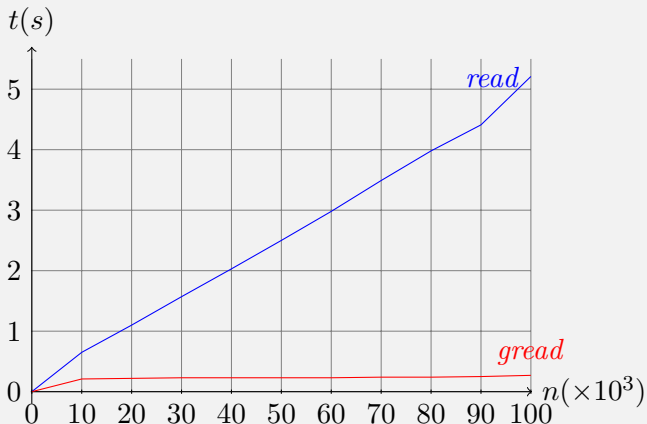


Figure: Execution times of reading $C \rightarrow (C \rightarrow \dots)$



Reading a Large Data Type



Overhead is very small, and that thanks to the use of the UU-parsers also parse times do hardly increase.



Why is this so complicated ...

1. The problem is complicated



Why is this so complicated ...

1. The problem is complicated
2. We do in 350 lines more than Bison (10.000 lines) is doing



Why is this so complicated ...

1. The problem is complicated
2. We do in 350 lines more than Bison (10.000 lines) is doing
3. Extra constructors are needed because we need existentials



Why is this so complicated ...

1. The problem is complicated
2. We do in 350 lines more than Bison (10.000 lines) is doing
3. Extra constructors are needed because we need existentials
4. If we have lazy evaluation, we also want it at the type level!



Why is this so complicated ...

1. The problem is complicated
2. We do in 350 lines more than Bison (10.000 lines) is doing
3. Extra constructors are needed because we need existentials
4. If we have lazy evaluation, we also want it at the type level!

```
f :: ∀a.(a → ∃b (b, a, b → b → Int))  
let (b, a, g) = f b  
in g b a
```



To Take Home

- ▶ The transformation library has been used unmodified for all the transformations
- ▶ The library can be used for any collection of abstract syntax trees, which contain references to each other, and of which the structure has to be inspected

