

Parser Memoisation

Marcelo Barbosa de Sousa January 24, 2011

Roadmap

Motivation

YAPL - Yet Another Parsing Library

Memoisation

Conclusion





Roadmap

Motivation

YAPL - Yet Another Parsing Library

Memoisation

Conclusion





Research Problem

Currently all parsing libraries will parse a piece of text twice with the same non-terminal if the grammar has various way of reaching the same point in the input. One approach to solving this problem is to use memoisation techniques. It is however not entirely clear how to do this magnificant typeful way.

Research Problem

Currently all parsing libraries will parse a piece of text twice with the same non-terminal if the grammar has various way of reaching the same point in the input. One approach to solving this problem is to use memoisation techniques. It is however not entirely clear how to do this in a online, typeful way.



Parser Combinators





Context free grammars

$$< N, \Sigma, P, S >$$

- ► N : Set of non-terminals
- $ightharpoonup \Sigma$: Set of terminals, e.g. *Char*
- ▶ P : Set of productions (A \in N, $\alpha \in V^*$)
- ightharpoonup S : Set of start symbols, S \in N
- ▶ $V = N \cup \Sigma$: Set of symbols
- $\triangleright N \cap \Sigma = \emptyset$





Parser Type

```
type Parser s \ a = [s] \rightarrow atype Parser s \ a = [s] \rightarrow (a, [s])type Parser s \ a = [s] \rightarrow [a]type Parser s \ a = [s] \rightarrow [(a, [s])]type Parser a = String \rightarrow [(a, String)]
```

Basic Parser Combinators Types

$$pSym :: Eq \ s \Rightarrow s \rightarrow Parser \ s \ s$$

$$pSym \ c = \lambda inp \rightarrow \mathbf{case} \ inp \ \mathbf{of}$$

$$(s:ss) \mid c \equiv s \rightarrow [(s,ss)]$$

$$otherwise \rightarrow []$$

 $pRet :: a \rightarrow Parser \ s \ a$ $pRet \ x = \lambda inp \rightarrow [(x, inp)]$

$$pSym :: Eq \ s \Rightarrow s \rightarrow Parser \ s \ pSym \ c = \lambda inp \rightarrow case \ inp \ of \ (s:ss) \mid c \equiv s \rightarrow [(s,ss)] \ otherwise \rightarrow []$$

$$pRet :: a \rightarrow Parser \ s \ a$$

 $pRet \ x = \lambda inp \rightarrow [(x, inp)]$

infixr
$$3 < | >$$

 $(< | >) :: Parser s a \rightarrow Parser s a \rightarrow Parser s a$
 $p < | > q = \lambda inp \rightarrow p inp + q inp$

```
infixl 5 < * > (< * >) :: Parser s (b \rightarrow a) \rightarrow Parser s b \rightarrow Parser s a <math>p < * > q = \lambda inp \rightarrow [(b2a \ b, rr) \mid (b2a, r) \leftarrow p \ inp, (b, rr) \leftarrow q \ r]
```

```
infixr 3 < | >

(< | >) :: Parser s a \rightarrow Parser s a \rightarrow Parser s a

p < | > q = \lambda inp \rightarrow p inp + q inp
```

infixl
$$5 < * >$$
 $(< * >) :: Parser s (b \rightarrow a) \rightarrow Parser s b \rightarrow Parser s a$
 $p < * > q = \lambda inp \rightarrow [(b2a \ b, rr) \mid (b2a, r) \leftarrow p \ inp,$
 $(b, rr) \leftarrow q \ r]$

Derived Parser Combinators

infix
$$7 < \$ >$$
 $(< \$ >) :: (b \rightarrow a) \rightarrow Parser \ s \ b \rightarrow Parser \ s \ a$
 $f < \$ > q = pRet \ f < * > q$

infixl 3 'opt' opt :: Parser $s \ a \rightarrow a \rightarrow Parser \ s \ a$ p 'opt' $v = p < | > pReturn \ v$

pMany :: Parser $s \ a \rightarrow$ Parser $s \ [a]$ pMany p = (:) < \$ > p < * > pMany p 'opt' []



Grammars are not always left factorized which results in inefficient parsing.

$$s$$
 :: Parser Char String $s = f < | > g$

$$f = fsem < \$ > pSym 'a' < * > p < * > q$$

$$g = h < * > p < * > r$$

 $h = gsem < $ > pSym$ 'a'

Grammars are not always left factorized which results in inefficient parsing.

$$s$$
 :: Parser Char String $s = f < | > g$

$$f = fsem < \$ > pSym$$
'a'< $* > p < * > q$

g = h < * > p < * > rh = gsem < \$ > pSym'a'

Grammars are not always left factorized which results in inefficient parsing.

$$s$$
 :: Parser Char String $s = f < | > g$

$$f = fsem < \$ > pSym$$
 'a' $< * > p < * > q$

$$g = h < * > p < * > r$$

 $h = gsem < $ > pSym$ 'a'

$$s$$
 :: Parser Char String $s = f < | > g$

$$f = fsem < \$ > pSym$$
'a' $< * > p < * > q$

$$g = h < * > p < * > r$$

 $h = gsem < \$ > pSym$ 'a'

The inefficiency results of our implementation of the choice combinator.

In the case of parser s, pSym'a' < * > p could be shared among both branches.



Roadmap

Motivation

YAPL - Yet Another Parsing Library

Memoisation

Conclusion

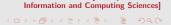




Goals

- ► Left-factorization for free
- ► Full control on parsing, online parsing, error control
- ▶ Same representation as previous parser combinators





Defunctionalizing Parsers

We can represent parsers with GADTs:

```
data Parser :: * \rightarrow *where\rightarrow Parser CharSym :: Char\rightarrow Parser CharRet :: a\rightarrow Parser aAlt :: Parser a\rightarrow Parser a \rightarrow Parser aSeq :: Parser (b \rightarrow a) \rightarrow Parser b \rightarrow Parser a
```

Continuations as a Stack of Parsers

We can represent the parsers yet to be analyzed as a stack.

```
data Pending :: * \rightarrow *where
Done :: Pending ()
Stack :: Parser a \rightarrow Pending b \rightarrow Pending (a,b)
```

- ▶ Done represents the end of our parser.
- Stack contains the current parser to be processed and the rest of the stack.

```
Ret (:[]) 'Seq' Sym 'a' \rightsquigarrow Stack (Ret (:[])) (Stack (Sym 'a') Done) :: Pending (a \rightarrow [a], (Char, ()))
```



Continuations as a Stack of Parsers

We can represent the parsers yet to be analyzed as a stack.

```
data Pending :: * \rightarrow * where
Done :: Pending ()
Stack :: Parser a \rightarrow Pending b \rightarrow Pending (a,b)
```

- ▶ Done represents the end of our parser.
- Stack contains the current parser to be processed and the rest of the stack.

```
Ret (:[]) 'Seq' Sym 'a' \rightarrow Stack (Ret (:[])) (Stack (Sym 'a') Done) :: Pending (a \rightarrow [a], (Char, ()))
```



Continuations as a Stack of Parsers

We can represent the parsers yet to be analyzed as a stack.

```
data Pending :: * \rightarrow *where
Done :: Pending ()
Stack :: Parser a \rightarrow Pending b \rightarrow Pending (a,b)
```

- ▶ Done represents the end of our parser.
- Stack contains the current parser to be processed and the rest of the stack.

```
Ret (:[]) 'Seq' Sym 'a' \rightsquigarrow Stack (Ret (:[])) (Stack (Sym 'a') Done) :: Pending (a \rightarrow [a], (Char, ()))
```

4 D > 4 A > 4 E > 4 E > E 9 Q P

data State
$$a = \forall b \circ State \ (b \rightarrow a) \ (Pending \ b)$$

type States $a = [State \ a]$

The use of existential types in *State* is crucial as we will see next.

Our interface function *runParser* simply builds the first state given a parser.

```
runParser :: Parser \ a \rightarrow [Char] \rightarrow [a]

runParser \ p = parse \ [State \ fst \ (Stack \ p \ Done)]
```

We can define the parse function iteratively:

```
parse :: States a \rightarrow [Char] \rightarrow [a]
parse states [] = evalStates states
parse states (x:xs) = (transition states x) 'parse' xs
```

For simplicity, I've consider [a] as the result of *evalStates*

```
evalStates :: States \ a \rightarrow [a]
evalStates \ [] = []
evalStates \ ((State \ f \ Done) : rest) = (f \ ()) : evalStates rest
evalStates \ (\_: rest) = error "parse error"
```

The parsing result type is flexible, and in *evalStates* we can do some error reporting although limited.



We can define the parse function iteratively:

```
egin{array}{ll} \textit{parse} :: \textit{States } a 
ightarrow [\textit{Char}] 
ightarrow [a] \\ \textit{parse states} [] &= \textit{evalStates states} \\ \textit{parse states} (x:xs) &= (\textit{transition states } x) '\textit{parse'} xs \end{array}
```

For simplicity, I've consider [a] as the result of evalStates

```
evalStates :: States \ a \rightarrow [a]
evalStates \ [] = []
evalStates \ ((State \ f \ Done) : rest) = (f \ ()) : evalStates \ rest
evalStates \ (\_: rest) = error "parse error"
```

The parsing result type is flexible, and in *evalStates* we can do some error reporting although limited.



```
transition :: States a \rightarrow Char \rightarrow States \ a

transition states char =

let unFoldNtStates = concatMap unFoldNtAtHeads states

in foldl (reduce char) [] unFoldNtStates
```

In each *transition* we unfold non-terminals, which are parsers and put them on the stack until we have on top of the stack a parser that consumes the current input position.

For simplicity we consider that the input is a String.





Reducing Terminals

```
reduce :: Char \rightarrow States a \rightarrow State a \rightarrow States a reduce char states (State f (Stack (Sym c) rest))

| c \equiv char = (State (\lambdarest \rightarrow f (char, rest)) rest) : states | otherwise = states reduce \_ states s@(State f Done) = s : states reduce \_ states \_ = states
```

At this point we are consuming the input character and generating a new list of states. We can introduce error reporting as well.

```
\begin{array}{l} \textit{unFoldNtAtHeads} :: \textit{State } a \rightarrow \textit{States } a \\ \textit{unFoldNtAtHeads } s@(\textit{State } r \ \textit{Done}) &= [s] \\ \textit{unFoldNtAtHeads } s@(\textit{State } r \ (\textit{Stack } (\textit{Sym } c) \ \textit{rest})) = [s] \end{array}
```

```
unFoldNtAtHeads (State r (Stack (Ret f) rest)) = unFoldNtAtHeads (State (\lambda rest \rightarrow r (f, rest)) rest)
```

```
unFoldNtAtHeads (State r (Stack (Seq p q) rest)) = unFoldNtAtHeads (State (\lambda(pr, (qr, rest)) \rightarrow r ((pr qr), rest) (Stack p (Stack q rest)))
```

```
unFoldNtAtHeads (State r (Stack (Alt p q) rest)) =

let statesp = unFoldNtAtHeads $ State r (Stack p rest

statesq = unFoldNtAtHeads $ State r (Stack q rest
```



```
unFoldNtAtHeads :: State \ a \rightarrow States \ a unFoldNtAtHeads \ s@(State \ r \ Done) = [s] unFoldNtAtHeads \ s@(State \ r \ (Stack \ (Sym \ c) \ rest)) = [s]
```

```
unFoldNtAtHeads (State r (Stack (Ret f) rest)) = unFoldNtAtHeads (State (\lambda rest \rightarrow r (f, rest)) rest)
```

```
unFoldNtAtHeads (State r (Stack (Seq p q) rest)) = unFoldNtAtHeads (State (\lambda(pr, (qr, rest)) \rightarrow r ((pr qr), rest) (Stack p (Stack q rest)))
```

unFoldNtAtHeads (State r (Stack (Alt p q) rest)) =
let statesp = unFoldNtAtHeads \$ State r (Stack p rest)
statesq = unFoldNtAtHeads \$ State r (Stack q rest)



4 D > 4 A > 4 E > 4 E > E 9 Q P

```
unFoldNtAtHeads :: State \ a \rightarrow States \ a unFoldNtAtHeads \ s@(State \ r \ Done) = [s] unFoldNtAtHeads \ s@(State \ r \ (Stack \ (Sym \ c) \ rest)) = [s]
```

```
unFoldNtAtHeads\ (State\ r\ (Stack\ (Ret\ f)\ rest)) = unFoldNtAtHeads\ (State\ (\lambda rest\ 	o r\ (f,rest))\ rest)
```

```
\begin{array}{l} \textit{unFoldNtAtHeads} \; (\textit{State} \; r \; (\textit{Stack} \; (\textit{Seq} \; p \; q) \; \textit{rest})) = \\ \textit{unFoldNtAtHeads} \; (\textit{State} \; (\lambda(\textit{pr}, (\textit{qr}, \textit{rest})) \rightarrow r \; ((\textit{pr} \; \textit{qr}), \textit{rest})) \\ & (\textit{Stack} \; p \; (\textit{Stack} \; q \; \textit{rest}))) \end{array}
```

unFoldNtAtHeads (State r (Stack (Alt p q) rest)) =

let statesp = unFoldNtAtHeads \$ State r (Stack p rest)

statesq = unFoldNtAtHeads \$ State r (Stack q rest)



```
unFoldNtAtHeads :: State \ a \rightarrow States \ a

unFoldNtAtHeads \ s@(State \ r \ Done) = [s]

unFoldNtAtHeads \ s@(State \ r \ (Stack \ (Sym \ c) \ rest)) = [s]
```

```
unFoldNtAtHeads (State r (Stack (Ret f) rest)) = unFoldNtAtHeads (State (\lambda rest \rightarrow r (f, rest)) rest)
```

```
unFoldNtAtHeads\ (State\ r\ (Stack\ (Seq\ p\ q)\ rest)) = \\ unFoldNtAtHeads\ (State\ (\lambda(pr,(qr,rest)) \to r\ ((pr\ qr),rest)) \\ (Stack\ p\ (Stack\ q\ rest)))
```

```
unFoldNtAtHeads (State r (Stack (Alt p q) rest)) =
  let statesp = unFoldNtAtHeads $ State r (Stack p rest)
     statesq = unFoldNtAtHeads \$ State r (Stack q rest)
in statesp ++ statesa
```



Faculty of Science Information and Computing Sciences

Running example

```
pChar :: Char \rightarrow Parser String
pChar c = Ret(:[]) 'Seq' Sym c
```

```
> runParser (pChar 'a') "a"
["a"]
```

```
> runParser (pChar 'a') "ab"
["a"]
```

```
> runParser (pChar 'a' 'Alt' pChar 'b') "a"
["a"]
> runParser (pChar 'a' 'Alt' pChar 'b') "b"
["b"]
```



Running example

```
pChar :: Char \rightarrow Parser String

pChar c = Ret (:[]) 'Seq' Sym c
```

```
> runParser (pChar 'a') "a"
["a"]
```

```
> runParser (pChar 'a') "ab"
["a"]
```

```
> runParser (pChar 'a' 'Alt' pChar 'b') "a"
["a"]
> runParser (pChar 'a' 'Alt' pChar 'b') "b"
["b"]
```



Running example

```
pChar :: Char \rightarrow Parser String

pChar c = Ret (:[]) 'Seq' Sym c
```

```
> runParser (pChar 'a') "a"
["a"]
```

```
> runParser (pChar 'a') "ab"
["a"]
```

```
> runParser (pChar 'a' 'Alt' pChar 'b') "a"
["a"]
> runParser (pChar 'a' 'Alt' pChar 'b') "b"
["b"]
```



```
pChar :: Char \rightarrow Parser String
pChar c = Ret(:[]) 'Seq' Sym c
```

```
> runParser (pChar 'a') "a"
["a"]
```

```
> runParser (pChar 'a') "ab"
["a"]
```

```
> runParser (pChar 'a' 'Alt' pChar 'b') "a"
["a"]
> runParser (pChar 'a' 'Alt' pChar 'b') "b"
["b"]
```



```
> runParser (pChar 'a') "b"
[]
```

```
> runParser (pChar 'a') ""
*** Exception: parse error
```

```
> runParser (Ret "a") ""
*** Exception: parse error
```



```
> runParser (pChar 'a') "b"
[]
```

```
> runParser (pChar 'a') ""
*** Exception: parse error
```

```
> runParser (Ret "a") ""
*** Exception: parse error
```





```
> runParser (pChar 'a') "b"
[]
```

```
> runParser (pChar 'a') ""
*** Exception: parse error
```

```
> runParser (Ret "a") ""
*** Exception: parse error
```





The Ret Problem

Since the Ret combinator does not require any input to succeed we need an extra step on evaluation:

```
evalStates ((State f (Stack (Ret x) Done)) : rest) = (f (x, ())) : evalStates rest
```

```
> runParser (Ret "a") ""
["a"]
```

```
> runParser (Alt (Ret "b") (Ret "a")) ""
["b","a"]
```



The Ret Problem

Since the Ret combinator does not require any input to succeed we need an extra step on evaluation:

```
evalStates ((State f (Stack (Ret x) Done)) : rest) = (f (x, ())) : evalStates rest
```

```
> runParser (Ret "a") ""
["a"]
```

```
> runParser (Alt (Ret "b") (Ret "a")) ""
["b", "a"]
```

It's easy to extend the library with other simple combinators such as pSatisfay.

We need to extend our Parser GADT

$$Sat :: (Char \rightarrow Bool) \rightarrow Parser\ Char$$

Since Sat consumes input we do not unfold anything

$$unFoldNtAtHeads\ s@(State\ r\ (Stack\ (Sat\ f)\ rest)) = [s]$$

Finally we need to create a new state in the reduce step

```
reduce char states (State f (Stack (Sat s) rest))

| s \ char = (State \ (\lambda rest \rightarrow f \ (char, rest)) \ rest) : states

| \ otherwise = states
```



It's easy to extend the library with other simple combinators such as pSatisfay.

We need to extend our Parser GADT

 $Sat :: (Char \rightarrow Bool) \rightarrow Parser Char$

Since Sat consumes input we do not unfold anything

 $unFoldNtAtHeads\ s@(State\ r\ (Stack\ (Sat\ f)\ rest)) = [s]$

Finally we need to create a new state in the reduce step

reduce char states (State f (Stack (Sat s) rest)) $| s \ char = (State \ (\lambda rest \rightarrow f \ (char, rest)) \ rest) : states$ $| \ otherwise = states$



It's easy to extend the library with other simple combinators such as pSatisfay.

We need to extend our Parser GADT

$$Sat :: (Char \rightarrow Bool) \rightarrow Parser Char$$

Since Sat consumes input we do not unfold anything

$$unFoldNtAtHeads\ s@(State\ r\ (Stack\ (Sat\ f)\ rest)) = [s]$$

Finally we need to create a new state in the reduce step

```
reduce char states (State f (Stack (Sat s) rest))

\mid s char = (State (\lambda rest \rightarrow f (char, rest)) rest) : states

\mid otherwise = states
```

```
pDigit :: Parser Char

pDigit = Sat (\lambda x \rightarrow `0` \leqslant x \land x \leqslant `9`)

pDigitAsInt :: Parser Int

pDigitAsInt = Ret (\lambda c \rightarrow fromEnum \ c - fromEnum \ `0`)

'Seq' pDigit
```

```
> runParser pDigitAsInt "1"
[1]
```

```
> runParser pDigitAsInt "92"
[9]
```



```
pDigit :: Parser Char

pDigit = Sat (\lambda x \rightarrow `0` \leqslant x \land x \leqslant `9`)

pDigitAsInt :: Parser Int

pDigitAsInt = Ret (\lambda c \rightarrow fromEnum \ c - fromEnum \ `0`)

'Seq' pDigit
```

```
> runParser pDigitAsInt "1"
[1]
```

> runParser pDigitAsInt "92"
[9]



We can represent our initial parsers combinators with our new library

```
pReturn :: a \rightarrow Parser a

pReturn = Ret

pSym :: Char \rightarrow Parser Char

pSym = Sym
```

$$(<|>)$$
:: Parser $a \rightarrow$ Parser $a \rightarrow$ Parser $a \rightarrow$ Parser $a \rightarrow$ $p < |> q = p$ 'Alt' q

(<*>):: Parser $(b \rightarrow a) \rightarrow$ Parser $b \rightarrow$ Parser $a \rightarrow b \rightarrow$ Parser $b \rightarrow$ Parser $a \rightarrow b \rightarrow$ Parser $b \rightarrow$ Parser $b \rightarrow$ Parser $a \rightarrow$ Parser $b \rightarrow$ Parser



We can represent our initial parsers combinators with our new library

```
pReturn :: a \rightarrow Parser a

pReturn = Ret

pSym :: Char \rightarrow Parser Char

pSym = Sym
```

$$(<|>) :: Parser a \rightarrow Parser a \rightarrow Parser a$$

 $p < |> q = p$ 'Alt' q

(<*>):: Parser $(b \rightarrow a) \rightarrow$ Parser $b \rightarrow$ Parser $a \rightarrow b \rightarrow$ Parser $a \rightarrow$



We can represent our initial parsers combinators with our new library

```
pReturn :: a \rightarrow Parser a

pReturn = Ret

pSym :: Char \rightarrow Parser Char

pSym = Sym
```

$$(<|>)$$
:: Parser $a \rightarrow$ Parser

$$(<*>)$$
:: Parser $(b \rightarrow a) \rightarrow$ Parser $b \rightarrow$ Parser $a p < * > q = p$ 'Seq' q



```
pManyA :: Parser String
pManyA = pMany (pSym `a`)
pManyAb :: Parser String
pManyAb = (\lambda as b \rightarrow as ++ [b])
< \$ > pManyA < * > pSym `b`
```

```
> runParser pManyA ""
[""]
```

```
> runParser pManyA "a"
["","a"]
```



```
pManyA :: Parser String
pManyA = pMany (pSym 'a')
pManyAb :: Parser String
pManyAb = (\lambda as b \rightarrow as + [b])
< \$ > pManyA < * > pSym 'b'
```

```
> runParser pManyA ""
[""]
```

```
> runParser pManyA "a"
["","a"]
```

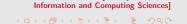


```
> runParser pManyA "aaaa"
["a","aaa","aaaa","aa",""]
```

```
> runParser pManyAb "aaaa"
*** Exception: parse error
```

```
> runParser pManyAb "aaaab"
["aaaab"]
```





State of art of our library

- suited for left-factorization
- easily extensible
- online parsing
- easy implementable control error
- easy to use





Roadmap

Motivation

YAPL - Yet Another Parsing Library

Memoisation

Conclusion





The Equality Problem





Equality on Parsers

Since we have a clear representation of parsers we define equality parsers in a natural way.

instance
$$Eq \ a \Rightarrow Eq \ (Parser \ a)$$
 where
 $Sym \ c \equiv Sym \ d = c \equiv d$
 $Ret \ a \equiv Ret \ b = a \equiv b$
 $Alt \ p \ q \equiv Alt \ r \ s = p \equiv r \land q \equiv s$
 $Seq \ f \ p \equiv Seq \ g \ q = f \equiv g \land p \equiv q$

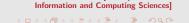
Unfortunately this won't type check, because in our *Parser* representation we loose the type of the type variable *b* in *Seq*.

Recursive descent parsing

Even if it type checked the following evaluation would not terminate.

```
pManyA :: Parser String
pManyA = pMany (pSym 'a')
```

Having recursive parsers difficulties equality.



Stable Pointers

A stable pointer is a reference to a Haskell expression that is guaranteed not to be affected by garbage collection, i.e., it will neither be deallocated nor will the value of the stable pointer itself change during garbage collection (ordinary references may be relocated during garbage collection). Consequently, stable pointers can be passed to foreign code, which can treat it as an opaque reference to a Haskell value.

Equality on Parsers with Stable Pointers

There is a predefined Eq instance for (StablePtr a) which helps with the equality of stable pointers, but unfortunately stable pointers live in the IO world. Furthermore the comparison of parsers in different states requires us to help the compiler to recognize that two parsers have the same type.

```
data Equal :: * \rightarrow * \rightarrow *where Eq :: Equal a a
```

```
eq:: Parser a \rightarrow Parser \ b \rightarrow Maybe \ (Equal \ a \ b)
eq p \ q = \mathbf{let} \ pptr = unsafePerformIO \$ \ castParser \ p
qptr = unsafePerformIO \$ \ newStablePtr \ q
\mathbf{in} \ \mathbf{if} \ pptr \equiv qptr
\mathbf{then} \ Just \$ \ unsafeCoerce \ Eq
\mathbf{else} \ Nothing
```



Equality on Parsers with Stable Pointers

There is a predefined Eq instance for (StablePtr a) which helps with the equality of stable pointers, but unfortunately stable pointers live in the IO world. Furthermore the comparison of parsers in different states requires us to help the compiler to recognize that two parsers have the same type.

```
data Equal :: * \rightarrow * \rightarrow *where Eq :: Equal a a
```

```
eq:: Parser a \rightarrow Parser b \rightarrow Maybe (Equal a b)
eq p q = let pptr = unsafePerformIO $ castParser <math>p
qptr = unsafePerformIO $ newStablePtr <math>q
in if pptr \equiv qptr
then Just $ unsafeCoerce Eq
else Nothing
```



[Faculty of Science Information and Computing Sciences]

4 D > 4 A > 4 E > 4 E > E 9 Q P

Equality on Parsers with Stable Pointers

We can now define the Eq instance for Parser.

```
instance Eq (Parser a) where p \equiv q = \mathbf{case} \ p'eq' \ q \ \mathbf{of}
Nothing \rightarrow False
Just Eq \rightarrow True
```

We want to be able to merge states that have common parsers in their stack.

First we modify our transition function to be able to unfold and merge states

```
transition :: States \ a 	o Char 	o States \ a transition \ states \ char = let \ unFoldNtStates = unFoldAndMerge \ states in \ foldl \ (reduce \ char) \ [\ ] \ unFoldNtStates
```





```
unFoldAndMerge :: States a \rightarrow States a
unFoldAndMerge states
| termination states = states
      otherwise = unFoldAndMerge $ unFoldNtAtHeads states
termination :: States a \rightarrow Bool
termination xs = all isFolded xs
isFolded :: State \ a \rightarrow Bool
isFolded (State r Done) = True isFolded (State r (Stack (Sym c) rest)) = True isFolded _ = False
```



To be able to merge we need to add new constructors to our Pending structure

Split :: Pending $r1 \rightarrow$ Pending $r2 \rightarrow$ Pending (Either $r1 \ r2$)

and generalize our semantic function

```
data State a = \forall b \circ State (Func b a) (Pending b)
data Func b a where
Single :: (b \rightarrow a) \qquad \rightarrow Func \ b \ a
Two :: (b \rightarrow a) \rightarrow (b \rightarrow a) \rightarrow Func \ b \ a
```





To be able to merge we need to add new constructors to our Pending structure

Split :: Pending $r1 \rightarrow$ Pending $r2 \rightarrow$ Pending (Either $r1 \ r2$)

and generalize our semantic function

data State $a = \forall b \circ State$ (Func b a) (Pending b)

data Func
$$b$$
 a **where**

$$Single :: (b \rightarrow a) \longrightarrow Func \ b \ a$$

$$Two :: (b \rightarrow a) \rightarrow (b \rightarrow a) \rightarrow Func \ b \ a$$

$$(b \rightarrow a) \rightarrow (b \rightarrow a) \rightarrow Func \ b$$

We need to modify unFoldNtAtHeads to implement sharing

```
unFoldNtAtHeads :: States \ a \rightarrow States \ a
unFoldNtAtHeads \ [\ ] = [\ ]
```

```
unFoldNtAtHeads (s@(State func (Stack (Ret f) rest)) : states) =
let ns = unFoldRetAtHead s
in unFoldNtAtHeads $ ns : states
```

```
unFoldRetAtHead :: State a \rightarrow State a

unFoldRetAtHead (State (Single r) (Stack (Ret f) rest)) = State (Single (\lambdarest \rightarrow r (f, rest))) rest

unFoldRetAtHead (State (Two f1 f2) (Stack (Ret f) rest)) = State (Two (\lambdar1 \rightarrow f1 (f, r1)) (\lambdar2 \rightarrow f2 (f, r2))) rest
```



We need to modify unFoldNtAtHeads to implement sharing

```
unFoldNtAtHeads :: States \ a \rightarrow States \ a unFoldNtAtHeads \ [\ ] = [\ ]
```

```
unFoldNtAtHeads (s@(State func (Stack (Ret f) rest)): states) = let \ ns = unFoldRetAtHeads s in \ unFoldNtAtHeads s s: states
```

```
unFoldRetAtHead :: State a \rightarrow State a

unFoldRetAtHead (State (Single r) (Stack (Ret f) rest)) = State (Single (\lambda rest \rightarrow r (f, rest))) rest

unFoldRetAtHead (State (Two f1 f2) (Stack (Ret f) rest)) = State (Two (\lambda r1 \rightarrow f1 (f, r1)) (\lambda r2 \rightarrow f2 (f, r2))) rest
```



We need to modify unFoldNtAtHeads to implement sharing

```
unFoldNtAtHeads :: States \ a \rightarrow States \ a unFoldNtAtHeads \ [\ ] = [\ ]
```

```
unFoldNtAtHeads (s@(State func (Stack (Ret f) rest)): states) = let ns = unFoldRetAtHeads s in unFoldNtAtHeads $ ns: states
```

```
unFoldRetAtHead :: State a \rightarrow State \ a

unFoldRetAtHead (State (Single r) (Stack (Ret f) rest)) = State (Single (\lambda rest \rightarrow r (f, rest))) rest

unFoldRetAtHead (State (Two f1\ f2) (Stack (Ret f) rest)) = State (Two (\lambda r1 \rightarrow f1\ (f,r1)) (\lambda r2 \rightarrow f2\ (f,r2))) rest
```



4 D > 4 A > 4 E > 4 E > E 9 Q P

We can now share equal choices

```
unFoldNtAtHeads (s@(State func (Stack (Alt p q) rest)): states):
  let optState = unFoldAltAtHead s
  in optState + unFoldNtAtHeads states
unFoldAlt :: Parser a \rightarrow [Parser a]
unFoldAlt (Alt p q) = unFoldAlt p + unFoldAlt q
unFoldAlt p
unFoldAltAtHead :: State a \rightarrow States a
unFoldAltAtHead (State r (Stack altParser rest)) =
  let lstAltParsers = unFoldAlt altParser
      uniqueAltParsers = nub\ lstAltParsers
  in map (\lambda p \rightarrow State\ r\ (Stack\ p\ rest)) uniqueAltParsers
```

The interesting case and real value of sharing occurs when we sequence

```
unFoldNtAtHeads (s@(State func (Stack (Seq p q) rest)) : states) :
let (optState, nst) = unFoldSeqAtHead s states
    optStates = unFoldNtAtHeads nst
in optState : optStates
```



The interesting case and real value of sharing occurs when we sequence

```
unFoldSeqAtHead :: State a \rightarrow States a \rightarrow (State a, States a)
unFoldSeqAtHead\ v@(State\ (Single\ r)\ (Stack\ (Seq\ p\ q)\ r1))
   (k@(State\ (Single\ s)\ (Stack\ (Seq\ t\ u)\ r2)):states) =
  case p'eq' t of
     Just Eq \rightarrow let opts = State
                          (Two (\lambda(pr, Left (qr, r1)) \rightarrow r (pr qr, r1))
                                 (\lambda(pr, Right(qu, r2)) \rightarrow s(pr qu, r2))
                          (Stack p (Split (Stack q r1) (Stack u r2)))
                    in (opts, states)
     Nothing \rightarrow let (nst, nsts) = unFoldSeqAtHead v states
                    in (nst, k : nsts)
```

Once we encounter Split we split the current state

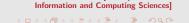
```
unFoldNtAtHeads (s@(State (Two fl fr) (Split pl pr)): states) = 

\mathbf{let} sleft = State (Single (\lambda r1 \rightarrow fl (Left r1))) pl

sright = State (Single (\lambda r2 \rightarrow fr (Right r2))) pr

\mathbf{in} sleft: sright: states
```





```
[Stack s Done]

ightharpoonup [Stack (f < * > Sym 'y') Done, Stack (f < * > Sym 'z') Done]

<math>
ightharpoonup [Stack f (Split (Stack (Sym 'y') Done) (Stack (Sym 'z') Done))]

<math>
ightharpoonup [Split (Stack (Sym 'y') Done) (Stack (Sym 'z') Done))]

<math>
ightharpoonup [Stack (Sym 'y') Done, Stack (Sym 'z') Done]

<math>
ightharpoonup [Done]
```



```
unFoldSegAtHead :: State a \rightarrow States \ a \rightarrow (State \ a, States \ a)
unFoldSeqAtHead v@(State (Single r) (Stack (Seq p q) r1))
   (k@(State\ (Single\ s)\ (Stack\ (Seq\ t\ u)\ r2)):states) =
  case p'eq' t of
     Just Eq \rightarrow let opts = State
                          (Two (\lambda(pr, Left (qr, r1)) \rightarrow r (pr qr, r1))
                                 (\lambda(pr,Right(qu,r2)) \rightarrow s(pr qu,r2))
                          (Stack p (Split (Stack q r1) (Stack u r2)))
                    in (opts, states)
     Nothing \rightarrow let (nst, nsts) = unFoldSeqAtHead v states
                    in (nst, k : nsts)
```

Here we can also check for the equality of $\it q$ and $\it u$ and create a state where we represent future sharing between that state and



◆□▶◆御▶◆壹▶◆壹▶ 壹 夕久◎

Extend Pending again

LazyShare :: Parser p o Pending r o Pending (Share p r)

We need to semantically distinguish between *Stack* and *LazyShare*

data *Share* p r **where** *Share* p r \rightarrow *Share* p r

$$xOrz = Sat ((\equiv 'x') \lor (\equiv 'z'))$$

 $s = f < \$ > xOrz < * > q < * > r$
 $< | >$
 $g < \$ > Sym 'x' < * > q < * > s$

```
[Stack s Done]
> [Stack (f <$> xOrz <*> q <*> r) Done,
   Stack (g <$> Sym 'x' <*> q <*> s) Done]
> [Stack (f <$> xOrz <*> q) (Stack r Done),
   Stack (g <$> Sym 'x' <*> q) (Stack s Done)]
> [Stack (f <$> xOrz) (Sharing q (Stack r Done)),
   Stack (g <$> Sym 'x') (Sharing q (Stack s Done))]
> [Stack f (Stack xOrz (Sharing q (Stack r Done))),
   Stack g (Stack (Sym 'x') (Sharing q (Stack s Done)
> [Stack xOrz (Sharing q (Stack r Done)),
   Stack (Sym 'x') (Sharing q (Stack s Done))]
> [Sharing q (Stack r Done), Sharing q (Stack s Done)
> [Stack q (Split (Stack r Done) (Stack s Done))]
```

> [Split (Stack r Done) (Stack s Done)]

> [Stack r Done, Stack s Done]





Generalization

So far we have achieve sharing between two branches of choice, however generalizing it to N branches should be easily done as long as we can map a function in *Func* to a *Pending*.

Split :: Pending $r1 \rightarrow$ Pending $r2 \rightarrow$ Pending (Pair $r1 \ r2$)

data Pair r1 r2 **where** Pair :: $r1 \rightarrow r2 \rightarrow Pair r1 r2$

Our new State representation could be

data State $a = \forall b \circ State [(b \rightarrow a)]$ (Pending b)

Roadmap

Motivation

YAPL - Yet Another Parsing Library

Memoisation

Conclusion





Conclusion

Advantages

- Sharing
- Natural languages grammars
- Full parsing control (online parsing, error reporting)
- Easy to use
- Almost well typed

Disadvantages

- Inefficiency
- Parser equality
- Some combinators might be difficult to implement efficiently



Conclusion

Advantages

- Sharing
- Natural languages grammars
- Full parsing control (online parsing, error reporting)
- Easy to use
- ► Almost well typed

Disadvantages

- ► Inefficiency
- ► Parser equality
- Some combinators might be difficult to implement efficiently



Future work

- Generalization for N branches
- Results
- ► Try other approaches to parser equality
- ► Left recursive grammars
- Monadic parsers



