



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]

Advanced Functional Programming

2012-2013, periode 2

Doaitse Swierstra

Department of Information and Computing Sciences
Utrecht University

December 11, 2012

10. Advanced parser Combinators



10.1 Problems with “List of Successes”



Recap: parser Combinators

The naïve implementation of parser combinators uses the list-of-successes method, which is a combination of a state monad and a list monad:

$$\begin{aligned} \langle * \rangle &:: \text{Parser } (b \rightarrow a) \rightarrow \text{Parser } b \rightarrow \text{Parser } a \\ p \langle * \rangle q &= \lambda \text{inp} \rightarrow [(b2a \ b, \text{qrest}) \mid (b2a, \text{prest}) \leftarrow p \ \text{inp} \\ &\quad , (b, \text{qrest}) \leftarrow q \ \text{prest} \\ &\quad] \\ \langle || \rangle &:: \text{Parser } a \rightarrow \text{Parser } a \rightarrow \text{Parser } a \\ p \langle || \rangle q &= \lambda \text{inp} \rightarrow p \ \text{inp} \text{ ++ } q \ \text{inp} \end{aligned}$$


Problems with Erroneous Input

- ▶ If your input does not conform to the language recognized by the parser all you get as a result is: `[]`.



Problems with Erroneous Input

- ▶ If your input does not conform to the language recognized by the parser all you get as a result is: [].
- ▶ It may take quite a while before you get this negative result, since the backtracking may try all other alternatives at all positions.



Problems with Erroneous Input

- ▶ If your input does not conform to the language recognized by the parser all you get as a result is: `[]`.
- ▶ It may take quite a while before you get this negative result, since the backtracking may try all other alternatives at all positions.
- ▶ There is no indication of where things went wrong.



Problems with Space Consumption

- ▶ A **complete result** has to be constructed before any part of it is returned
- ▶ The **complete input** is present in memory as long as no parse has been found
- ▶ Efficiency may depend critically on the ordering of the alternatives, and thus on how the grammar was written

For all of these problems we have found solutions.



10.2 History Parsers



Replace depth-first by breath-first

We introduce a Steps data type which contains a computed result (using a GADT and an existential type, to which we will come back later).

data Steps a **where**

Step :: Progress \rightarrow Steps a \rightarrow Steps a

Apply :: $\forall a b. (b \rightarrow a) \rightarrow$ Steps b \rightarrow Steps a

Fail :: ...



Replace depth-first by breath-first

We introduce a Steps data type which contains a computed result (using a GADT and an existential type, to which we will come back later).

data Steps a **where**

Step :: Progress \rightarrow Steps a \rightarrow Steps a

Apply :: $\forall a b. (b \rightarrow a) \rightarrow$ Steps b \rightarrow Steps a

Fail :: ...

The Progress field describes how much progress we made in the input (i.e. how much of the input was consumed by this step)



Computing a result

We compute a result on the fly, and change the parser type into a “continuation monad”:

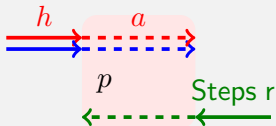
newtype HP st a
= HP ($\forall r. (a \rightarrow st \rightarrow \text{Steps } r) \rightarrow st \rightarrow \text{Steps } r$)



Computing a result

We compute a result on the fly, and change the parser type into a “continuation monad”:

newtype HP st a
= HP ($\forall r. (a \rightarrow st \rightarrow \text{Steps } r) \rightarrow st \rightarrow \text{Steps } r$)



We define the function which compares two alternatives.

$$\text{best}' :: \text{Steps } b \rightarrow \text{Steps } b \rightarrow \text{Steps } b$$

$$\text{Fail } \dots \text{ 'best}' \quad \dots \quad = \text{Fail } \dots$$

$$\text{Fail } \dots \text{ 'best}' \quad r \quad = r$$

$$l \quad \text{'best}' \quad \text{Fail } \dots = l$$

$$\text{Step } n \ l \ \text{'best}' \ \text{Step } m \ r$$

$$| \ n = m = \text{Step } n \ (l \ \text{'best}' \ r)$$

$$| \ n < m = \text{Step } n \ (l \ \text{'best}' \ \text{Step } (m - n) \ r)$$

$$| \ n > m = \text{Step } m \ (\text{Step } (n - m) \ l \ \text{'best}' \ r)$$


History parsers are Functor and Applicative

instance Functor (T st) **where**

fmap f (HP ph) = HP ($\lambda k \rightarrow \text{ph } (k \circ f)$)

instance Applicative (HP state) **where**

HP ph $\langle * \rangle$ ~ (HP qh)

= HP ($\lambda k \rightarrow \text{ph } (\lambda pr \rightarrow \text{qh } (\lambda qr \rightarrow k (pr qr))))$)

pure a = HP ($\$a$)

instance Alternative (T state) **where**

HP ph $\langle | \rangle$ HP qh = HP ($\lambda k \text{ inp} \rightarrow \text{ph } k \text{ inp}$ ‘best’ qh k inp)

empty = HP ($\lambda k \text{ inp} \rightarrow \text{noAlts}$)



History parsers are Functor and Applicative

instance Functor (T st) **where**

$$\text{fmap } f \text{ (HP ph)} = \text{HP } (\lambda k \rightarrow \text{ph } (k \circ f))$$

instance Applicative (HP state) **where**

$$\text{HP ph } \langle * \rangle \sim (\text{HP qh})$$

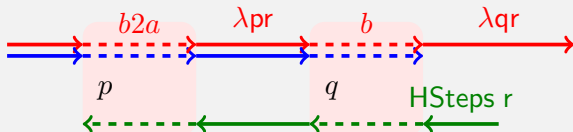
$$= \text{HP } (\lambda k \rightarrow \text{ph } (\lambda pr \rightarrow \text{qh } (\lambda qr \rightarrow k \text{ (pr qr)})))$$

$$\text{pure } a = \text{HP } (\$a)$$

instance Alternative (T state) **where**

$$\text{HP ph } \langle | \rangle \text{ HP qh} = \text{HP } (\lambda k \text{ inp} \rightarrow \text{ph } k \text{ inp 'best' qh } k \text{ inp})$$

$$\text{empty} = \text{HP } (\lambda k \text{ inp} \rightarrow \text{noAlts})$$



History parsers are Functor and Applicative

instance Functor (T st) **where**

$\text{fmap } f \text{ (HP ph)} = \text{HP } (\lambda k \rightarrow \text{ph } (k \circ f))$

instance Applicative (HP state) **where**

$\text{HP ph } \langle * \rangle \sim (\text{HP } qh)$

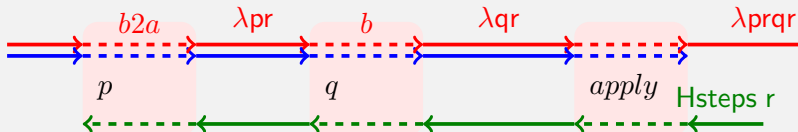
$= \text{HP } (\lambda k \rightarrow \text{ph } (\lambda pr \rightarrow qh (\lambda qr \rightarrow k (pr qr))))$

$\text{pure } a = \text{HP } (\$a)$

instance Alternative (T state) **where**

$\text{HP ph } \langle | \rangle \text{ HP } qh = \text{HP } (\lambda k \text{ inp} \rightarrow \text{ph } k \text{ inp 'best' } qh k \text{ inp})$

$\text{empty} = \text{HP } (\lambda k \text{ inp} \rightarrow \text{noAlts})$



10.3 Online Result Construction



Online results

One of the problems which remains is that we only have access to the result once we have found a complete parse.



Online results

One of the problems which remains is that we only have access to the result once we have found a complete parse.

- ▶ for our just introduced parsers this is obvious

$$p \langle * \rangle q = \lambda \text{inp} \rightarrow [(b2a \ b, \text{rr}) \mid (b2a, \text{prest}) \leftarrow p \ \text{inp} \\ , (b, \ \text{qrest}) \leftarrow q \ \text{prest}]$$

We only get the first element of the list of results once q has found a match!



Online results

One of the problems which remains is that we only have access to the result once we have found a complete parse.

- ▶ for our just introduced parsers this is obvious
- ▶ but this also holds for the “list-of-successes” method; it is caused by the pattern-matching in the sequential composition

$$p \langle * \rangle q = \lambda \text{inp} \rightarrow [(b2a \ b, \text{rr}) \mid (b2a, \text{prest}) \leftarrow p \ \text{inp} \\ , (b, \ \text{qrest}) \leftarrow q \ \text{prest}]$$

We only get the first element of the list of results once q has found a match!



Change of Specification

In principle the non-online behaviour is correct: we ask for a complete result, and we can only get a result once we have found at least one complete parse!



Change of Specification

In principle the non-online behaviour is correct: we ask for a complete result, and we can only get a result once we have found at least one complete parse!

We observe that, while parsing according to our breadth-first strategy, once we have only **one living alternative left** we could just as well return the result corresponding to the recognised part!



Change of Specification

In principle the non-online behaviour is correct: we ask for a complete result, and we can only get a result once we have found at least one complete parse!

We observe that, while parsing according to our breadth-first strategy, once we have only **one living alternative left** we could just as well return the result corresponding to the recognised part!

This is especially useful if we incorporate error-correction in such a way that we are guaranteed to get at least one “possibly successfully corrected” parse.



Future Based Parsers

newtype FP st a = FP ($\forall r. (st \rightarrow Steps \quad r) \rightarrow$
 $st \rightarrow Steps (a, r)$
)

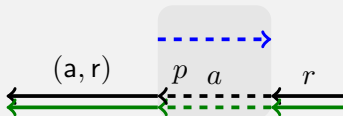
We merge fragments of the result we are constructing with the progress information:



Future Based Parsers

newtype FP st a = FP ($\forall r. (st \rightarrow Steps \quad r) \rightarrow$
 $st \rightarrow Steps \quad (a, r)$
)

We merge fragments of the result we are constructing with the progress information:



We have to make sure that if we compare two alternatives we have progress information at the head:

```

norm :: Steps a → Steps a
norm (Apply f (Step p l)) = Step p (Apply f l)
norm (Apply f (Fail ...)) = Fail ...
norm (Apply f (Apply g l)) = norm (Apply (f ∘ g) l)
norm steps                    = steps

x 'best' y = norm x 'best' norm y

best' :: Steps b → Steps b → Steps b

```



FP is Applicative

instance Applicative (FP state) **where**

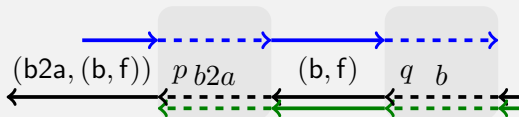
$$\begin{aligned} \text{FP pf } \langle * \rangle \sim (\text{FP qf}) &= \text{FP } ((\text{apply} \circ) \circ (\text{pf} \circ \text{qf})) \\ \text{pure a} &= \text{FP } ((\text{push a}) \circ) \end{aligned}$$

instance Alternative (FP state) **where**

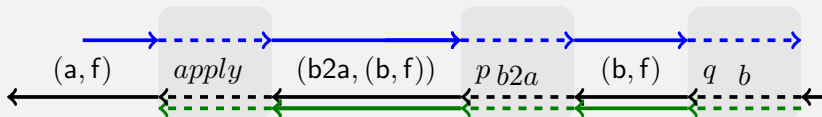
$$\begin{aligned} \text{FP pf } \langle | \rangle \text{ FP qf} &= \text{FP } (\lambda k \text{ inp} \rightarrow \text{pf } k \text{ inp 'best' qf } k \text{ inp}) \\ \text{empty} &= \text{FP } (\lambda k \text{ inp} \rightarrow \text{noAlts}) \end{aligned}$$



Sequential composition for FParser



Sequential composition for FParser



FParsec is ISParser

```
pSym a = FP (λk inp →  
  case inp of (s : ss) → if s == a then addStep ∘ push s $ k ss  
              else Fail ...  
  []           → Fail ...  
  )
```



Helper code

$\text{eval} :: \text{Steps } r \rightarrow r$
 $\text{eval } (\text{Step } n \ l) = (\text{eval } l)$
 $\text{eval } (\text{Fail } ss \ ls) = \dots$
 $\text{eval } (\text{Apply } f \ l) = f (\text{eval } l)$

$\text{push} \quad :: v \rightarrow \text{Steps } r \rightarrow \text{Steps } (v, r)$
 $\text{push } v = \text{Apply } (\lambda r \rightarrow (v, r))$

$\text{apply} :: \text{Steps } (b \rightarrow a, (b, r)) \rightarrow \text{Steps } (a, r)$
 $\text{apply} = \text{Apply } (\lambda (b2a, \sim(b, r)) \rightarrow (b2a \ b, r))$



Helper code

$\text{eval} :: \text{Steps } r \rightarrow r$
 $\text{eval } (\text{Step } n \ l) = (\text{eval } l)$
 $\text{eval } (\text{Fail } ss \ ls) = \dots$
 $\text{eval } (\text{Apply } f \ l) = f (\text{eval } l)$

$\text{push} \quad :: v \rightarrow \text{Steps } r \rightarrow \text{Steps } (v, r)$
 $\text{push } v = \text{Apply } (\lambda r \rightarrow (v, r))$

$\text{apply} :: \text{Steps } (b \rightarrow a, (b, r)) \rightarrow \text{Steps } (a, r)$
 $\text{apply} = \text{Apply } (\lambda (b2a, \sim(b, r)) \rightarrow (b2a \ b, r))$

Notice the \sim in `apply`. This makes that the function can already produce something!



10.4 A Monadic Interface



Monadic Interface: Parsing XML

Using a Monadic interface we can e.g. check an XML file for well balanced tags:

data XML = Tag t [XML]

pMany p = (:) <\$> p <*> *pMany* p <|> *pSucceed* []

pXML = **do** t ← *pOpenTag*

Tag t <\$> *pMany* *pXML* <*> *pCloseTag* t



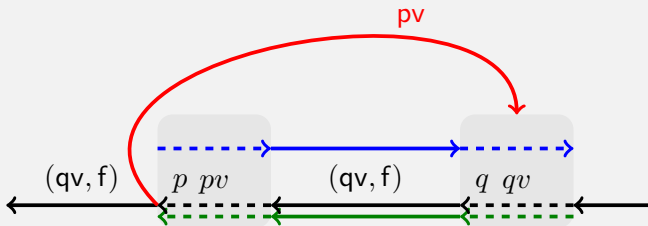
Our first attempt" FP

instance Monad FP s where

$p \gg= q = \lambda k i \rightarrow \text{let steps} = p (q \text{ } pv \text{ } k) i$
 $(pv, -) = \text{eval steps}$

in Apply snd steps

return $v = p \text{Succeed } v$



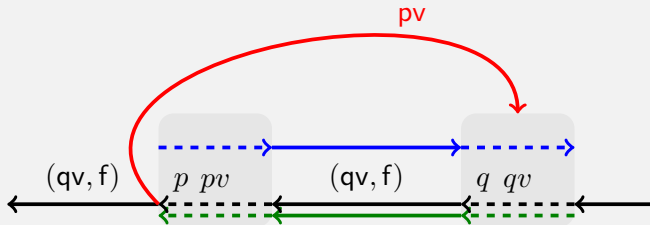
Our first attempt" FP

instance Monad FP s where

$p \gg= q = \lambda k i \rightarrow \text{let steps} = p (q \text{ } pv \text{ } k) i$
 $(pv, -) = \text{eval steps}$

in Apply snd steps

return $v = p \text{Succeed } v$



Unfortunately this is not correct. This may lead to a black hole, since the value pv may not be available yet in q , when needed.

in producing Steps. The function norm may push the value

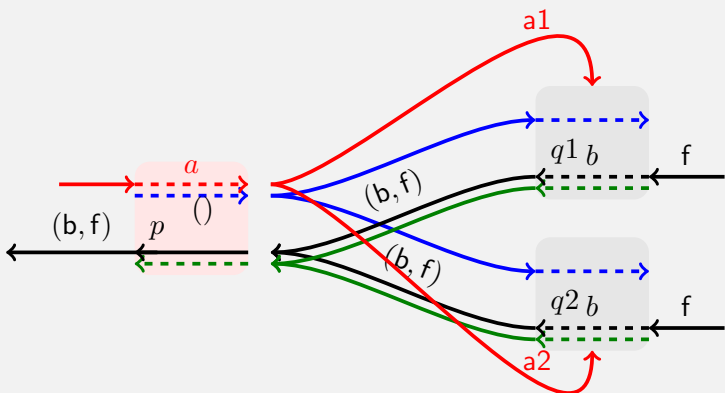
information behind the unneeded information



Solution: Combining HP and FP

$(\gg) :: \text{HP st } a \rightarrow (a \rightarrow \text{FP st } b) \rightarrow \text{FP st } b$

$p \gg q = \text{FP } (\lambda k \text{ st } \rightarrow p (\lambda pv \text{ st}' \rightarrow q \text{ } pv \text{ } k \text{ st}')) \text{ st}$



Making the solution into a Monad

Our next kind of parser is a tupling between a history based and a future based parser:

```
data Parser s a = P (HP s a) (FP s a)
```

```
instance Applicative (Parser s) where
```

```
(P hp fp) <*> ~ (P hq fq) = P (hp <*> hq) (fp <*> fq)
```

```
(P hp fp) <|> (P hq fq) = P (hp <|> hq) (fp <|> fq)
```

```
pSucceed a = P (pSucceed a) (pSucceed a)
```

```
pFail = P pFail pFail
```



The Monadic Interface Code

instance Monad (*Parser* s) **where**

$$\begin{aligned} & (\text{P (HP } p) _) \gg= qq \\ & = \text{P (HP } (\lambda k \text{ st} \rightarrow p (\lambda a \text{ st}' \rightarrow \text{unHP (qq a) k st}') \text{ st}))} \\ & \quad (\text{FP } (\lambda k \text{ st} \rightarrow p (\lambda a \text{ st}' \rightarrow \text{unFP (qq a) k st}') \text{ st})) \\ & \quad \text{where unHP (P (HP h) _) = h} \\ & \quad \quad \text{unFP (P _ (FP f) _) = f} \\ & \text{return } x = \text{P (pSucceed } x) (\text{pSucceed } x) \end{aligned}$$

Note that from left hand side of the bind we always take the history based parser, whereas for the right hand side we have two cases to take care of.



Further optimisations

Once we have started to tuple various variants of parsers we might just as well:



Further optimisations

Once we have started to tuple various variants of parsers we might just as well:

- ▶ also tuple a pure recogniser, so we can avoid construction of results which will be discarded anyway



Further optimisations

Once we have started to tuple various variants of parsers we might just as well:

- ▶ also tuple a pure recogniser, so we can avoid construction of results which will be discarded anyway
- ▶ tuple a possibly empty parser, which is needed for an efficient implementation of the permutation parser with components that may be empty



Further optimisations

Once we have started to tuple various variants of parsers we might just as well:

- ▶ also tuple a pure recogniser, so we can avoid construction of results which will be discarded anyway
- ▶ tuple a possibly empty parser, which is needed for an efficient implementation of the permutation parser with components that may be empty
- ▶ a list of possible starter symbols to be used in error messages



Further optimisations

Once we have started to tuple various variants of parsers we might just as well:

- ▶ also tuple a pure recogniser, so we can avoid construction of results which will be discarded anyway
- ▶ tuple a possibly empty parser, which is needed for an efficient implementation of the permutation parser with components that may be empty
- ▶ a list of possible starter symbols to be used in error messages
- ▶ ...



10.5 Error Correction



Error correction

We can extend the system with an error correcting mechanism.

- ▶ we may delete a symbol, at a certain cost
- ▶ we may insert a symbol, at a certain cost
- ▶ the function best does not select the longest sequence of steps, but the cheapest
- ▶ limited look-ahead is needed in order to get fast parsers



The correction function $pSym$

We show a simplified error correcting parser:

```
data Steps result = Shift (Steps result)
                  | Fail   (Steps result)
                  | Done
```

$pSym$ a =

FP \$ **let** pSym'

= λk input \rightarrow

case input **of**

inp@(b : bs) \rightarrow **if** a == b

then Step \circ push b \$ k bs

else Fail \circ push a \$ k bs

'best'

Fail (pSym' k bs)

\rightarrow Fail \circ push a \$ k input

[]

in pSym'



Refinement of Error-correcting Process

1. We may associate a cost with each insertion or deletion step, so we can take the “cheapest future”; some symbols are unlikely to have been forgotten.
2. Limited look-ahead in order to speed-up correction process
3. Store a report about the corrections taken in the state
4. Collect a list of expected symbols, in order to generate nice error messages.
5. Use an abstract interpretation to find a non-recursive alternative, in order to avoid infinite insertions.



Computing the minimal length of an alternative

In each tuple which represents a parser we incorporate a value of type Nat:

```
data Nat = Zero
```

```
  | Succ Nat deriving Show
```

```
nat_min :: Nat → Nat → Int → (Nat, Bool)
```

```
nat_min _      Zero  _ = (Zero, False)
```

```
nat_min Zero   _     _ = (Zero, True)
```

```
nat_min l      Infinite _ = (l, True)
```

```
nat_min (Succ ll) (Succ rr) n
```

```
  = if n > 1000 then error "problem with comparing length
```

```
    else let (v, b) = nat_min ll rr (n + 1)
```

```
      in (Succ v, b))
```

```
nat_add Zero   r = r
```

```
nat_add (Succ l) r = Succ (nat_add l r)
```



The Actual Parser Types

data P st a

= P (T st a) -- HP, FP and recogniser

(Maybe (T st a)) -- non-empty parsers

Nat -- minimal length

(Maybe a) -- possibly empty



The Actual Parser Types

data P st a
= P (T st a) -- HP, FP and recogniser
(Maybe (T st a)) -- non-empty parsers
Nat -- minimal length
(Maybe a) -- possibly empty

And the parsing triple:

data T st a
= T -- history
($\forall r. (a \rightarrow st \rightarrow Steps\ r) \rightarrow st \rightarrow Steps\ r$)
-- future
($\forall r. (st \rightarrow Steps\ r) \rightarrow st \rightarrow Steps\ (a, r)$)
-- recogniser
($\forall r. (st \rightarrow Steps\ r) \rightarrow st \rightarrow Steps\ r$)



Dealing with Fail

We have been a bit sloppy about failing parsers. We now give the full Fail-alternative of the Steps a type:

```
type Syms = [String]
```

```
data Steps a where
```

```
  Step ::      Progress → Steps a                → Steps a
```

```
  Apply :: ∀a b.(b → a) → Steps b                → Steps a
```

```
  Fail  ::      Syms    → [Syms → (Int, Steps a)] → Steps a
```

The Strings field keeps track of symbols which were expected.



Dealing with Fail

We have been a bit sloppy about failing parsers. We now give the full Fail-alternative of the Steps a type:

```
type Syms = [String]
```

```
data Steps a where
```

```
Step :: Progress → Steps a → Steps a
```

```
Apply :: ∀a b. (b → a) → Steps b → Steps a
```

```
Fail :: Syms → [Syms → (Int, Steps a)] → Steps a
```

The Strings field keeps track of symbols which were expected. The are collected in the function best

```
best' :: Steps b → Steps b → Steps b
```

```
Fail sl ll 'best' Fail sr rr = Fail (sl ++ sr) (ll ++ rr)
```

```
Fail _ _ 'best' r = r
```

```
l 'best' Fail _ _ = l
```



Getting rid of Fail

In case a repair was really necessary the function eval will encounter a Fail in the list of steps:

1. all the expected symbols are passed to all the alternatives
2. the resulting tree is examined upto a certain depth
3. the cheapest branch is taken



Getting rid of Fail

In case a repair was really necessary the function eval will encounter a Fail in the list of steps:

1. all the expected symbols are passed to all the alternatives
2. the resulting tree is examined upto a certain depth
3. the cheapest branch is taken

```
eval (Fail expected ls)  
  = eval (getCheapest 3 (map ($expected) ls))  
eval ...
```

