



Universiteit Utrecht

[Faculty of Science  
Information and Computing Sciences]

# Talen en Compilers

2019 - 2020, period 2

Jurriaan Hage

Department of Information and Computing Sciences  
Utrecht University

2019-11-18

# 3. Parser combinators



# This lecture

## Parser combinators

Recap

Parsing

Developing parser combinators



## 3.1 Recap



# Parsing problem

Given a grammar  $G$  and a string  $s$ , the **parsing problem** is to decide whether or not  $s \in L(G)$ .

Furthermore, if  $s \in L(G)$ , we want evidence/proof/an explanation why this is the case, usually in the form of a parse tree.



# Parse trees in Haskell

$S \rightarrow S-D \mid D$	$D \rightarrow 0 \mid 1$	<b>data</b> S = Minus S D   SingleDigit D
		<b>data</b> D = Zero   One

Concrete syntax: context-free grammar.

Abstract syntax: (Haskell datatype), does no longer contain information about terminals that can easily be reconstructed.



Context-free grammars can be used to describe lots of interesting languages.

Several grammars can describe the same language, not all of them being equally suited as a starting point for parsing.

Ambiguity is an example of an undesirable property of grammars.



## 3.2 Parsing





# Approaches to parsing

## Parser generators

- ▶ External program
- ▶ based on a bottom-up algorithm, usually LL or LR
- ▶ complex theory
- ▶ limited look-ahead, usually one token
- ▶ only built-in abstractions
- ▶ generated parsers are extremely fast

## Parser combinators

- ▶ Library
- ▶ based on a top-down algorithm
- ▶ underlying theory is simple
- ▶ in principle unlimited look-ahead
- ▶ user-definable abstractions
- ▶ fast as long as certain constructs are not used

Both approaches place certain (but different) constraints on the grammars being used.



# Approaches to parsing – contd.

In the first part of the course, we will work with combinators.



# Approaches to parsing – contd.

In the first part of the course, we will work with combinators.

Towards the end of the course, we will learn the theory of parser generators.



# Approaches to parsing – contd.

In the first part of the course, we will work with combinators.

Towards the end of the course, we will learn the theory of parser generators.

In the practicum tasks, you will use both parser generators and parser combinators.



## Aside: Combinators

The term combinator denotes a self-contained function in **lambda calculus**, the formal system that is the basis of Haskell and other functional programming languages.

**Parser combinators** are thus a set of (small) library functions that can be used to construct parsers.



# Lexing and parsing

Often, parsing is split into a two-phase process:



# Lexing and parsing

Often, parsing is split into a two-phase process:

## Lexing

In a first phase, whitespace and comments are removed and the input is organized into a sequence of **tokens** – small entities that belong together such as keywords, identifiers or operators.



# Lexing and parsing

Often, parsing is split into a two-phase process:

## Lexing

In a first phase, whitespace and comments are removed and the input is organized into a sequence of **tokens** – small entities that belong together such as keywords, identifiers or operators.

## Parsing

In the second phase, an abstract syntax tree is constructed from the list of tokens rather than from the original list of characters.





## Lexing and parsing – contd.

In the world of generators, lexing and parsing is often performed by separate generators. In Haskell, for example: Alex (lexer) and Happy (parser). For C: flex (lexer), yacc/bison (parser).



## Lexing and parsing – contd.

In the world of generators, lexing and parsing is often performed by separate generators. In Haskell, for example: Alex (lexer) and Happy (parser). For C: flex (lexer), yacc/bison (parser).

With parser combinators, there are different options:

- ▶ use only one phase,
- ▶ use the same parser combinators for both phases,
- ▶ use dedicated lexer combinators for lexing,
- ▶ use a hand-written special-purpose lexer,
- ▶ combine a lexer generator with parser combinators.



## 3.3 Developing parser combinators



# Words of warning

We are going to **develop** a suitable type of parsers.

Along the path, we will make several suboptimal or even wrong attempts.



# First attempt: a predicate on strings

| **type** Parser<sub>1</sub> = String → Bool



# First attempt: a predicate on strings

| **type** Parser<sub>1</sub> = String → Bool

With this type, we can write very simple parsers:

| manyLetters<sub>1</sub> :: Parser<sub>1</sub>  
manyLetters<sub>1</sub> xs = all isLetter xs  
someDigits<sub>1</sub> :: Parser<sub>1</sub>  
someDigits<sub>1</sub> xs = all isDigit xs ∧ not (null xs)



# First attempt – contd.

| **type** Parser<sub>1</sub> = String → Bool

Disadvantages:

- ▶ Only yes or no as answer.
- ▶ Works only on strings.
- ▶ Difficult to combine.



## First attempt – contd.

| **type** Parser<sub>1</sub> = String → Bool

Disadvantages:

- ▶ Only yes or no as answer.
- ▶ Works only on strings.
- ▶ Difficult to combine.

We first look into combining parsers, later at the other points.





# Motivation: sequencing parsers

Assume we want to combine:

```
manyLetters1 :: Parser1  
someDigits1  :: Parser1
```

and parse many letters followed by some digits.



# Motivation: sequencing parsers

Assume we want to combine:

```
manyLetters1 :: Parser1  
someDigits1  :: Parser1
```

and parse many letters followed by some digits.

We cannot, because:

- ▶ both parsers work on the complete input string
- ▶ there is no way to split the input,
- ▶ we do not (in general) know where to split the input without running the first parser.



# Second attempt: returning the remaining string

## Idea

- ▶ Parsers can consume an initial part of the input.
- ▶ Parsers only look at the initial part of the input.
- ▶ Parsers return the rest of the string.



# Second attempt: returning the remaining string

## Idea

- ▶ Parsers can consume an initial part of the input.
- ▶ Parsers only look at the initial part of the input.
- ▶ Parsers return the rest of the string.

What type to choose?

| **type** Parser<sub>2</sub> = String → ...



# Second attempt: returning the remaining string

## Idea

- ▶ Parsers can consume an initial part of the input.
- ▶ Parsers only look at the initial part of the input.
- ▶ Parsers return the rest of the string.

What type to choose?

| **type** `Parser2 = String → ...`

We need the remaining string only if parsing was successful:

| **data** Maybe a = Nothing | Just a



# Second attempt: returning the remaining string

## Idea

- ▶ Parsers can consume an initial part of the input.
- ▶ Parsers only look at the initial part of the input.
- ▶ Parsers return the rest of the string.

What type to choose?

| **type** `Parser2 = String → Maybe String`

We need the remaining string only if parsing was successful:

| **data** `Maybe a = Nothing | Just a`



## Example

We now have to write the parsers such that they work on the initial part of the string:

```
manyLetters2 :: Parser2
manyLetters2 xs = Just (dropWhile isLetter xs)
```



## Example

We now have to write the parsers such that they work on the initial part of the string:

```
manyLetters2 :: Parser2
manyLetters2 xs = Just (dropWhile isLetter xs)
```

Note that `manyLetters2` cannot fail.





## Example

We now have to write the parsers such that they work on the initial part of the string:

```
manyLetters2 :: Parser2
manyLetters2 xs = Just (dropWhile isLetter xs)
```

Note that manyLetters<sub>2</sub> cannot fail.

```
someDigits2 :: Parser2
someDigits2 xs = case span isDigit xs of
    ([], _) → Nothing
    (_, ys) → Just ys
```



## Example – contd.

We can now sequence `manyLetters2` and `someDigits2`:

```
lettersThenDigits2 :: Parser2
lettersThenDigits2 xs =
  case manyLetters2 xs of
    Nothing → Nothing
    Just ys → someDigits2 ys
```



## Example – contd.

We can now sequence `manyLetters2` and `someDigits2`:

```
lettersThenDigits2 :: Parser2
lettersThenDigits2 xs =
  case manyLetters2 xs of
    Nothing → Nothing
    Just ys → someDigits2 ys
```

We can abstract from the sequencing operation.



## Example – contd.

We can now sequence `manyLetters2` and `someDigits2`:

```
lettersThenDigits2 :: Parser2
lettersThenDigits2 xs =
  case manyLetters2 xs of
    Nothing → Nothing
    Just ys → someDigits2 ys
```

We can abstract from the sequencing operation.

```
(<*>) :: Parser2 → Parser2 → Parser2
(p <*> q) xs =
  case p xs of
    Nothing → Nothing
    Just ys → q ys
```



## Example – contd.

We can now sequence `manyLetters2` and `someDigits2`:

```
lettersThenDigits2 :: Parser2
lettersThenDigits2 xs =
  case manyLetters2 xs of
    Nothing → Nothing
    Just ys → someDigits2 ys
```

```
lettersThenDigits2 :: Parser2
lettersThenDigits2 =
  manyLetters2
  <*> someDigits2
```

We can abstract from the sequencing operation.

```
(<*>) :: Parser2 → Parser2 → Parser2
(p <*> q) xs =
  case p xs of
    Nothing → Nothing
    Just ys → q ys
```



# The end of the input

The `lettersThenDigits2` parser works as follows:

<code>lettersThenDigits<sub>2</sub> "abc123"</code>	evaluates to
<code>lettersThenDigits<sub>2</sub> "abc"</code>	evaluates to
<code>lettersThenDigits<sub>2</sub> "123"</code>	evaluates to
<code>lettersThenDigits<sub>2</sub> "a1x"</code>	evaluates to



# The end of the input

The `lettersThenDigits2` parser works as follows:

<code>lettersThenDigits<sub>2</sub> "abc123"</code>	evaluates to	Just ""
<code>lettersThenDigits<sub>2</sub> "abc"</code>	evaluates to	Nothing
<code>lettersThenDigits<sub>2</sub> "123"</code>	evaluates to	Just ""
<code>lettersThenDigits<sub>2</sub> "a1x"</code>	evaluates to	Just "x"



# The end of the input

The `lettersThenDigits2` parser works as follows:

<code>lettersThenDigits<sub>2</sub> "abc123"</code>	evaluates to	<code>Just ""</code>
<code>lettersThenDigits<sub>2</sub> "abc"</code>	evaluates to	<code>Nothing</code>
<code>lettersThenDigits<sub>2</sub> "123"</code>	evaluates to	<code>Just ""</code>
<code>lettersThenDigits<sub>2</sub> "a1x"</code>	evaluates to	<code>Just "x"</code>

We define a special parser that expects the input to be empty:

```
eof2 :: Parser2
eof2 [] = Just []
eof2 _ = Nothing
```





# The end of the input

The `lettersThenDigits2` parser works as follows:

<code>lettersThenDigits<sub>2</sub> "abc123"</code>	evaluates to	<code>Just ""</code>
<code>lettersThenDigits<sub>2</sub> "abc"</code>	evaluates to	<code>Nothing</code>
<code>lettersThenDigits<sub>2</sub> "123"</code>	evaluates to	<code>Just ""</code>
<code>lettersThenDigits<sub>2</sub> "a1x"</code>	evaluates to	<code>Just "x"</code>

We define a special parser that expects the input to be empty:

```
eof2 :: Parser2
eof2 [] = Just []
eof2 _ = Nothing
```

Now we can reject "a1x":

```
lettersThenDigits'2 = manyLetters2 <*> someDigits2 <*> eof2
```



## Another example

Consider the grammar

$$S \rightarrow \text{Letter}^* a$$



## Another example

Consider the grammar

|  $S \rightarrow \text{Letter}^* a$

Is this grammar ambiguous?



## Another example

Consider the grammar

|  $S \rightarrow \text{Letter}^* a$

Is this grammar ambiguous?

No, but it is problematic to parse with our current approach.



## Another example

Consider the grammar

|  $S \rightarrow \text{Letter}^* a$

Is this grammar ambiguous?

No, but it is problematic to parse with our current approach.

We can easily define a parser for a single a:

|  $\text{singleA}_2 :: \text{Parser}_2$   
|  $\text{singleA}_2 ('a' : xs) = \text{Just } xs$   
|  $\text{singleA}_2 \_ = \text{Nothing}$

Can you now see the problem?



# Ambiguity revisited

|  $(\text{manyLetters}_2 \langle * \rangle \text{singleA}_2)$  "cba" evaluates to Nothing

There are multiple prefixes of "cba" that can be seen as a sequence of letters, yet  $\text{manyLetters}_2$  is greedy and returns only one.

Such cases of ambiguity can arise during the parsing process even if the grammar as a whole is unambiguous.



# Ambiguity revisited

| (`manyLetters2`  $\langle * \rangle$  `singleA2`) "cba" evaluates to Nothing

There are multiple prefixes of "cba" that can be seen as a sequence of letters, yet `manyLetters2` is greedy and returns only one.

Such cases of ambiguity can arise during the parsing process even if the grammar as a whole is unambiguous.

## Our solution

Let parsers return multiple results instead of just one.

- ▶ Allows us to deal with the above case and also ambiguous grammars.
- ▶ Potential source of inefficiency.



## Third attempt: multiple results

What type to choose?

| `type Parser3 = String → ...`





## Third attempt: multiple results

What type to choose?

| **type** Parser<sub>3</sub> = String → ...

We can use a list. Failure is now represented as the empty list. Successful results are represented by their corresponding remaining strings.



## Third attempt: multiple results

What type to choose?

**|** `type Parser3 = String → [String]`

We can use a list. Failure is now represented as the empty list. Successful results are represented by their corresponding remaining strings.

The technique of using a list of successful results as a return value is called **list of successes** method.



# Choice and sequence

The new parser type gives us an easy way to write down a choice between two parsers:

$$\begin{aligned} & (<|>) :: \text{Parser}_3 \rightarrow \text{Parser}_3 \rightarrow \text{Parser}_3 \\ & (p <|> q) \text{ xs} = p \text{ xs} \text{ ++ } q \text{ xs} \end{aligned}$$



# Choice and sequence

The new parser type gives us an easy way to write down a choice between two parsers:

$$\begin{array}{l} | \quad (<|>) :: \text{Parser}_3 \rightarrow \text{Parser}_3 \rightarrow \text{Parser}_3 \\ | \quad (p <|> q) \text{ xs} = p \text{ xs} \text{ ++ } q \text{ xs} \end{array}$$

On the other hand, sequencing becomes a bit more difficult, because we have to deal with multiple results:

$$\begin{array}{l} | \quad (<*>) :: \text{Parser}_3 \rightarrow \text{Parser}_3 \rightarrow \text{Parser}_3 \\ | \quad (p <*> q) \text{ xs} = [zs \mid ys \leftarrow p \text{ xs}, zs \leftarrow q \text{ ys}] \end{array}$$



# Choice and sequence

The new parser type gives us an easy way to write down a choice between two parsers:

$$\begin{aligned} & (\langle | \rangle) :: \text{Parser}_3 \rightarrow \text{Parser}_3 \rightarrow \text{Parser}_3 \\ & (p \langle | \rangle q) \text{ xs} = p \text{ xs} \text{ ++ } q \text{ xs} \end{aligned}$$

On the other hand, sequencing becomes a bit more difficult, because we have to deal with multiple results:

$$\begin{aligned} & (\langle * \rangle) :: \text{Parser}_3 \rightarrow \text{Parser}_3 \rightarrow \text{Parser}_3 \\ & (p \langle * \rangle q) \text{ xs} = [zs \mid ys \leftarrow p \text{ xs}, zs \leftarrow q \text{ ys}] \end{aligned}$$

We define that  $(\langle * \rangle)$  binds stronger than  $(\langle | \rangle)$ :

$$\begin{aligned} & \text{infixl } 4 \langle * \rangle \\ & \text{infixl } 3 \langle | \rangle \end{aligned}$$



# Revisiting the examples

We can build  $\text{manyLetters}_3$  out of smaller blocks!

|  $\text{ManyLetters} \rightarrow \text{Letter ManyLetters} \mid \varepsilon$



# Revisiting the examples

We can build  $\text{manyLetters}_3$  out of smaller blocks!

|  $\text{ManyLetters} \rightarrow \text{Letter ManyLetters} \mid \varepsilon$

We can easily define parsers for  $\varepsilon$  and Letter:

|  $\text{epsilon}_3 :: \text{Parser}_3$

|  $\text{epsilon}_3 \text{ xs} = [\text{xs}]$

|  $\text{letter}_3 :: \text{Parser}_3$

|  $\text{letter}_3 (x : \text{xs}) \mid \text{isLetter } x = [\text{xs}]$

|  $\text{letter}_3 \_ = []$



## Revisiting the examples

We can build  $\text{manyLetters}_3$  out of smaller blocks!

$\text{ManyLetters} \rightarrow \text{Letter ManyLetters} \mid \varepsilon$

We can easily define parsers for  $\varepsilon$  and Letter:

$\text{epsilon}_3 :: \text{Parser}_3$

$\text{epsilon}_3 \text{ xs} = [\text{xs}]$

$\text{letter}_3 :: \text{Parser}_3$

$\text{letter}_3 (x : \text{xs}) \mid \text{isLetter } x = [\text{xs}]$

$\text{letter}_3 \_ = []$

Now we can define a parser for  $\text{ManyLetters}$ :

$\text{manyLetters}_3 :: \text{Parser}_3$

$\text{manyLetters}_3 = \text{letter}_3 \langle * \rangle \text{manyLetters}_3 \langle | \rangle \text{epsilon}_3$





# More abstraction

$\text{satisfy}_3 :: (\text{Char} \rightarrow \text{Bool}) \rightarrow \text{Parser}_3$

$\text{satisfy}_3 p (x : xs) \mid p x = [xs]$

$\text{satisfy}_3 \_ \_ = []$

$\text{letter}_3 = \text{satisfy}_3 \text{isLetter}$

$\text{digit}_3 = \text{satisfy}_3 \text{isDigit}$



## More abstraction

$\text{satisfy}_3 :: (\text{Char} \rightarrow \text{Bool}) \rightarrow \text{Parser}_3$

$\text{satisfy}_3 p (x : xs) \mid p x = [xs]$

$\text{satisfy}_3 \_ \_ = []$

$\text{letter}_3 = \text{satisfy}_3 \text{isLetter}$

$\text{digit}_3 = \text{satisfy}_3 \text{isDigit}$

$\text{many}_3 :: \text{Parser}_3 \rightarrow \text{Parser}_3$

$\text{many}_3 p = p \langle * \rangle \text{many}_3 p \langle | \rangle \text{epsilon}_3$

$\text{some}_3 :: \text{Parser}_3 \rightarrow \text{Parser}_3$

$\text{some}_3 p = p \langle * \rangle \text{some}_3 p \langle | \rangle p$



## More abstraction

$\text{satisfy}_3 :: (\text{Char} \rightarrow \text{Bool}) \rightarrow \text{Parser}_3$

$\text{satisfy}_3 p (x : xs) \mid p x = [xs]$

$\text{satisfy}_3 \_ \_ = []$

$\text{letter}_3 = \text{satisfy}_3 \text{isLetter}$

$\text{digit}_3 = \text{satisfy}_3 \text{isDigit}$

$\text{many}_3 :: \text{Parser}_3 \rightarrow \text{Parser}_3$

$\text{many}_3 p = p \langle * \rangle \text{many}_3 p \langle | \rangle \text{epsilon}_3$

$\text{some}_3 :: \text{Parser}_3 \rightarrow \text{Parser}_3$

$\text{some}_3 p = p \langle * \rangle \text{some}_3 p \langle | \rangle p$

$\text{manyLetters}_3 = \text{many}_3 \text{letter}_3$

$\text{someDigits}_3 = \text{some}_3 \text{digit}_3$

$\text{lettersThenDigits}_3 = \text{manyLetters}_3 \langle * \rangle \text{someDigits}_3$



## Another example

$$\begin{array}{l} | \quad I \rightarrow 0 \mid 1 \mid B \\ | \quad B \rightarrow [ E ] \\ | \quad E \rightarrow I , E \mid I \end{array}$$


## Another example

I	→	0   1   B	[0, [[1,0], [0,1,1]]]
B	→	[ E ]	1
E	→	I , E   I	[[[[0,1,0,1]]]]



## Another example

I	→ 0   1   B	[0, [[1,0], [0,1,1]]]
B	→ [ E ]	1
E	→ I , E   I	[[[[0,1,0,1]]]]

We need one additional abstraction:

symbol <sub>3</sub>	:: Char → Parser <sub>3</sub>
symbol <sub>3</sub> x	= satisfy <sub>3</sub> (== x)



## Another example

I	→ 0   1   B	[0, [[1,0], [0,1,1]]]
B	→ [ E ]	1
E	→ I , E   I	[[[[0,1,0,1]]]]

We need one additional abstraction:

symbol <sub>3</sub>	:: Char → Parser <sub>3</sub>
symbol <sub>3</sub> x	= satisfy <sub>3</sub> (== x)

The rest is entirely systematic:

i, b, e	:: Parser <sub>3</sub>
i	= symbol <sub>3</sub> '0' < > symbol <sub>3</sub> '1' < > b
b	= symbol <sub>3</sub> '[' <*> e <*> symbol <sub>3</sub> ']'
e	= i <*> symbol <sub>3</sub> ',' <*> e < > i



# Intermediate summary

We have

- ▶ a small library of basic parser combinators,
- ▶ parsers for larger grammars can be constructed easily,
- ▶ new abstractions can be defined,
- ▶ we can follow the grammar structure in order to build a parser systematically.





# Intermediate summary

We have

- ▶ a small library of basic parser combinators,
- ▶ parsers for larger grammars can be constructed easily,
- ▶ new abstractions can be defined,
- ▶ we can follow the grammar structure in order to build a parser systematically.

Still problematic:

- ▶ Only yes or no as answer.
- ▶ Works only on strings.



# Intermediate summary

We have

- ▶ a small library of basic parser combinators,
- ▶ parsers for larger grammars can be constructed easily,
- ▶ new abstractions can be defined,
- ▶ we can follow the grammar structure in order to build a parser systematically.

Still problematic:

- ▶ Only yes or no as answer.
- ▶ Works only on strings.

Let us address the answers next.



## Fourth step: adding results

Last lecture, we have seen that we can represent parse trees as values of specifically defined Haskell datatypes.



## Fourth step: adding results

Last lecture, we have seen that we can represent parse trees as values of specifically defined Haskell datatypes.

Therefore, it is clear that different parsers should return different types of results.



## Fourth step: adding results

Last lecture, we have seen that we can represent parse trees as values of specifically defined Haskell datatypes.

Therefore, it is clear that different parsers should return different types of results.

We parameterize the type of parsers over the type of the result. For each successful parse, we now return the result and the remaining string:

**type** Parser<sub>4</sub> r = String → [(r, String)]



# Simple parsers with results

$\text{epsilon}_4 :: \text{Parser}_4 ()$

$\text{epsilon}_4 \text{ xs} = [(() , \text{xs})]$

$\text{satisfy}_4 :: (\text{Char} \rightarrow \text{Bool}) \rightarrow \text{Parser}_4 \text{ Char}$

$\text{satisfy}_4 \text{ p } (x : \text{xs}) \mid \text{p } x = [(x, \text{xs})]$

$\text{satisfy}_4 \text{ -- --} = []$



# Simple parsers with results

```
epsilon4 :: Parser4 ()  
epsilon4 xs = [((()), xs)]  
  
satisfy4 :: (Char → Bool) → Parser4 Char  
satisfy4 p (x : xs) | p x = [(x, xs)]  
satisfy4 _ _ = []
```

As before (except for the types):

```
letter4, digit4 :: Parser4 Char  
letter4 = satisfy4 isLetter  
digit4 = satisfy4 isDigit  
  
symbol4 :: Char → Parser4 Char  
symbol4 x = satisfy4 (== x)
```



## Choice with results

We can easily combine parsers with  $(\langle | \rangle)$  if they have the same result type:

$$\begin{aligned} & (\langle | \rangle) :: \text{Parser}_4 a \rightarrow \text{Parser}_4 a \rightarrow \text{Parser}_4 a \\ & (p \langle | \rangle q) xs = p xs \text{ ++ } q xs \end{aligned}$$

(Definition is unchanged.)





# Choice with results

We can easily combine parsers with  $(\langle | \rangle)$  if they have the same result type:

$$\begin{array}{l} (\langle | \rangle) :: \text{Parser}_4 a \rightarrow \text{Parser}_4 a \rightarrow \text{Parser}_4 a \\ (\text{p } \langle | \rangle \text{ q}) \text{ xs} = \text{p xs} \text{ ++ } \text{q xs} \end{array}$$

(Definition is unchanged.)

## Question

What if the parsers have different result types?



# Changing the result of a parser

We define a new function

$$\begin{aligned} | \quad (<\$>) &:: (a \rightarrow b) \rightarrow \text{Parser}_4 a \rightarrow \text{Parser}_4 b \\ | \quad (f <\$> p) \text{ xs} &= [(f r, ys) \mid (r, ys) \leftarrow p \text{ xs}] \end{aligned}$$

that changes the results of a parser. It has the same priority as ( $<*>$ ):

$$| \quad \text{infixl } 4 \text{ } (<\$>)$$

This function is similar to map for lists:

$$| \quad \text{map} :: (a \rightarrow b) \rightarrow [a] \rightarrow [b]$$


# Example: bits

| Bit  $\rightarrow$  0 | 1



## Example: bits

| Bit  $\rightarrow$  0 | 1

| **data** Bit = Zero | One



## Example: bits

| Bit  $\rightarrow$  0 | 1

| **data** Bit = Zero | One

Parser:

```
bit :: Parser4 Bit
bit =
  <|>
  symbol4 '0'
  symbol4 '1'
```

Does not produce a Bit without adapting the results.



## Example: bits

| Bit  $\rightarrow$  0 | 1

| **data** Bit = Zero | One

Parser:

```
bit :: Parser4 Bit
bit = const Zero <$> symbol4 '0'
    <|> const One <$> symbol4 '1'
```

Does produce a Bit adapting the results.



## Example: bits

| Bit  $\rightarrow$  0 | 1

| **data** Bit = Zero | One

Parser:

| bit :: Parser<sub>4</sub> Bit  
| bit = const Zero <\$> symbol<sub>4</sub> '0'  
| <|> const One <\$> symbol<sub>4</sub> '1'

Does produce a Bit adapting the results.

Recall:

| const :: a  $\rightarrow$  b  $\rightarrow$  a  
| const x y = x



# Combining results

How does ( $\langle * \rangle$ ) work in the presence of results?





# Combining results

How does  $\langle * \rangle$  work in the presence of results?

One option is to return a pair of results:

$\langle * \rangle :: \text{Parser}_4\ a \rightarrow \text{Parser}_4\ b \rightarrow \text{Parser}_4\ (a, b)$

$(p \langle * \rangle q)\ xs = [((r, s), zs) \mid (r, ys) \leftarrow p\ xs, (s, zs) \leftarrow q\ ys]$



# Combining results

How does  $\langle * \rangle$  work in the presence of results?

One option is to return a pair of results:

$$\begin{aligned} \langle * \rangle &:: \text{Parser}_4 a \rightarrow \text{Parser}_4 b \rightarrow \text{Parser}_4 (a, b) \\ (p \langle * \rangle q) \text{ xs} &= [((r, s), zs) \mid (r, ys) \leftarrow p \text{ xs}, (s, zs) \leftarrow q \text{ ys}] \end{aligned}$$

Unfortunately, this is inconvenient for long sequences:

$$\begin{aligned} &\text{letter}_4 \langle * \rangle \text{letter}_4 \langle * \rangle \text{letter}_4 \langle * \rangle \text{letter}_4 \langle * \rangle \text{letter}_4 \\ &:: \text{Parser} (((((\text{Char}, \text{Char}), \text{Char}), \text{Char}), \text{Char}), \text{Char}) \end{aligned}$$


# Combining results

How does ( $\langle * \rangle$ ) work in the presence of results?

One option is to return a pair of results:

$$\begin{aligned} \langle * \rangle &:: \text{Parser}_4 a \rightarrow \text{Parser}_4 b \rightarrow \text{Parser}_4 (a, b) \\ (p \langle * \rangle q) \text{ xs} &= [((r, s), zs) \mid (r, ys) \leftarrow p \text{ xs}, (s, zs) \leftarrow q \text{ ys}] \end{aligned}$$

Unfortunately, this is inconvenient for long sequences:

$$\begin{aligned} &\text{letter}_4 \langle * \rangle \text{letter}_4 \langle * \rangle \text{letter}_4 \langle * \rangle \text{letter}_4 \langle * \rangle \text{letter}_4 \\ &:: \text{Parser} (((((\text{Char}, \text{Char}), \text{Char}), \text{Char}), \text{Char}), \text{Char}) \end{aligned}$$

We have to pattern match on these nested pairs in a subsequent function applied via ( $\langle \$ \rangle$ ). But since we are applying ( $\langle \$ \rangle$ ) anyway, there is a better option.



## Combining results – contd.

We use the following definition instead:

$$\begin{aligned} & (\langle * \rangle) :: \text{Parser}_4 (a \rightarrow b) \rightarrow \text{Parser}_4 a \rightarrow \text{Parser}_4 b \\ & (p \langle * \rangle q) xs = [(f r, zs) \mid (f, ys) \leftarrow p xs, (r, zs) \leftarrow q ys] \end{aligned}$$

Now  $(\langle * \rangle)$  is like function application lifted to parsers.



# Example: Dutch postal codes

| PostCode → Dig Dig Dig Dig Letter Letter



## Example: Dutch postal codes

| PostCode → Dig Dig Dig Dig Letter Letter

Haskell abstract syntax:

```
| data PostCode = Code Dig Dig Dig Dig Letter Letter
| type Dig      = Char -- convenient, but not precise
| type Letter   = Char -- convenient, but not precise
```



## Example: Dutch postal codes

PostCode  $\rightarrow$  Dig Dig Dig Dig Letter Letter

Haskell abstract syntax:

```
data PostCode = Code Dig Dig Dig Dig Letter Letter
type Dig      = Char -- convenient, but not precise
type Letter   = Char -- convenient, but not precise
```

Parser:

```
postCode :: Parser4 PostCode
postCode = Code <$> digit4 <*> digit4 <*> digit4 <*> digit4
           <*> letter4 <*> letter4
```

Why is this function type-correct?



## Example: Dutch post codes – contd.

Both operators associate to the left, so postCode is in fact:

```
postCode =  
(((((((Code <$> digit4) <*> digit4) <*> digit4) <*> digit4)  
      <*> letter4) <*> letter4)
```





## Example: Dutch post codes – contd.

Both operators associate to the left, so postCode is in fact:

```
postCode =  
  ((((((Code <$> digit4) <*> digit4) <*> digit4) <*> digit4)  
    <*> letter4) <*> letter4)
```

Now consider the types:

Code ::

Dig → Dig → Dig → Dig → Letter → Letter → PostCode

Code <\$> digit<sub>4</sub> ::

Parser<sub>4</sub> (Dig → Dig → Dig → Letter → Letter → PostCode)

(Code <\$> digit<sub>4</sub>) <\*> digit<sub>4</sub> ::

Parser<sub>4</sub> (Dig → Dig → Letter → Letter → PostCode)

...



# Are we done yet?

With  $\text{Parser}_4$ , we have completed all the hard work. All that remains are some final touches.

## Other symbol types

Nothing in our parser design depends on the fact that we are working on strings. All we need is a list of symbols as an input, so we can move to

| **type**  $\text{Parser}_5$   $s$   $r = [s] \rightarrow [(r, [s])]$

Some of the types change. For example:

|  $\text{satisfy}_5 :: (s \rightarrow \text{Bool}) \rightarrow \text{Parser}_5$   $s$   $s$



# Are we done yet? – contd.

Not in the lecture notes, but recommended:

## Making parsers abstract

It is better to hide the implementation of parsers:

```
newtype Parser6 s r = Parser ([s] → [(r, [s])])  
runParser (Parser p) = p
```

Allows us to replace the implementation with a better one later.



# Summary

Despite the long development, the final version is still simple:

**newtype** Parser<sub>6</sub> s r = Parser ([s] → [(r, [s])])

We have combinators representing the constructs of grammars:

- ▶ parsing individual symbols,
- ▶ choice, sequence,
- ▶ empty strings,
- ▶ repetition.

Furthermore, we can produce results and modify intermediate results.



# Quiz

1. Given a parser  $p$  of the right type:
2.  $\mid$  `natural :: Parser Char Int`
3.  $\mid$  `option :: Parser Char a  $\rightarrow$  a  $\rightarrow$  Parser Char a`



# Next lecture

- ▶ Summary of the interface of the parser combinators.
- ▶ Constructing parsers from grammars.
- ▶ Pitfalls and limitations.
- ▶ Grammar transformations.

