# Talen en Compilers

## 2010/2011, periode 2

Johan Jeuring

Department of Information and Computing Sciences
Utrecht University

November 24, 2010

**Universiteit Utrecht**

# 5. Grammar and parser design

# This lecture

Grammar and parser design

Grammar

Lexing / scanning and parsing

Semantics

Loose ends

Universiteit Utrecht

# Example: Travel schemes

```
       Groningen    8:37
 9:44  Zwolle       9:49
10:15  Utrecht     10:21
11:05  Den Haag
```

Universiteit Utrecht

# Example: Travel schemes

```
        Groningen    8:37
 9:44   Zwolle       9:49
10:15   Utrecht     10:21
11:05   Den Haag
```

Some questions to ask for a given travel scheme:

- ▶ What is the net travel time?
- ▶ What is the total waiting time?
- ▶ What is the minimal change time?
- ▶ . . .

Universiteit Utrecht

# Designing a grammar and a parser

- ▶ Give some example inputs.
- ▶ Construct a grammar from example inputs.
- ▶ Test the grammar on the example inputs.
- ▶ Analyze the grammar.
- ▶ Possibly transform the grammar.
- ▶ Decide on the types.
- ▶ Construct a basic parser.
- ▶ Define semantic functions.
- ▶ Check the results.

# 5.1 Grammar

# Step 1: Give some example inputs

```
        Groningen    8:37
  9:44  Zwolle       9:49
10:15   Utrecht     10:21
11:05   Den Haag

        Zwolle
```

# Step 2: Constructing a grammar

```
        Groningen   8:37
 9:44   Zwolle      9:49
10:15   Utrecht    10:21
11:05   Den Haag
```

$$
\begin{aligned}
\text{TS} \quad &\rightarrow \text{TS Departure Arrival TS} \\
&\quad | \;\; \text{Station} \\
\text{Station} \quad &\rightarrow \text{Identifier}^{+} \\
\text{Departure} &\rightarrow \text{Time} \\
\text{Arrival} \quad &\rightarrow \text{Time} \\
\text{Time} \quad &\rightarrow \text{Nat} : \text{Nat}
\end{aligned}
$$

Universiteit Utrecht

# Step 3: Testing the grammar

```
        Groningen    8:37
 9:44   Zwolle       9:49
10:15   Utrecht     10:21
11:05   Den Haag


        Zwolle
```

| | |
|---|---|
| TS | → TS Departure Arrival TS |
| | \| Station |
| Station | → Identifier$^+$ |
| Departure | → Time |
| Arrival | → Time |
| Time | → Nat : Nat |

Universiteit Utrecht

# Step 4: Analyzing the grammar

TS        $\rightarrow$ TS Departure Arrival TS
          | Station
Station   $\rightarrow$ Identifier$^+$
Departure $\rightarrow$ Time
Arrival   $\rightarrow$ Time
Time      $\rightarrow$ Nat : Nat

Universiteit Utrecht

# Step 4: Analyzing the grammar

$$
\begin{array}{ll}
\text{TS} & \rightarrow \text{TS Departure Arrival TS} \\
 & \mid \text{Station} \\
\text{Station} & \rightarrow \text{Identifier}^+ \\
\text{Departure} & \rightarrow \text{Time} \\
\text{Arrival} & \rightarrow \text{Time} \\
\text{Time} & \rightarrow \text{Nat:Nat}
\end{array}
$$

## Observations

- ▶ The grammar is not explicit about the use of whitespace.
- ▶ A single station is a valid travel scheme according to the grammar.
- ▶ The grammar for times is imprecise.
- ▶ The grammar is left-recursive.

Universiteit Utrecht

TS → Station Departure
    (Arrival Station Departure)*
    Arrival Station
  | Station

Universiteit Utrecht

# Another grammar for the language

TS → Station Departure
       (Arrival Station Departure)*
       Arrival Station
    | Station

- ▶ Has a different focus.
- ▶ Not left-recursive, but can be left-factored.

# Step 5: Transforming the grammar

TS → TS Departure Arrival TS
    |  Station

Universiteit Utrecht

# Step 5: Transforming the grammar

TS → TS Departure Arrival TS
    | Station

The symbols Departure Arrival are an associative separator:

TS → Station Departure Arrival TS
    | Station

Universiteit Utrecht

# Step 5: Transforming the grammar

TS → TS Departure Arrival TS
    | Station

The symbols Departure Arrival are an associative separator:

TS → Station Departure Arrival TS
    | Station

Simplification:

TS → (Station Departure Arrival)$^*$ Station

Universiteit Utrecht

# Step 5: Transforming the grammar

TS → TS Departure Arrival TS
   | Station

The symbols Departure Arrival are an associative separator:

TS → Station Departure Arrival TS
   | Station

Simplification:

TS → (Station Departure Arrival)$^*$ Station

Abstraction:

TS  → Leg$^*$ Station
Leg → Station Departure Arrival

Universiteit Utrecht

# Step 6: Deciding on the types

Which grammar do we use as a basis for the abstract syntax?

Which grammar do we use as a basis for the parser?

Universiteit Utrecht

# Step 6: Deciding on the types

Which grammar do we use as a basis for the abstract syntax?

Which grammar do we use as a basis for the parser?

- ▶ No need to use the same grammar for both purposes.
- ▶ Main criteria for abstract syntax: readability, possibility to define semantic functions.
- ▶ Main criterion for parser: efficiency, i.e., left-factored and no left-recursion.
- ▶ If the underlying grammars for abstract syntax and grammar are different, more work is required in the semantic functions.

Universiteit Utrecht

# Deciding on the type – contd.

In our case, the grammar

TS  $\to$ Leg* Station
Leg $\to$ Station Departure Arrival

is both readable and suitable for a parser.

**Universiteit Utrecht**

# Deciding on the type – contd.

In our case, the grammar

TS $\to$ Leg* Station
Leg $\to$ Station Departure Arrival

is both readable and suitable for a parser.

**data** TS         = TS [Leg] Station
**data** Leg        = Leg Station Departure Arrival
**type** Station    = String
**type** Departure  = Time
**type** Arrival    = Time
**data** Time       = Time Hours Minutes
**type** Hours      = Int
**type** Minutes    = Int

# 5.2 Lexing / scanning and parsing

# Dealing with whitespace

In the grammar, we left the handling of spaces implicit.

Intuitively, any two **tokens** (semantically connected sequences of characters) can be separated by spaces.

Universiteit Utrecht

# Dealing with whitespace

In the grammar, we left the handling of spaces implicit.

Intuitively, any two **tokens** (semantically connected sequences of characters) can be separated by spaces.

There are (at least) three possiblities to deal with spaces:

- ► write a scanner by hand that produces tokens,
- ► construct a scanner using parser combinators that produces tokens,
- ► deal with spaces in the parser itself.

We have a look at the latter two options.

Universiteit Utrecht

# A parser for whitespace

First attempt:

> spaces :: Parser Char String
> spaces = many (satisfy isSpace)

# A parser for whitespace

First attempt:

> spaces :: Parser Char String
> spaces = many (satisfy isSpace)

## Question

What about

> moreSpaces :: Parser Char String
> moreSpaces = (++) <$> spaces <*> spaces

?

# Observations regarding whitespace

- Where to call spaces? It is best to handle whitespace systematically, in order to prevent uses of spaces all over the place.
- Multiple sequential uses of spaces lead to ambiguity.
- Most often, we want to discard **all** whitespace in a particular place, i.e., we are not interested in the partial results returned by many.

Universiteit Utrecht

# Whitespace policy

We can systematically handle whitespace if we agree that

- ▶ whenever a parser consumes a complete token from the input, it is responsible for consuming subsequent whitespace,
- ▶ the parser for the start symbol starts by consuming initial whitespace,
- ▶ no other parser consumes any whitespace.

Universiteit Utrecht

# Greedy parsers

In order to prevent partial results, we extend our **primitive** parser combinators with a greedy choice operator.

```
(≪|>) :: Parser s a → Parser s a → Parser s a
Parser p ≪|> Parser q = Parser (λxs → let r = p xs
                                        in if null r
                                           then q xs
                                           else  r)
```

Universiteit Utrecht

# Greedy parsers

In order to prevent partial results, we extend our **primitive** parser combinators with a greedy choice operator.

$(\ll|>) :: \text{Parser s a} \to \text{Parser s a} \to \text{Parser s a}$
$\text{Parser p} \ll|> \text{Parser q} = \text{Parser } (\lambda xs \to \textbf{let } r = p \ xs$
$$\textbf{in if } \text{null } r$$
$$\textbf{then } q \ xs$$
$$\textbf{else } \ r)$$

We can then define greedy versions of many and some:

$\text{greedy}, \text{greedy}_1 :: \text{Parser s a} \to \text{Parser s [a]}$
$\text{greedy} \ \ p = (:) <\$> p <\!*\!> \text{greedy } p \ll|> \text{succeed } []$
$\text{greedy}_1 \ p = (:) <\$> p <\!*\!> \text{greedy } p$

Universiteit Utrecht

# Parsing spaces greedily

```
spaces :: Parser Char String
spaces = greedy (satisfy isSpace)
```

Universiteit Utrecht

# Parsing spaces greedily

spaces :: Parser Char String
spaces = greedy (satisfy isSpace)

Now

moreSpaces :: Parser Char String
moreSpaces = (++) <$> spaces <*> spaces

is not even all that problematic.

Universiteit Utrecht

# Parsing spaces greedily

spaces :: Parser Char String
spaces = greedy (satisfy isSpace)

Now

moreSpaces :: Parser Char String
moreSpaces = (++) <$> spaces <*> spaces

is not even all that problematic.

In general, use of greedy can improve parser efficiency. But careful in cases where backtracing may be needed!

# Step 7: A direct parser for travel schemes

We have two kinds of tokens that occur in travel schemes:

- ▶ station names,
- ▶ times.

Universiteit Utrecht

# Step 7: A direct parser for travel schemes

We have two kinds of tokens that occur in travel schemes:

- ▶ station names,
- ▶ times.

Let us first write parsers for these.

Grammar:

Station $\rightarrow$ Identifier$^+$

Parser:

station :: Parser Char Station
station = unwords <$> greedy$_1$ (identifier <* spaces)

Note how we consume spaces **in the end**.

Universiteit Utrecht

[Faculty of **Science**
Information and Computing Sciences]

# Tokens in travel schemes

Time → Nat : Nat

Haskell:

```
data Time    = Time Hours Minutes
type Hours   = Int
type Minutes = Int
```

Parser:

```
time :: Parser Char Time
time =
    Time <$> natural <* symbol ':' <*> natural <* spaces
```

We choose not to allow spaces between the numbers and the colon.

# A direct parser for travel schemes

TS  → Leg* Station
Leg → Station Departure Arrival

Haskell:

**data** TS  = TS  [Leg] Station
**data** Leg = Leg Station Departure Arrival

# A direct parser for travel schemes

TS $\rightarrow$ Leg* Station
Leg $\rightarrow$ Station Departure Arrival

Haskell:

```
data TS = TS [Leg] Station
data Leg = Leg Station Departure Arrival
```

Parser:

```
ts :: Parser Char TS
ts = TS <$> many leg <*> station

leg :: Parser Char Leg
leg = Leg <$> station <*> time <*> time
```

Universiteit Utrecht

# A direct parser for travel schemes

TS &rarr; Leg* Station
Leg &rarr; Station Departure Arrival

Haskell:

```
data TS = TS [Leg] Station
data Leg = Leg Station Departure Arrival
```

Parser:

```
ts :: Parser Char TS
ts = TS <$> many leg <*> station

leg :: Parser Char Leg
leg = Leg <$> station <*> time <*> time
```

No space handling required!

Universiteit Utrecht

# Running the parser

As agreed, we parse initial spaces in the top-level parser:

```
start :: Parser Char TS
start = spaces *> ts <* eof
```

We can now run parse start on an input string and get a result.

Universiteit Utrecht

# Using an explicit scanner

If we want to define a separate scanner, we start by defining a datatype for Tokens:

```
data Token = TStation Station
           | TTime    Time
```

# Using an explicit scanner

If we want to define a separate scanner, we start by defining a datatype for Tokens:

```
data Token = TStation Station
           | TTime    Time
```

We adapt the parsers for stations and times to produce tokens:

```
tstation :: Parser Char Token
tstation = TStation . unwords <$> greedy₁ (identifier <* spaces)
ttime :: Parser Char Token
ttime =
  TTime <$>
  (Time <$> natural <* symbol ':' <*> natural <* spaces)
```

# Using an explicit scanner – contd.

The scanner parses any number of tokens:

> anyToken :: Parser Char Token
> anyToken = tstation $<|>$ ttime
>
> scan :: Parser Char [Token]
> scan = spaces $*>$ greedy anyToken $<*$ eof

In the subsequent parsing phase, we then work with Token as the type of symbols.

Universiteit Utrecht

# Parsing a particular type of tokens

```
station :: Parser Token Station
station = fromStation <$> satisfy isStation

isStation :: Token → Bool
isStation (TStation _) = True
isStation _            = False

fromStation :: Token → Station
fromStation (TStation x) = x
fromStation _            = error "fromStation"
```

# Parsing a particular type of tokens

```
station :: Parser Token Station
station = fromStation <$> satisfy isStation

isStation :: Token → Bool
isStation (TStation _) = True
isStation _            = False

fromStation :: Token → Station
fromStation (TStation x) = x
fromStation _            = error "fromStation"
```

Similarly:

```
time :: Parser Token Time
time = fromTime <$> satisfy isTime
```

Universiteit Utrecht

# A token parser for travel schemes

The rest of the parser is unchanged apart from the symbol type:

ts :: Parser Token TS
ts = TS <$> many leg <*> station

leg :: Parser Token Leg
leg = Leg <$> station <*> time <*> time

Universiteit Utrecht

# Use a scanner or not?

Whether to use a scanner or not is mostly a matter of taste.

Universiteit Utrecht

# Use a scanner or not?

Whether to use a scanner or not is mostly a matter of taste.

## Advantages

- ▶ Easier to keep whitespace handling localized.
- ▶ Easier to maintain decent efficiency.
- ▶ Easier to produce good error messages.

Universiteit Utrecht

# Use a scanner or not?

Whether to use a scanner or not is mostly a matter of taste.

## Advantages

- ▶ Easier to keep whitespace handling localized.
- ▶ Easier to maintain decent efficiency.
- ▶ Easier to produce good error messages.

## Disadvantages

- ▶ Extra work required.
- ▶ All of the advantages can also be gained without a separate scanner.

Universiteit Utrecht

# 5.3 Semantics

# Step 8: Adding semantic functions

On the abstract syntax, we can define semantic functions.

Semantic functions typically follow the datatypes closely:

```
data TS        = TS [Leg] Station
data Leg       = Leg Station Departure Arrival
type Station   = String
type Departure = Time
type Arrival   = Time
data Time      = Time Hours Minutes
type Hours     = Int
type Minutes   = Int
```

Universiteit Utrecht

# Net travel time

netTravelTimeTS :: TS → Minutes
netTravelTimeTS (TS ls _) = sum (map netTravelTimeLeg ls)

netTravelTimeLeg :: Leg → Minutes
netTravelTimeLeg (Leg _ dep arr) = arr 'timeDiff' dep

# Net travel time

netTravelTimeTS :: TS $\rightarrow$ Minutes
netTravelTimeTS (TS ls _) = sum (map netTravelTimeLeg ls)
netTravelTimeLeg :: Leg $\rightarrow$ Minutes
netTravelTimeLeg (Leg _ dep arr) = arr 'timeDiff' dep

The rest is dealing with times.

Universiteit Utrecht

We assume that no two times are more than a day apart:

$$
\begin{aligned}
&\text{timeDiff :: Time} \rightarrow \text{Time} \rightarrow \text{Minutes} \\
&\text{timeDiff (Time } h_1\ m_1) \text{ (Time } h_2\ m_2) \\
&\quad | \ h_2 > h_1 \vee (h_2 \text{ == } h_1 \wedge m_2 > m_1) \\
&\qquad\qquad = (h_2 \qquad - h_1) * 60 + m_2 - m_1 \\
&\quad | \ \text{otherwise} = (h_2 + 24 - h_1) * 60 + m_2 - m_1
\end{aligned}
$$

# Waiting time

waitingTimeTS :: TS → Minutes
waitingTimeTS (TS ls _) = waitingTimeLegs ls

waitingTimeLegs :: [Leg] → Minutes
waitingTimeLegs (Leg _ _ arr : ls@(Leg _ dep _ : _)) =
    dep 'timeDiff' arr + waitingTimeLegs ls
waitingTimeLegs _ = 0

**Universiteit Utrecht**

waitingTimeTS :: TS → Minutes
waitingTimeTS (TS ls _) = waitingTimeLegs ls

waitingTimeLegs :: [Leg] → Minutes
waitingTimeLegs (Leg _ _ arr : ls@(Leg _ dep _ : _)) =
    dep 'timeDiff' arr + waitingTimeLegs ls
waitingTimeLegs _ = 0

Minimal waiting time is a variation of this theme: first compute
a list of waiting times per station, then compute the minimum
thereof.

# Eliminating the abstract syntax tree

If we are only interested in one particular semantic function, there is no need to compute an abstract syntax tree first – we can plug the semantic function directly into the parser.

Universiteit Utrecht

# Eliminating the abstract syntax tree

If we are only interested in one particular semantic function, there is no need to compute an abstract syntax tree first – we can plug the semantic function directly into the parser.

For net travel time:

```
ts :: Parser Token Minutes
ts = sum <$> many leg <* station

leg :: Parser Token Minutes
leg = flip timeDiff <$ station <*> time <*> time
```

# Step 9: Checking and improving

Of course, we should test the parser on lots of examples, and make adaptions if necessary.

**Universiteit Utrecht**

# 5.4 Loose ends

# Parsing times

Both syntax and parser for times are too liberal.

$\text{Time} \rightarrow \text{Nat} : \text{Nat}$

Universiteit Utrecht

# Parsing times

Both syntax and parser for times are too liberal.

$$\text{Time} \rightarrow \text{Nat} : \text{Nat}$$

Two options:

- ▶ Adapt the grammar, and hence the parser – morally correct, but tedious.
- ▶ A more pragmatic solution: first parse liberally, then check and perhaps reject afterwards.

# Another sequencing combinator

We cannot reject a result using $(<\$>)$:

$(<\$>) :: (a \rightarrow b) \rightarrow \text{Parser s a} \rightarrow \text{Parser s b}$

Does not work:

hours :: Parser Char Hours
hours $= (\lambda x \rightarrow \textbf{if } x < 24 \textbf{ then} \ldots \textbf{else} \ldots ) <\$>$ natural

# Another sequencing combinator

We cannot reject a result using $(<\$>)$:

$(<\$>) :: (a \rightarrow b) \rightarrow \text{Parser s a} \rightarrow \text{Parser s b}$

Does not work:

hours :: Parser Char Hours
hours $= (\lambda x \rightarrow \textbf{if } x < 24 \textbf{ then} \ldots \textbf{else} \ldots) <\$>$ natural

What if we could build a new parser based on a previous result?

Universiteit Utrecht

# Another sequencing combinator – contd.

We introduce $(\gg\!\!=)$ – pronounced "bind". This is another **primitive** parser combinator:

```
(≫=) :: Parser s a → (a → Parser s b) → Parser s b
Parser p ≫= f =
    Parser (λxs → [(s, zs) | (r, ys) ← p xs,
                            (s, zs) ← runParser (f r) ys])
```

Universiteit Utrecht

# Another sequencing combinator – contd.

We introduce $(\gg\!=)$ – pronounced "bind". This is another **primitive** parser combinator:

$(\gg\!=) :: \mathsf{Parser}\ s\ a \to (a \to \mathsf{Parser}\ s\ b) \to \mathsf{Parser}\ s\ b$
$\mathsf{Parser}\ p \gg\!= f =$
   $\mathsf{Parser}\ (\lambda xs \to [(s, zs) \mid (r, ys) \leftarrow p\ xs,$
                              $(s, zs) \leftarrow \mathsf{runParser}\ (f\ r)\ ys])$

Now:

$\mathsf{hours} :: \mathsf{Parser}\ \mathsf{Char}\ \mathsf{Hours}$
$\mathsf{hours} = \mathsf{natural} \gg\!= \lambda x \to$
        **if** $x < 24$ **then** $\mathsf{succeed}\ x$ **else** $\mathsf{empty}$

Universiteit Utrecht

### Question

What is the difference between the following two keyword parsers?

$$\text{keyword}_1, \text{keyword}_2 :: \text{String} \rightarrow \text{Parser Char String}$$
$$\text{keyword}_1 \quad = \text{token}$$
$$\text{keyword}_2 \ xs = \text{greedy isLetter} \gg= \lambda ys \rightarrow$$
$$\qquad\qquad \textbf{if } xs \text{ == } ys \textbf{ then } \text{succeed } ys \textbf{ else } \text{empty}$$

# Using bind

## Question

What is the difference between the following two keyword parsers?

$$\text{keyword}_1, \text{keyword}_2 :: \text{String} \rightarrow \text{Parser Char String}$$
$$\text{keyword}_1 \quad = \text{token}$$
$$\text{keyword}_2 \; \text{xs} = \text{greedy isLetter} \gg\!\!= \lambda \text{ys} \rightarrow$$
$$\quad\quad\quad \textbf{if } \text{xs} == \text{ys } \textbf{then } \text{succeed ys } \textbf{else } \text{empty}$$

## Hint

Consider keyword "let" $*\!\!>$ spaces $*\!\!>$ identifier and the string "letx".

Universiteit Utrecht

# Applicative functors

The operations parsers support a very common operations –
many other types support the same interface(s):

```
class Functor f where
    fmap  :: (a → b) → f a → f b

(<$>) = fmap

class Functor f ⇒ Applicative f where
    pure   :: a → f a
    (<*>) :: f (a → b) → f a → f b

class Applicative f ⇒ Alternative f where
    empty :: f a
    (<|>) :: f a → f a → f a
```

Universiteit Utrecht

# Monads

**class** Monad m **where**
   return :: a → m a
   (≫=) :: m a → (a → m b) → m b

In contrast to applicative functors, you have probably seen
monads before.

Universiteit Utrecht

# Monads

class Monad m where
 return :: a → m a
 (≫=) :: m a → (a → m b) → m b

In contrast to applicative functors, you have probably seen
monads before.

More about applicative functors and monads in the master
course on **Advanced Functional Programming**.

Universiteit Utrecht

# A common pitfall

### Question

What happens here?

```
many (option (symbol 'x') ' ')
```

## Question

What happens here?

| many (option (symbol 'x') ' ')

In general, be very careful not to call many on anything that can succeed on the empty string.

# A common pitfall

## Question

What happens here?

```
many (option (symbol 'x') ' ')
```

In general, be very careful not to call many on anything that can succeed on the empty string.

In particular, many (many p) will always go wrong.