**Universiteit Utrecht**

# Talen en Compilers

## 2010/2011, periode 2

Johan Jeuring

Department of Information and Computing Sciences
Utrecht University

November 24, 2010

# 4. Parser combinators and grammar transformations

# This lecture

Parser combinators and grammar transformations

Parser combinators summary

Grammar transformations

Separators and operators

Universiteit Utrecht

# 4.1 Parser combinators summary

# Parser combinator interface

These are from the last lecture:

**type** Parser s r $=$ [s] $\rightarrow$ [(r, [s])]
epsilon :: Parser s ()
$(<|>)$ :: Parser s a $\rightarrow$ Parser s a $\rightarrow$ Parser s a
$(<\!\!*\!\!>)$ :: Parser s (a $\rightarrow$ b) $\rightarrow$ Parser s a $\rightarrow$ Parser s b
$(<\!\$\!>)$ :: (a $\rightarrow$ b) $\rightarrow$ Parser s a $\rightarrow$ Parser s b
satisfy :: (s $\rightarrow$ Bool) $\rightarrow$ Parser s s

Universiteit Utrecht

# Parser combinator interface

These are from the last lecture:

**type** Parser s r    **(abstract)**

epsilon :: Parser s ()

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

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

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

satisfy :: (s $\rightarrow$ Bool) $\rightarrow$ Parser s s

# Parser combinator interface

These are from the last lecture:

**type** Parser s r   **(abstract)**
epsilon :: Parser s ()
$(<|>)$ :: Parser s a $\rightarrow$ Parser s a $\rightarrow$ Parser s a
$(<*>)$ :: Parser s (a $\rightarrow$ b) $\rightarrow$ Parser s a $\rightarrow$ Parser s b
$(<\$>)$ :: (a $\rightarrow$ b) $\rightarrow$ Parser s a $\rightarrow$ Parser s b
satisfy :: (s $\rightarrow$ Bool) $\rightarrow$ Parser s s

And the parser for the empty language (also called failp):

empty :: Parser s a
empty = Parser (const [])

# Derived parser combinators

We have seen more functions, but these can be defined in terms of the basic combinators:

symbol :: Eq s $\Rightarrow$ s $\rightarrow$ Parser s s
symbol = satisfy (== x)

many :: Parser s a $\rightarrow$ Parser s [a]
many p = (:) <$> p <*> many p <|> const [] <$> epsilon

some :: Parser s a $\rightarrow$ Parser s [a]    -- also called many$_1$
some p = (:) <$> p <*> many p

Universiteit Utrecht

# Derived parser combinators

We have seen more functions, but these can be defined in terms of the basic combinators:

symbol :: Eq s $\Rightarrow$ s $\rightarrow$ Parser s s
symbol = satisfy (== x)

many :: Parser s a $\rightarrow$ Parser s [a]
many p = (:) <$> p <*> many p <|> const [] <$> epsilon

some :: Parser s a $\rightarrow$ Parser s [a]   -- also called $many_1$
some p = (:) <$> p <*> many p

Similarly:

option :: Parser s a $\rightarrow$ a $\rightarrow$ Parser s a
option p default = p <|> const default <$> epsilon

Universiteit Utrecht

# Example: matching parentheses

Consider this grammar:

$S \rightarrow ( S ) S \mid \varepsilon$

Universiteit Utrecht

# Example: matching parentheses

Consider this grammar:

$$S \rightarrow ( S ) S \mid \varepsilon$$

Haskell datatype:

```
data Parens = Match Parens Parens | Empty
```

Universiteit Utrecht

# Example: matching parentheses

Consider this grammar:

$S \rightarrow (\ S\ )\ S \mid \varepsilon$

Haskell datatype:

```
data Parens = Match Parens Parens | Empty
```

Parser:

```
parens :: Parser Char Parens
parens =
          <$> symbol '(' <*> parens <*> symbol ')'
          <*> parens
      <|>             <$> epsilon
```

Universiteit Utrecht

# Example: matching parentheses

Consider this grammar:

$$S \rightarrow ( S ) S \mid \varepsilon$$

Haskell datatype:

```
data Parens = Match Parens Parens | Empty
```

Parser:

```
parens :: Parser Char Parens
parens  =         (λ_ x _ y → Match x y)
           <$> symbol '(' <*> parens <*> symbol ')'
           <*>  parens
       <|> const Empty <$> epsilon
```

Universiteit Utrecht

# More derived combinators

We often need to fill in a result for $\varepsilon$:

```
succeed :: a → Parser s a
succeed x = const x <$> epsilon
```

Universiteit Utrecht

# More derived combinators

We often need to fill in a result for $\varepsilon$:

```
succeed :: a → Parser s a
succeed x = const x <$> epsilon
```

We often do not need all the results of a sequence:

```
(<$) :: a → Parser b → Parser a
x <$ p = const x <$> p

(<*) :: Parser a → Parser b → Parser a
p <* q = const <$> p <*> q

(*>) :: Parser a → Parser b → Parser b
p *> q = flip const <$> p <*> q
```

# Matched parentheses again

We can now improve the parser from

```
parens :: Parser Char Parens
parens  =           (λ_ x _ y → Match x y)
            <$> symbol '(' <*> parens <*> symbol ')'
            <*> parens
        <|> const Empty <$> epsilon
```

to

```
parens :: Parser Char Parens
parens =
   Match <$ symbol '(' <*> parens <* symbol ')' <*> parens
    <|> succeed Empty
```

# 4.2 Grammar transformations

# Removing duplicate productions

Example:

$A \rightarrow u \mid u \mid v$

can be transformed into

$A \rightarrow u \mid v$

Universiteit Utrecht

# Removing duplicate productions

Example:

$A \rightarrow u \mid u \mid v$

can be transformed into

$A \rightarrow u \mid v$

Parser:

$a = u <|> u <|> v$

becomes

$a = u <|> v$

Universiteit Utrecht

# Removing duplicate productions

Example:

$A \rightarrow u \mid u \mid v$

can be transformed into

$A \rightarrow u \mid v$

Parser:

$a = u <|> u <|> v$

becomes

$a = u <|> v$

- ▶ Removes a source of ambiguity.
- ▶ Simplifies the grammar and the code.
- ▶ Improves efficiency of the parser.

Universiteit Utrecht

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

# Inlining nonterminals

Example:

$A \rightarrow uBv \mid z$
$B \rightarrow x \mid w$

can be transformed into

$A \rightarrow uxv \mid uwv \mid z$

$B \rightarrow x \mid w$

Universiteit Utrecht

# Inlining nonterminals

Example:

$A \rightarrow uBv \mid z$
$B \rightarrow x \mid w$

can be transformed into

$A \rightarrow uxv \mid uwv \mid z$

$B \rightarrow x \mid w$

Parser:

$a = u <*> b <*> v <|> z$
$b = x <|> w$

becomes

$a = \quad\quad u <*> x <*> v$
$\quad\quad <|> u <*> w <*> v <|> z$
$b = x <|> w$

Universiteit Utrecht

# Inlining nonterminals

Example:

$A \rightarrow uBv \mid z$
$B \rightarrow x \mid w$

can be transformed into

$A \rightarrow uxv \mid uwv \mid z$

$B \rightarrow x \mid w$

Parser:

a = u <*> b <*> v <|> z
b = x <|> w

becomes

a =        u <*> x <*> v
           <|> u <*> w <*> v <|> z
b = x <|> w

- ▶ Mainly attractive if the nonterminal is used in only a few places, and after inlining becomes unreachable.
- ▶ The reverse transformation – introducing a new nonterminal – can also be useful.
- ▶ No effect on efficiency of the parser.

Universiteit Utrecht

# Removing unreachable productions

Example:

$A \rightarrow uxv \mid uwv \mid z$

$B \rightarrow x \mid w$

can be transformed into

$A \rightarrow uxv \mid uwv \mid z$

Universiteit Utrecht

# Removing unreachable productions

Example:

$A \rightarrow uxv \mid uwv \mid z$

$B \rightarrow x \mid w$

can be transformed into

$A \rightarrow uxv \mid uwv \mid z$

Parser:

$a = \quad u <\!*\!> x <\!*\!> v$
$\quad\quad <\!|\!> u <\!*\!> w <\!*\!> v <\!|\!> z$
$b = x <\!|\!> w$

becomes

$a = \quad u <\!*\!> x <\!*\!> v$
$\quad\quad <\!|\!> u <\!*\!> w <\!*\!> v <\!|\!> z$

Universiteit Utrecht

# Removing unreachable productions

Example:

$A \rightarrow$ uxv | uwv | z

$B \rightarrow$ x | w

can be transformed into

$A \rightarrow$ uxv | uwv | z

Parser:

```
a =      u <*> x <*> v
     <|> u <*> w <*> v <|> z
b = x <|> w
```

becomes

```
a =      u <*> x <*> v
     <|> u <*> w <*> v <|> z
```

- ► Only if B is unreachable and not needed as a (secondary) start symbol.
- ► Simplifies the grammar.
- ► Corresponds to dead code removal.

Universiteit Utrecht

# Left factoring

Example:

A → xy | xz | v

can be transformed into

A → xQ | v
Q → y | z

Universiteit Utrecht

# Left factoring

Example:

$A \to xy \mid xz \mid v$

can be transformed into

$A \to xQ \mid v$
$Q \to y \mid z$

Parser:

$a = x <\!*\!> y <\!|\!> x <\!*\!> z <\!|\!> v$

becomes

$a = x <\!*\!> q <\!|\!> v$
$q = y <\!|\!> z$

Universiteit Utrecht

# Left factoring

Example:

Parser:

$A \rightarrow xy \mid xz \mid v$

$a = x <\!*\!> y <\!|\!> x <\!*\!> z <\!|\!> v$

can be transformed into

becomes

$A \rightarrow xQ \mid v$
$Q \rightarrow y \mid z$

$a = x <\!*\!> q <\!|\!> v$
$q = y <\!|\!> z$

▶ Note that $x$ can be an arbitrarily long sequence of symbols.
   The longer the sequence, and the more alternatives have
   the same prefix, the more useful this transformation is.

▶ What is the effect on the parsers?

# Left factoring – contd.

Consider the grammar:

$S \rightarrow xSy \mid xSx \mid x$

Let us develop a parser for this grammar.

Universiteit Utrecht

# Left factoring – contd.

Consider the grammar:

$$S \rightarrow xSy \mid xSx \mid x$$

Let us develop a parser for this grammar.

After left factoring, we obtain:

$$S \rightarrow xT$$
$$T \rightarrow Sy \mid Sx \mid \varepsilon$$

Universiteit Utrecht

# Left factoring – contd.

Consider the grammar:

$S \rightarrow xSy \mid xSx \mid x$

Let us develop a parser for this grammar.

After left factoring, we obtain:

$S \rightarrow xT$
$T \rightarrow Sy \mid Sx \mid \varepsilon$

Left factoring again:

$S \rightarrow xT$
$T \rightarrow SU \mid \varepsilon$
$U \rightarrow y \mid x$

Universiteit Utrecht

# Left factoring – contd.

Consider the grammar:

$$S \rightarrow xSy \mid xSx \mid x$$

Let us develop a parser for this grammar.

After left factoring, we obtain:

$$S \rightarrow xT$$
$$T \rightarrow Sy \mid Sx \mid \varepsilon$$

Left factoring again:

$$S \rightarrow xT$$
$$T \rightarrow SU \mid \varepsilon$$
$$U \rightarrow y \mid x$$

- ► Left factoring corresponds to an optimization of the parser.
- ► Depending on the grammar and the parser combinators used, it can be absolutely essential.

Universiteit Utrecht

# Left recursion

A production is called **left-recursive** if the right hand side starts with the nonterminal of the left hand side.

Example:

$A \rightarrow Az$

Universiteit Utrecht

# Left recursion

A production is called **left-recursive** if the right hand side starts with the nonterminal of the left hand side.

Example:

$A \to Az$

A grammar is called left-recursive if $A \Rightarrow^+ Az$ for some nonterminal A of the grammar.

Universiteit Utrecht

# Left recursion

A production is called **left-recursive** if the right hand side starts with the nonterminal of the left hand side.

Example:

$A \rightarrow Az$

A grammar is called left-recursive if $A \Rightarrow^{+} Az$ for some nonterminal A of the grammar.

## Question

Can a grammar be left-recursive if it does not have any left-recursive productions?

Universiteit Utrecht

# Left recursion

A production is called **left-recursive** if the right hand side starts with the nonterminal of the left hand side.

Example:

$$A \rightarrow Az$$

A grammar is called left-recursive if $A \Rightarrow^{+} Az$ for some nonterminal A of the grammar.

## Question

Can a grammar be left-recursive if it does not have any left-recursive productions?

Yes, grammars can be indirectly left-recursive.

Universiteit Utrecht

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

# Left recursion and parsers

The production

$A \rightarrow Az$

corresponds to a parser

$a = a <\!\!*\!\!> z$

What happens here?

Universiteit Utrecht

# Left recursion and parsers

The production

$$A \rightarrow Az$$

corresponds to a parser

$$a = a <\!\!*\!\!> z$$

What happens here?

- ▶ Parsers systematically derived from left-recursive grammars loop.
- ▶ Removing left recursion is thus essential if we want a combinator parser.

# Removing left recursion

Transforming a (directly) left-recursive nonterminal A such that the left recursion is removed is relatively simple.

Universiteit Utrecht

# Removing left recursion

Transforming a (directly) left-recursive nonterminal A such that the left recursion is removed is relatively simple.

First, split the productions for A into left-recursive and others:

$$A \rightarrow Ax_1 \mid Ax_2 \mid \ldots \mid A\,x_n$$
$$A \rightarrow y_1 \mid y_2 \mid \ldots \mid y_m \quad \text{(none of the } y_i \text{ start with A)}$$

Universiteit Utrecht

# Removing left recursion

Transforming a (directly) left-recursive nonterminal A such that the left recursion is removed is relatively simple.

First, split the productions for A into left-recursive and others:

$$A \rightarrow Ax_1 \mid Ax_2 \mid \ldots \mid A\, x_n$$
$$A \rightarrow y_1 \mid y_2 \mid \ldots \mid y_m \quad \text{(none of the } y_i \text{ start with A)}$$

This grammar can be transformed to:

$$A \rightarrow y_1 \mid y_1Z \mid y_2 \mid y_2Z \mid \ldots \mid y_m \mid y_mZ$$
$$Z \rightarrow x_1 \mid x_1Z \mid x_2 \mid x_2Z \mid \ldots \mid x_n \mid x_nZ$$

Universiteit Utrecht

# Example: Removing left recursion

Consider:

$$S \rightarrow SS$$
$$S \rightarrow \mathtt{s}$$

One left-recursive production, one other – already split.

Universiteit Utrecht

# Example: Removing left recursion

Consider:

$$S \rightarrow SS$$
$$S \rightarrow \mathtt{s}$$

One left-recursive production, one other – already split.

Applying the transformation yields:

$$S \rightarrow \mathtt{s} \mid \mathtt{s}Z$$
$$Z \rightarrow S \mid SZ$$

Universiteit Utrecht

# 4.3 Separators and operators

# Associative separator/operator

Consider:

| Decls → Decls ; Decls
| Decls → Decl

Universiteit Utrecht

# Associative separator/operator

Consider:

> Decls → Decls ; Decls
> Decls → Decl

This grammar is left-recursive and ambiguous.

Universiteit Utrecht

# Associative separator/operator

Consider:

Decls → Decls ; Decls
Decls → Decl

This grammar is left-recursive and ambiguous.

From a language specification point, one can argue that ambiguity is not problematic if the intended meaning of the different parse trees is the same.

Universiteit Utrecht

# Associative separator/operator

Consider:

Decls → Decls ; Decls
Decls → Decl

This grammar is left-recursive and ambiguous.

From a language specification point, one can argue that ambiguity is not problematic if the intended meaning of the different parse trees is the same.

In this case, the meaning will be the same if ; is intended to be **associative**, i.e., if

$d_1 ; (d_2 ; d_3)$    and    $(d_1 ; d_2) ; d_3$

have the same meaning.

Universiteit Utrecht

# Associative separator/operator – contd.

If the operator is associative, we can freely choose how to remove the ambiguity by choosing either

Decls → Decl ; Decls
Decls → Decl

or

Decls → Decls ; Decl
Decls → Decl

# Associative separator/operator – contd.

If the operator is associative, we can freely choose how to remove the ambiguity by choosing either

> Decls → Decl ; Decls
> Decls → Decl

or

> Decls → Decls ; Decl
> Decls → Decl

Note that the former grammar can be left-factored, and the latter is still left-recursive (but they are no longer ambiguous).

Universiteit Utrecht

# Parsing separated sequences

Separated sequences occur often, so it makes sense to define an abstraction.

Universiteit Utrecht

# Parsing separated sequences

Separated sequences occur often, so it makes sense to define an abstraction.

Left-factoring

Decls → Decl ; Decls
Decls → Decl

yields

Decls → Decl Decls'
Decls' → ; Decls | $\varepsilon$

# Parsing separated sequences

Separated sequences occur often, so it makes sense to define an abstraction.

Left-factoring

Decls → Decl ; Decls
Decls → Decl

yields

Decls → Decl Decls′
Decls′ → ; Decls | $\varepsilon$

We can inline Decls in Decls′:

Decls → Decl Decls′
Decls′ → ; Decl Decls′ | $\varepsilon$

Universiteit Utrecht

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

# Parsing separated sequences – contd.

Decls → Decl Decls′
Decls′ → ; Decl Decls′ | $\varepsilon$

Here Decls′ is just a sequence:

Decls → Decl Decls′
Decls′ → (; Decl)*

Universiteit Utrecht

# Parsing separated sequences – contd.

Decls $\rightarrow$ Decl Decls'
Decls' $\rightarrow$ ; Decl Decls' $|\ \varepsilon$

Here Decls' is just a sequence:

Decls $\rightarrow$ Decl Decls'
Decls' $\rightarrow$ (; Decl)*

We can inline Decls' now:

Decls $\rightarrow$ Decl (; Decl)*
Decls' $\rightarrow$ (; Decl)*

Universiteit Utrecht

# Parsing separated sequences – contd.

Decls → Decl Decls′
Decls′ → ; Decl Decls′ | $\varepsilon$

Here Decls′ is just a sequence:

Decls → Decl Decls′
Decls′ → (; Decl)*

We can inline Decls′ now:

Decls → Decl (; Decl)*
Decls′ → (; Decl)*

and remove Decls′ because it is unreachable:

Decls → Decl (; Decl)*

Universiteit Utrecht

# Parsing separated sequences – contd.

Decls → Decl ( ; Decl)*

Abstracting from Decl and ; , we can construct a parser for separated sequences from this grammar:

listOf :: Parser s a → Parser s b → Parser s [a]
listOf p s = (:) <$> p <*> many (s *> p)

We drop the results from parsing with s and collect the results of the elements in a list.

Universiteit Utrecht

# Parsing separated sequences – contd.

Decls → Decl ( ; Decl)*

Abstracting from Decl and ; , we can construct a parser for separated sequences from this grammar:

listOf :: Parser s a → Parser s b → Parser s [a]
listOf p s = (:) <$> p <*> many (s *> p)

We drop the results from parsing with s and collect the results of the elements in a list.

## Question

What if the separators/operators are supposed to carry meaning?

Universiteit Utrecht

# Operator chains

Consider

E → E O E | Nat
O → + | –

# Operator chains

Consider

$$E \rightarrow E\ O\ E \mid Nat$$
$$O \rightarrow + \mid -$$

Here, we cannot ignore the meaning of the operators because we have to distinguish them.

Universiteit Utrecht

# Operator chains

Consider

$$E \rightarrow E \; O \; E \mid Nat$$
$$O \rightarrow + \mid -$$

Here, we cannot ignore the meaning of the operators because we have to distinguish them.

Also, '−' associates to the left, so we are dealing with a left-recursive grammar.

Universiteit Utrecht

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

# Operator chains

Consider

E → E O E | Nat
O → + | –

Here, we cannot ignore the meaning of the operators because we have to distinguish them.

Also, '–' associates to the left, so we are dealing with a left-recursive grammar.

We inline and remove O and obtain this Haskell datatype:

Universiteit Utrecht

# Operator chains

Consider

$$E \rightarrow E\ O\ E \mid Nat$$
$$O \rightarrow + \mid -$$

Here, we cannot ignore the meaning of the operators because we have to distinguish them.

Also, '−' associates to the left, so we are dealing with a left-recursive grammar.

We inline and remove O and obtain this Haskell datatype:

**data** E = Plus E E | Minus E E | Nat Int

Universiteit Utrecht

# Operator chains – contd.

We would like to parse

> 1+2−3+4

as

> $((\text{Nat } 1 \text{ 'Plus' Nat } 2) \text{ 'Minus' Nat } 3) \text{ 'Plus' Nat } 4$

Universiteit Utrecht

# Operator chains – contd.

We would like to parse

1+2−3+4

as

((Nat 1 'Plus' Nat 2) 'Minus' Nat 3) 'Plus' Nat 4

## Questions

What are the types of the following expressions?

Plus
('Plus' Nat 2)

Universiteit Utrecht

# Operator chains – contd.

We want:

$\big|$ ((Nat 1 'Plus' Nat 2) 'Minus' Nat 3) 'Plus' Nat 4

Universiteit Utrecht

We want:

((Nat 1 'Plus' Nat 2) 'Minus' Nat 3) 'Plus' Nat 4

What does the following evaluate to?

foldl (flip ($)) (Nat 1)
    [('Plus' Nat 2), ('Minus' Nat 3), ('Plus' Nat 4)]

# Operator chains – contd.

We want:

$$((\text{Nat } 1 \text{ 'Plus' Nat } 2) \text{ 'Minus' Nat } 3) \text{ 'Plus' Nat } 4$$

What does the following evaluate to?

$$\text{foldl (flip (\$)) (Nat } 1)$$
$$[(\text{'Plus' Nat } 2), (\text{'Minus' Nat } 3), (\text{'Plus' Nat } 4)]$$

We can obtain this result as follows:

```
chainl :: Parser s a → Parser s (a → a → a) → Parser s a
chainl p s = foldl (flip ($)) <$> p <*> many (flip <$> s <*> p)
```

# Operator chains – contd.

We want:

$((\text{Nat } 1 \text{ ‘Plus‘ Nat } 2) \text{ ‘Minus‘ Nat } 3) \text{ ‘Plus‘ Nat } 4$

What does the following evaluate to?

foldl (flip ($)) (Nat $1$)
    $[(\text{‘Plus‘ Nat } 2), (\text{‘Minus‘ Nat } 3), (\text{‘Plus‘ Nat } 4)]$

We can obtain this result as follows:

```
chainl :: Parser s a → Parser s (a → a → a) → Parser s a
chainl p s = foldl (flip ($)) <$> p <*> many (flip <$> s <*> p)

e = chainl (Nat <$> natural) o
o = Plus <$ symbol '+' <|> Minus <$ symbol '-'
```

Universiteit Utrecht

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

# Chain combinators

There are combinators for left-associative and right-associative chains:

chainl :: Parser s a → Parser s (a → a → a) → Parser s a
chainl p s = foldl (flip ($)) <$> p <*> many (flip <$> s <*> p)

chainr :: Parser s a → Parser s (a → a → a) → Parser s a
chainr p s = foldr ($) <$> many (flip ($) <$> p <*> s) <*> p

# Chain combinators

There are combinators for left-associative and right-associative chains:

chainl :: Parser s a → Parser s (a → a → a) → Parser s a
chainl p s = foldl (flip ($)) <$> p <*> many (flip <$> s <*> p)

chainr :: Parser s a → Parser s (a → a → a) → Parser s a
chainr p s = foldr ($) <$> many (flip ($) <$> p <*> s) <*> p

Universiteit Utrecht

# Chain combinators

There are combinators for left-associative and right-associative chains:

chainl :: Parser s a → Parser s (a → a → a) → Parser s a
chainl p s = foldl (flip ($)) <$> p <*> many (flip <$> s <*> p)
chainr :: Parser s a → Parser s (a → a → a) → Parser s a
chainr p s = foldr ($) <$> many (flip ($) <$> p <*> s) <*> p

The presence of chainl catches one of the most common cases of left recursion in grammars.

# Operator priorities

Consider:

$$E \rightarrow E + E$$
$$E \rightarrow E - E$$
$$E \rightarrow E * E$$
$$E \rightarrow ( \ E \ )$$
$$E \rightarrow Nat$$

Universiteit Utrecht

# Operator priorities

Consider:

$$E \rightarrow E + E$$
$$E \rightarrow E - E$$
$$E \rightarrow E * E$$
$$E \rightarrow ( E )$$
$$E \rightarrow Nat$$

This is a typical grammar for expressions with operators.

For the same reasons as before, it is ambiguous.

Universiteit Utrecht

# Operator priorities

Consider:

$E \rightarrow E + E$
$E \rightarrow E - E$
$E \rightarrow E * E$
$E \rightarrow ( E )$
$E \rightarrow Nat$

This is a typical grammar for expressions with operators.

For the same reasons as before, it is ambiguous.

Given the priorities of the operators and their associativitiy, we can transform this grammar such that the ambiguity is removed.

Universiteit Utrecht

# Operator priorities – contd.

The basic idea is to parse operators of different priorities sequencially.

Universiteit Utrecht

# Operator priorities – contd.

The basic idea is to parse operators of different priorities sequencially.

For each priority level i, we get

$E_i \rightarrow E_i \ Op_i \ E_{i+1} \mid E_{i+1}$     (for left-associative operators)

or

$E_i \rightarrow E_{i+1} \ Op_i \ E_i \mid E_{i+1}$     (for right-associative operators)

or

$E_i \rightarrow E_{i+1} \ Op_i \ E_{i+1} \mid E_{i+1}$     (for non-associative operators)

# Operator priorities – contd.

The basic idea is to parse operators of different priorities sequencially.

For each priority level i, we get

$E_i \rightarrow E_i \ Op_i \ E_{i+1} \mid E_{i+1}$     (for left-associative operators)

or

$E_i \rightarrow E_{i+1} \ Op_i \ E_i \mid E_{i+1}$     (for right-associative operators)

or

$E_i \rightarrow E_{i+1} \ Op_i \ E_{i+1} \mid E_{i+1}$     (for non-associative operators)

The highest level contains the remaining productions.

Universiteit Utrecht

# Operator priorities – contd.

The basic idea is to parse operators of different priorities sequencially.

For each priority level i, we get

$E_i \rightarrow E_i \ Op_i \ E_{i+1} \mid E_{i+1}$      (for left-associative operators)

or

$E_i \rightarrow E_{i+1} \ Op_i \ E_i \mid E_{i+1}$      (for right-associative operators)

or

$E_i \rightarrow E_{i+1} \ Op_i \ E_{i+1} \mid E_{i+1}$     (for non-associative operators)

The highest level contains the remaining productions.

All forms of brackets point to the outer (lowest) level of expressions.

Universiteit Utrecht

# Operator priorities – contd.

Applied to

$$
\begin{aligned}
E &\rightarrow E + E \\
E &\rightarrow E - E \\
E &\rightarrow E * E \\
E &\rightarrow ( \; E \; ) \\
E &\rightarrow Nat
\end{aligned}
$$

we obtain:

$$
\begin{aligned}
E_1 &\rightarrow E_1 \; Op_1 \; E_2 \mid E_2 \\
E_2 &\rightarrow E_2 \; Op_2 \; E_3 \mid E_3 \\
E_3 &\rightarrow ( \; E_1 \; ) \mid Nat \\
Op_1 &\rightarrow + \mid - \\
Op_2 &\rightarrow *
\end{aligned}
$$

# Parsers for operator expressions

Since the abstract syntax trees always make the nesting structure explicit, it typically makes sense to derive the Haskell datatype from the ambiguous grammar:

$E \rightarrow E + E$
$E \rightarrow E - E$
$E \rightarrow E * E$
$E \rightarrow ( E )$
$E \rightarrow Nat$

Universiteit Utrecht

# Parsers for operator expressions

Since the abstract syntax trees always make the nesting structure explicit, it typically makes sense to derive the Haskell datatype from the ambiguous grammar:

| | |
|---|---|
| E → E + E | **data** E = Plus    E E |
| E → E – E |      \| Minus  E E |
| E → E * E |      \| Times  E E |
| E → ( E ) |      \| Parens E |
| E → Nat |      \| Nat |

# Parsers for operator expressions

Since the abstract syntax trees always make the nesting structure explicit, it typically makes sense to derive the Haskell datatype from the ambiguous grammar:

| | |
|---|---|
| E → E + E | **data** E = Plus    E E |
| E → E – E |       \| Minus E E |
| E → E * E |       \| Times E E |
| E → ( E ) | |
| E → Nat |       \| Nat |

# Parsers for operator expressions

Since the abstract syntax trees always make the nesting structure explicit, it typically makes sense to derive the Haskell datatype from the ambiguous grammar:

| | |
|---|---|
| E → E + E | **data** E = Plus   E E |
| E → E – E | &#124; Minus E E |
| E → E * E | &#124; Times E E |
| E → ( E ) | |
| E → Nat | &#124; Nat |

We can now use chainl and chainr again for each of the levels.

Universiteit Utrecht

# Parsers for operator expressions – contd.

$$E_1 \rightarrow E_1\ Op_1\ E_2 \mid E_2$$
$$E_2 \rightarrow E_2\ Op_2\ E_3 \mid E_3$$
$$E_3 \rightarrow (\ E_1\ ) \mid Nat$$
$$Op_1 \rightarrow +\mid -$$
$$Op_2 \rightarrow *$$

```
data E = Plus   E E
       | Minus E E
       | Times E E
       | Nat   Int
```

# Parsers for operator expressions – contd.

$$E_1 \rightarrow E_1 \ Op_1 \ E_2 \mid E_2$$
$$E_2 \rightarrow E_2 \ Op_2 \ E_3 \mid E_3$$
$$E_3 \rightarrow (\ E_1 \ ) \mid Nat$$
$$Op_1 \rightarrow + \mid -$$
$$Op_2 \rightarrow *$$

**data** $E =$ Plus    E E
     | Minus E E
     | Times E E
     | Nat    Int

Parser:

$e_1, e_2, e_3$ :: Parser Char E
$e_1$ = chainl $e_2$ op$_1$
$e_2$ = chainl $e_3$ op$_2$
$e_3$ = parenthesised $e_1 <|>$ Nat $<\$>$ natural

op$_1$, op$_2$ :: Parser Char (E $\rightarrow$ E $\rightarrow$ E)
op$_1$ = Plus    $<\$$ symbol '+' $<|>$ Minus $<\$$ symbol '-'
op$_2$ = Times $<\$$ symbol '*'

Universiteit Utrecht

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

# A general operator parser

We can abstract even further from this pattern:

**type** Op a = (Char, a → a → a)

gen :: [Op a] → Parser Char a → Parser Char a

gen ops p =

   chainl p (choice (map $(\lambda(s, c) \rightarrow$ c $<\$$ symbol s) ops))

where choice combines a list of parsers using $(<|>)$.

# A general operator parser

We can abstract even further from this pattern:

> **type** Op a $=$ (Char, a $\rightarrow$ a $\rightarrow$ a)
> gen :: [Op a] $\rightarrow$ Parser Char a $\rightarrow$ Parser Char a
> gen ops p $=$
>     chainl p (choice (map ($\lambda$(s, c) $\rightarrow$ c $<\!\$$ symbol s) ops))

where choice combines a list of parsers using ($<\!|\!>$).

Now:

> $e_1 =$ gen [('+', Plus), ('-', Minus)] $e_2$
> $e_2 =$ gen [('*', Times)]                $e_3$

Universiteit Utrecht

$$e_1 = \text{gen } [(\text{'+'}, \text{Plus}), (\text{'-'}, \text{Minus})] \; e_2$$
$$e_2 = \text{gen } [(\text{'*'}, \text{Times})] \qquad\qquad e_3$$

Universiteit Utrecht

$$e_1 = \text{gen } [(\text{'+'}, \text{Plus}), (\text{'-'}, \text{Minus})] \; e_2$$
$$e_2 = \text{gen } [(\text{'*'}, \text{Times})] \qquad e_3$$

We do not even need the intermediate levels anymore:

$$e_1 = \text{foldr gen } e_3$$
$$[[(\text{'+'}, \text{Plus}), (\text{'-'}, \text{Minus})], [(\text{'*'}, \text{Times})]]$$

Universiteit Utrecht

# A general operator parser – contd.

$$e_1 = \text{gen } [(\text{'+'}, \text{Plus}), (\text{'-'}, \text{Minus})] \; e_2$$
$$e_2 = \text{gen } [(\text{'*'}, \text{Times})] \qquad\qquad e_3$$

We do not even need the intermediate levels anymore:

$$e_1 = \text{foldr gen } e_3$$
$$\qquad\quad [[(\text{'+'}, \text{Plus}), (\text{'-'}, \text{Minus})], [(\text{'*'}, \text{Times})]]$$

Remarks:

- ▶ Numeric levels not required, just the relative ordering.
- ▶ Extra functionality can be added (such as the possibility of right-associative or unary operators).
- ▶ User-defined abstractions are very useful.

Universiteit Utrecht

# Next lecture

Developing a larger parser from scratch – a case study.

Including common pitfalls and practical problems.

Universiteit Utrecht