

Principled Parsing for Indentation-Sensitive Languages

Revisiting Landin's Offside Rule

Michael D. Adams

Portland State University

Haskell Implementors Workshop, September 14, 2012

Indentation in Haskell

- If a `let`, `where`, `do`, or `of` keyword is not followed by the lexeme `{`, the token `{n}` is inserted after the keyword, where n is the indentation of the next lexeme if there is one, or 0 if the end of file has been reached.
- If the first lexeme of a module is not `{` or `module`, then it is preceded by `{n}` where n is the indentation of the lexeme.
- Where the start of a lexeme is preceded only by white space on the same line, this lexeme is preceded by `<n>`, where n is the indentation of the lexeme, provided that it is not, as a consequence of the first two rules, preceded by `{n}`.
(Haskell Report 2010)

Indentation in Haskell 2010

```
L (<n>:ts) (m:ms) = ';' : (L ts (m:ms))    if m = n
                  = '}' : (L (<n>:ts) ms)    if n < m
L (<n>:ts)      ms = L ts ms
L ({n}:ts) (m:ms) = '{' : (L ts (n:m:ms))  if n > m
L ({n}:ts)      [] = '{' : (L ts [n])       if n > 0
L ({n}:ts)      ms = '{' : '}' : (L (<n>:ts) ms)
L (}')':ts) (0:ms) = '}' : (L ts ms)
L (}')':ts)      ms = parse-error
L ('{:ts)      ms = '{' : (L ts (0:ms))
L ( t :ts) (m:ms) = '}' : (L (t:ts) ms)
                  if m ≠ 0 and parse-error(t)
L ( t :ts)      ms = t : (L ts ms)
L      []      [] = []
L      [] (m:ms) = '}' : L [] ms           if m ≠ 0
```

Magic!

- GLR Parsing
 - Filter out invalid indentations

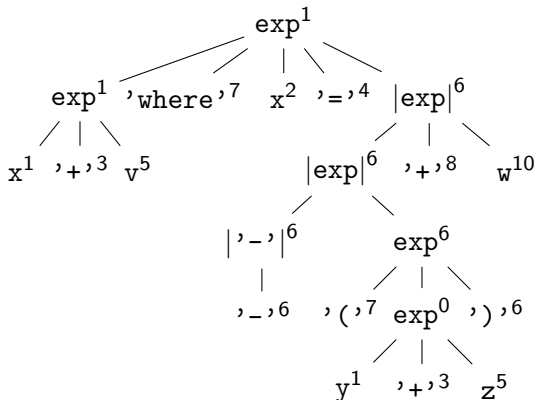
Sebastian Erdweg, Tillmann Rendel, Christian Kästner and Klaus Ostermann. Layout-sensitive Generalized Parsing. In Conference on Software Language Engineering (SLE), 2012.

- Different formalism
- LR Parsing
- LL Parsing
- No filtering
- Formal derivation

$x + v$ where

$x = -($

$y + z) + w$



Grammar

exp \rightarrow exp 'where' ID '=' | exp|

exp \rightarrow exp '+' exp

exp \rightarrow '-' exp

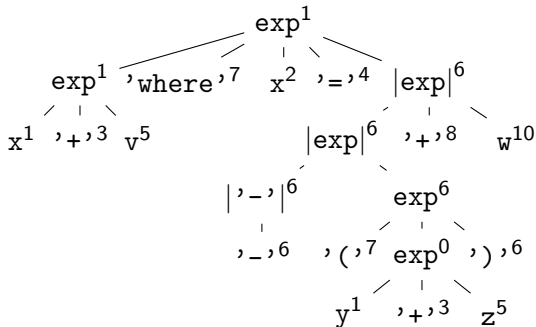
exp \rightarrow '(' exp ')'

exp \rightarrow ID

|exp| \rightarrow |exp| '+' exp

|exp| \rightarrow |'-'| exp

|'-'| \rightarrow '-'



Grammar

$\text{exp} \rightarrow \text{exp}^= \text{'where'}^{\geq} \text{ID}^{\geq} \text{'='}^{\geq} |\text{exp}|^{\geq}$

$\text{exp} \rightarrow \text{exp}^= \text{'+'}^{\geq} \text{exp}^=$

$\text{exp} \rightarrow \text{'-' }^{\geq} \text{exp}^=$

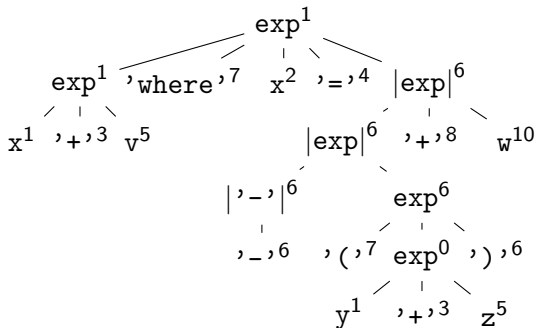
$\text{exp} \rightarrow \text{'(' }^{\geq} \text{exp}^{\otimes} \text{')' }^{\geq}$

$\text{exp} \rightarrow \text{ID}^{\geq}$

$|\text{exp}| \rightarrow |\text{exp}|^= \text{'+'}^{\geq} \text{exp}^=$

$|\text{exp}| \rightarrow |\text{'-' }|^= \text{exp}^=$

$|\text{'-' }| \rightarrow \text{'-' }^=$



Does it work?

Languages

Languages

- SRFI-49

Languages

- SRFI-49
- ISWIM

Languages

- SRFI-49
- ISWIM
- Miranda

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer
- Orwell

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer
- Orwell
- Curry

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer
- Orwell
- Curry
- Habit

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer
- Orwell
- Curry
- Habit
- Occam

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer
- Orwell
- Curry
- Habit
- Occam
- Python

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer
- Orwell
- Curry
- Habit
- Occam
- Python
- Haskell

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer
- Orwell
- Curry
- Habit
- Occam
- Python
- Haskell
- F# (?)

Languages

- SRFI-49
- ISWIM
- Miranda
- Gofer
- Orwell
- Curry
- Habit
- Occam
- Python
- Haskell
- F# (?)
- YAML (?)

- Math (!!!)
 - Lots of subtleties
- Implementation in Happy
 - Based on `haskell-src`
 - Parses all of `base` (with qualifications)
- Under review at POPL 2013
http://michaeldadams.org/papers/layout_parsing/

```
a = do b
      do { c
          d
        }
```



```
a = do b  
      do { c  
    }    d
```

```
a = do b  
      do { c  
    } d
```

```
a = do b
      do { c
}      d
```

```
a = do b
      do { c
} d do e
```

```
a = do b
      do { c
} d
```

```
a = do b
      do { c
}      d
```

```
a = do b
      do { c
} d do e
```

```
a = do b
      do { c
} d
```

```
a = do b
      do { c
} d do e
```

