# FROWN
# An LALR($k$) Parser Generator

RALF HINZE

Institute of Information and Computing Sciences
Utrecht University
Email: `ralf@cs.uu.nl`
Homepage: `http://www.cs.uu.nl/~ralf/`

September, 2001

(Pick the slides at `.../~ralf/talks.html#T30`.)

Joint work with Ross Paterson and Doaitse Swierstra.

# Outline

✖      Usage

✖      Recap: LR parsing

✖      Implementation

✖      Features

✖      Facts and figures

# Running example: well-formed parentheses

$$\textbf{data } \textit{Tree} \quad = \quad \textit{Node } [\textit{Tree}]$$

$$\textit{empty} \quad = \quad \textit{Node } [\,]$$

$$\textit{join } t \; (\textit{Node } us) \quad = \quad \textit{Node } (t : us)$$

$$\%\{$$

$$\textit{Terminal} \quad = \quad \text{'('} \mid \text{')'};$$

$$\textit{Nonterminal} \quad = \quad \textit{Expr}\{\textit{Tree}\};$$

$$\textit{Expr}\{\textit{join } t \; u\} \quad : \quad \text{'('}, \textit{Expr}\{t\}, \text{')'}, \textit{Expr}\{u\};$$

$$\{\textit{empty}\} \quad \mid \quad ;$$

$$\}\%$$

$$\textit{frown } ts \quad = \quad \textit{fail } \texttt{"syntax error"}$$

# Usage

Frown :-( is invoked as follows:

$$frown\ Paren.g$$

This generates a Haskell source file ($Paren.hs$) that contains (among other things) the desired parser:

$$expr\ ::\ (Monad\ m) \Rightarrow [Char] \rightarrow m\ Tree$$

Here, $Char$ is the type of terminals and $Tree$ is the type of semantic values associated with $Expr$.

# The standard example: arithmetic expressions

**type** $Op$ $=$ $Int \rightarrow Int \rightarrow Int$

$\%\{$

$Terminal$ $=$ $Nat\{Int\} \mid Add\{Op\} \mid Mul\{Op\} \mid L \mid R;$

$Nonterminal$ $=$ $Expr\{Int\} \mid Term\{Int\} \mid Factor\{Int\};$

$Expr\{v_1 \ \text{`}op\text{`} \ v_2\}$ $:$ $Expr\{v_1\}, Add\{op\}, Term\{v_2\};$

$\{e\}$ $\mid$ $Term\{e\};$

$Term\{v_1 \ \text{`}op\text{`} \ v_2\}$ $:$ $Term\{v_1\}, Mul\{op\}, Factor\{v_2\};$

$\{e\}$ $\mid$ $Factor\{e\};$

$Factor\{e\}$ $:$ $L, Expr\{e\}, R;$

$\{n\}$ $\mid$ $Nat\{n\};$

$\}\%$

$frown \ ts$ $=$ $fail \ (\texttt{"syntax error: "} \ \mathbin{+\!\!+} \ show \ ts)$

```haskell
data Terminal      =   Nat Int | Add Op | Mul Op | L | R
lexer              ::  String → [ Terminal ]
lexer [ ]          =   [ ]
lexer ( '+' : cs )  =   Add (+) : lexer cs
lexer ( '-' : cs )  =   Add (−) : lexer cs
lexer ( '*' : cs )  =   Mul (∗) : lexer cs
lexer ( '/' : cs )  =   Mul div : lexer cs
lexer ( '(' : cs )  =   L : lexer cs
lexer ( ')' : cs )  =   R : lexer cs
lexer ( c : cs )
      | isDigit c  =   let (n, cs') = span isDigit cs
                       in Nat (read (c : n)) : lexer cs'
      | otherwise  =   lexer cs
```

# Things to note

- The terminal symbols are arbitrary Haskell patterns (of the same *Terminal* type).

- Both terminal and nonterminal symbols may carry multiple semantic values (or no value).

- The parser generated for start symbol $Start\{\,T_1\,\}\dots\{\,T_n\,\}$ has type

$$start \quad :: \quad (Monad\ m) \Rightarrow [\,Terminal\,] \rightarrow m\ (T_1, \dots, T_n)$$

  A grammar may have several start symbols.

# Outline

✔  Usage

✘  Recap: LR parsing

✘  Implementation

✘  Features

✘  Facts and figures

# Shift-reduce parsing

The parsers that are generated by Frown :-( are so-called LALR($k$) parsers ('LA' $\cong$ lookahead, 'L' $\cong$ left-to-right scanning of input, 'R' $\cong$ constructing a rightmost derivation in reverse).

LR parsing is a general method of shift-reduce parsing.

*General idea:* reduce the input string to the start symbol of the grammar. Either *shift* a terminal symbol onto the stack or *reduce* the RHS of a production on top of the stack to the LHS.

*History:* LR parsers were first introduced by Knuth in 1965. When DeRemer devised the LALR method in 1969, the LR technique became the method of choice for parser generators.

# Running example

The grammar is first augmented by an EOF symbol (here '$'$'') and a new start symbol (here '$Start$').

```
%{
Terminal              =  '(' | ')' | '$';
Nonterminal           =  Start{ Tree } | Expr{ Tree };
[0]  Start{ t }       :  Expr{ t }, '$';
[1]  Expr{ join t u } :  '(', Expr{ t }, ')', Expr{ u };
[2]  Expr{ empty }    :  ;
}%
```

# A non-deterministic parser

$$\textbf{data } \textit{Stack} \;=\; \textit{Stack} \succ \textit{Symbol}$$

$$
\begin{aligned}
\textit{parse} &\;::\; (\textit{Monad } m) \Rightarrow \textit{Stack} \rightarrow [\,\textit{Terminal}\,] \rightarrow m \; \textit{Tree} \\
\textit{parse} &\;=\; \textit{shift} \mid \textit{reduce}_0 \mid \textit{reduce}_1 \mid \textit{reduce}_2
\end{aligned}
$$

$$
\begin{aligned}
\textit{shift } st \; (t : tr) &\;=\; \textit{parse} \; (st \succ t) \; tr \\
\textit{reduce}_0 \; (st \succ \textit{Expr}\{\,t\,\} \succ \text{'\$'}) &\;=\; \textit{return } t \\
\textit{reduce}_1 \; (st \succ \text{'('} \succ \textit{Expr}\{\,t\,\} \succ \text{')'} \succ \textit{Expr}\{\,u\,\}) & \\
&\;=\; \textit{parse} \; (st \succ \textit{Expr}\{\,\textit{join } t \; u\,\}) \\
\textit{reduce}_2 \; st &\;=\; \textit{parse} \; (st \succ \textit{Expr}\{\,\textit{empty}\,\})
\end{aligned}
$$

# A sample parse

$$
\begin{array}{lll}
& & \text{()()\$} \\
shift & \text{(} & \text{)()\$} \\
reduce_2 & \text{(}E & \text{)()\$} \\
shift & \text{(}E\text{)} & \text{()\$} \\
shift & \text{(}E\text{)(} & \text{)\$} \\
reduce_2 & \text{(}E\text{)(}E & \text{)\$} \\
shift & \text{(}E\text{)(}E\text{)} & \text{\$} \\
reduce_2 & \text{(}E\text{)(}E\text{)}E & \text{\$} \\
reduce_1 & \text{(}E\text{)}E & \text{\$} \\
reduce_1 & \text{(}E\text{)}E & \text{\$} \\
reduce_1 & E & \text{\$} \\
shift & E\text{\$} & \\
reduce_0 & S &
\end{array}
$$

# Recognition of handles

*Problem:* How can we decide which action to take? In particular, how can we efficiently determine which RHS resides on top of the stack?

☞ This is another language recognition problem!

A *handle* is the RHS of a production preceeded by a left context.

$$S \Longrightarrow_r^* \alpha N \omega \Longrightarrow_r \alpha \beta \omega \Longrightarrow_r^* \omega'$$

Here, $\alpha\beta$ is a handle of the *right-sentential form* $\alpha\beta\omega$ ($\omega$ and $\omega'$ contain only terminals).
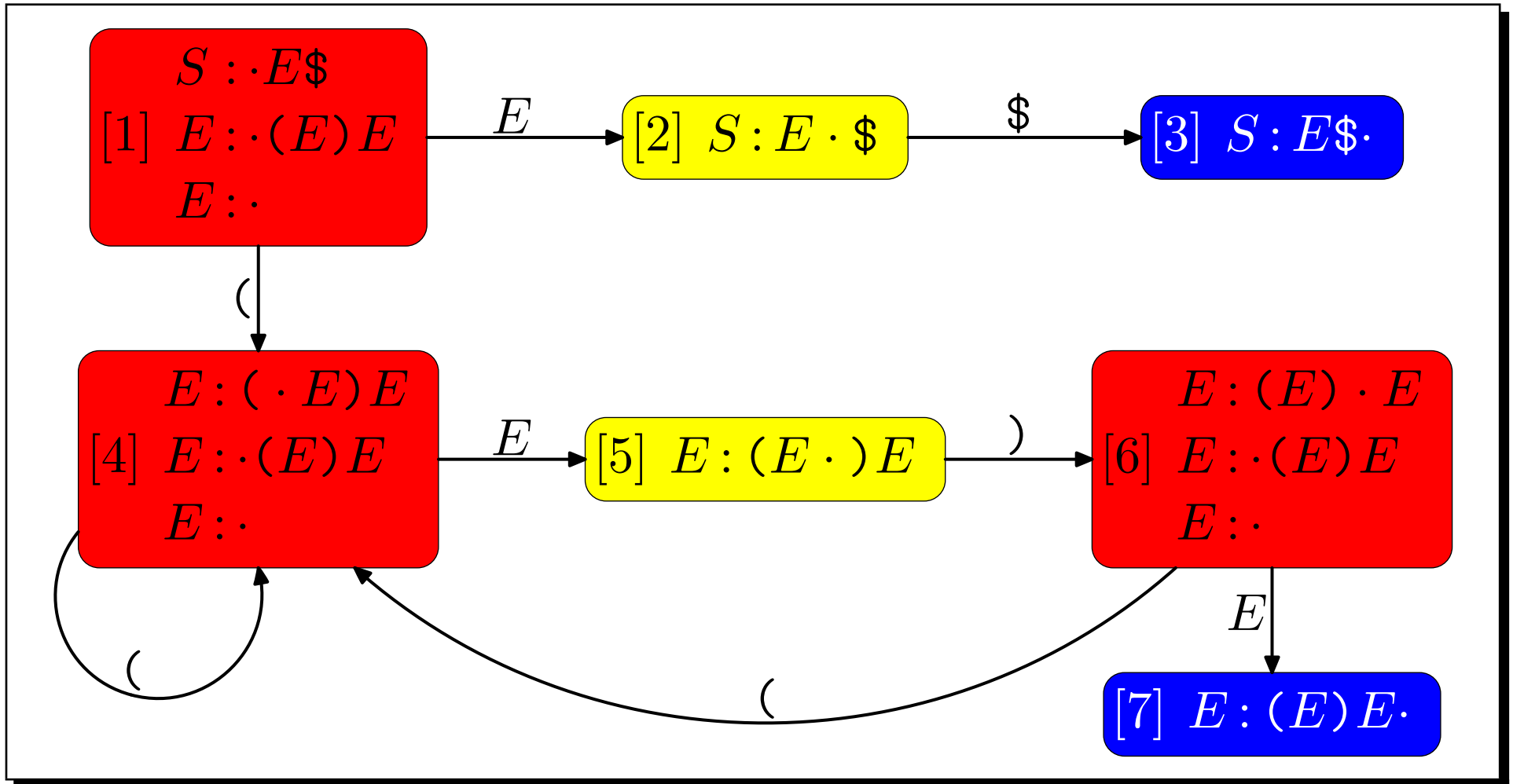
# Grammar of handles and left contexts

$\mathcal{H}(p)$ is the language of handles for production $p$; $\mathcal{L}(N)$ is the language of left contexts of $N$.

$$
\begin{array}{lll}
\mathcal{H}(0) & : & \mathcal{L}(\mathit{Start}), \mathit{Expr}, \text{'\$'}; \\
\mathcal{H}(1) & : & \mathcal{L}(\mathit{Expr}), \text{'('}, \mathit{Expr}, \text{')'}, \mathit{Expr}; \\
\mathcal{H}(2) & : & \mathcal{L}(\mathit{Expr}); \\[4pt]
\mathcal{L}(\mathit{Start}) & : & ; \\
\mathcal{L}(\mathit{Expr}) & : & \mathcal{L}(\mathit{Start}); \\
& | & \mathcal{L}(\mathit{Expr}), \text{'('}; \\
& | & \mathcal{L}(\mathit{Expr}), \text{'('}, \mathit{Expr}, \text{')'};
\end{array}
$$

☞ By construction, this grammar is left-recursive. Thus, $\mathcal{H}(p)$ and $\mathcal{L}(N)$ generate regular languages!

# LR(0) automaton

Regular languages can be recognized by *deterministic finite state machines* (set of states construction; $N : \alpha \cdot \beta$ is called an *item*).

# Outline

✔     Usage

✔     Recap: LR parsing

✘     Implementation

✘     Features

✘     Facts and figures

# Towards Haskell

The LR(0) automaton can be directly encoded as a functional program.

The stack records the transitions of the LR(0) automaton.

$$
\begin{aligned}
\textbf{data } GStack \;&=\; VStack \to State \\
\textbf{data } VStack \;&=\; GStack \succ Symbol
\end{aligned}
$$

**NB.** $s \succ t \to s'$ is meant to resemble the transition $s \xrightarrow{\;t\;} s'$.

# Shift and reduce states

State 2 is a *shift state*.

$$parse_2 \; st \; (\text{'\$'} : tr) \;\; = \;\; parse_3 \; (st \succ \text{'\$'} \to 3) \; tr$$

State 7 is a *reduce state*.

$$parse_7 \; (st \to 1 \succ \text{'('} \to 4 \succ E\{t\} \to 5 \succ \text{')'} \to 6 \succ E\{u\} \to 7)$$
$$= \;\; parse_2 \; (st \succ 1 \to E\{join \; t \; u\} \to 2)$$
$$parse_7 \; (st \to 4 \succ \text{'('} \to 4 \succ E\{t\} \to 5 \succ \text{')'} \to 6 \succ E\{u\} \to 7)$$
$$= \;\; parse_5 \; (st \succ 4 \to E\{join \; t \; u\} \to 5)$$
$$parse_7 \; (st \to 6 \succ \text{'('} \to 4 \succ E\{t\} \to 5 \succ \text{')'} \to 6 \succ E\{u\} \to 7)$$
$$= \;\; parse_7 \; (st \succ 6 \to E\{join \; t \; u\} \to 7)$$

# Conflicts

We have a *shift/reduce conflict* if a state contains both a shift and a reduce action (states 1, 4, and 6 in our running example).

☞ A shift/reduce conflict can possibly be resolved using *one* token of lookahead.

We have a *reduce/reduce conflict* if a state contains several reduce actions.

☞ A reduce/reduce conflict can possibly be resolved using $k$ tokens of lookahead.

# Computation of lookahead information

*Idea:* partially execute the LR(0) machine at compile time to determine the shifts that might follow a reduce action.

| | | |
|---|---|---|
| [1] | $E : 1 \cdot$ | $\{\$\}$ |
| [3] | $S : 1\,E\,2\,\$\,3 \cdot$ | $\{\}$ |
| [4] | $E : 4 \cdot$ | $\{\,)\,\}$ |
| [6] | $E : 6 \cdot$ | $\{\$,\,)\,\}$ |
| [7] | $E : 1\,(\,4\,E\,5\,)\,6\,E\,7 \cdot$ | $\{\$\}$ |
| | $E : 4\,(\,4\,E\,5\,)\,6\,E\,7 \cdot$ | $\{\,)\,\}$ |
| | $E : 6\,(\,4\,E\,5\,)\,6\,E\,7 \cdot$ | $\{\$,\,)\,\}$ |

**NB.** LALR($k$) parsers merge the lookahead information for each production. In other words, the parsers generated by Frown :-(  are slightly more general than LALR.

# In Haskell: representation of the stack

For each transition we introduce a constructor.

```
data Stack  =  Empty
            |  St_1_4 Stack
            |  St_1_2 Stack (Tree)
            |  St_2_3 Stack
            |  St_4_4 Stack
            |  St_4_5 Stack (Tree)
            |  St_5_6 Stack
            |  St_6_4 Stack
            |  St_6_7 Stack (Tree)
```

# In Haskell: LR(0) machine

$$expr \ tr \qquad\qquad\qquad\qquad = \ parse\_1 \ tr \ Empty$$

$$parse\_1 \ ts@[\,] \ st \qquad\qquad = \ parse\_2 \ ts \ (St\_1\_2 \ st \ (empty))$$

$$parse\_1 \ (\text{'}(\text{'} : tr) \ st \qquad\quad = \ parse\_4 \ tr \ (St\_1\_4 \ st)$$

$$parse\_1 \ ts \ st \qquad\qquad\quad = \ frown \ ts$$

$$parse\_2 \ tr@[\,] \ st \qquad\qquad = \ parse\_3 \ tr \ (St\_2\_3 \ st)$$

$$parse\_2 \ ts \ st \qquad\qquad\quad = \ frown \ ts$$

$$parse\_3 \ ts \ (St\_2\_3 \ (St\_1\_2 \ st \ (v0)))$$
$$\qquad\qquad\qquad\qquad\quad = \ return \ (v0)$$

$$parse\_4 \ (\text{'}(\text{'} : tr) \ st \qquad\quad = \ parse\_4 \ tr \ (St\_4\_4 \ st)$$

$$parse\_4 \ ts@(\text{'})\text{'} : tr) \ st \ = \ parse\_5 \ ts \ (St\_4\_5 \ st \ (empty))$$

$$parse\_4 \ ts \ st \qquad\qquad\quad = \ frown \ ts$$

$$parse\_5\ ('\!)'\!:tr)\ st\ =\ parse\_6\ tr\ (St\_5\_6\ st)$$

$$parse\_5\ ts\ st\ =\ frown\ ts$$

$$parse\_6\ ('\!('\!:tr)\ st\ =\ parse\_4\ tr\ (St\_6\_4\ st)$$

$$parse\_6\ ts\ st\ =\ parse\_7\ ts\ (St\_6\_7\ st\ (empty))$$

$$parse\_7\ ts\ (St\_6\_7\ (St\_5\_6\ (St\_4\_5\ (St\_1\_4\ st)\ (t)))\ (u))$$
$$=\ parse\_2\ ts\ (St\_1\_2\ st\ (join\ t\ u))$$

$$parse\_7\ ts\ (St\_6\_7\ (St\_5\_6\ (St\_4\_5\ (St\_4\_4\ st)\ (t)))\ (u))$$
$$=\ parse\_5\ ts\ (St\_4\_5\ st\ (join\ t\ u))$$

$$parse\_7\ ts\ (St\_6\_7\ (St\_5\_6\ (St\_4\_5\ (St\_6\_4\ st)\ (t)))\ (u))$$
$$=\ parse\_7\ ts\ (St\_6\_7\ st\ (join\ t\ u))$$

# Outline

✔      Usage

✔      Recap: LR parsing

✔      Implementation

✘      Features

✘      Facts and figures

# Features

- Multiple entry points (start symbols).

- Multiple attribute values.

- Precedences to resolve conflicts.

- `--lookahead=`$k$: Use additional lookahead to resolve reduce/reduce conflicts (only used where needed).

- `--backtrack`: Produce a backtracking parser.

$$start \ :: \ (MonadPlus \ m) \Rightarrow [\,Terminal\,] \rightarrow m \ (T_1, \ldots, T_n)$$

- `--monadic`: Monadic semantic actions.

- `--lexer`: Use monadic lexer instead of list of terminal symbols:

$$
\begin{aligned}
get &\quad :: \quad (Monad\ m) \Rightarrow m\ Terminal \\
start &\quad :: \quad (Monad\ m) \Rightarrow m\ (T_1, \dots, T_n)
\end{aligned}
$$

- `--expected`: In case of error pass the set of expected tokens to the error routine.

$$
frown \quad :: \quad (Monad\ m) \Rightarrow [\,Terminal\,] \rightarrow [\,Terminal\,] \rightarrow m\ a
$$

- Four different parser schemes (standard, LALR-like, stackless, combinator-based).

# Outline

✔      Usage

✔      Recap: LR parsing

✔      Implementation

✔      Features

✘      Facts and figures

# Facts and figures: expression parser

The expression grammar has 5 terminals, 4 nonterminals, and 7 productions. The LR automaton has 13 states and 23 transitions.

Running time $(expr\ (\texttt{"a+b"} + concat\ (replicate\ n\ \texttt{"*c+d"})))$.

|               | 1.000 | 10.000 | 100.000 |
|---------------|-------|--------|---------|
| `happy`       | 0.01  | 0.29   | 5.20    |
| `happy -a`    | 0.02  | 0.40   | 5.92    |
| `happy -c -g` | 0.01  | 0.25   | 4.44    |
| `happy -a -c -g` | 0.01 | 0.23  | 3.85    |
| `frown`       | 0.01  | 0.10   | 1.98    |
| `frown -oc`   | 0.01  | 0.11   | 2.01    |
| `frown -s`    | 0.01  | 0.10   | 2.03    |

# Facts and figures: Haskell parser

The Haskell grammar has 61 terminals, 121 nonterminals, and 277 productions. The LR automaton has 490 states and 2961 transitions.

Compilation time.

|              | time | .hs   | space |         |
|--------------|------|-------|-------|---------|
| happy        |      | 210K  | 150M  | -fno-cpr |
| happy -c -g  | 4.5  | 246K  | 180M  | -fno-cpr |
| happy -a -c -g | 4.5 | 164K  | 100M  |         |
| frown -oc    | 4.8  | 300K  | 160M  |         |

|              | .o     | a.out  | stripped |
|--------------|--------|--------|----------|
| happy        | 1.511K | 2.536K | 1.246K   |
| happy -c -g  | 1.093K | 2.202K | 1.134K   |
| happy -a -c -g | 346K | 1.711K | 894K     |
| frown -oc    | 1.489K | 2.547K | 1.362K   |

Running time.

| | 18K | 300K | 1653K |
|---|---|---|---|
| lex only | 0.01 | 0.55 | 4.8 |
| happy | 0.23 | 2.68 | 12.1 |
| happy -c -g | 0.14 | 2.23 | 11.6 |
| happy -a -c -g | 0.07 | 1.83 | 10.6 |
| frown -oc | 0.14 | 1.95 | 8.5 |

# Outline

✔      Usage

✔      Recap: LR parsing

✔      Implementation

✔      Features

✔      Facts and figures

# Conclusion

- Encoding the LR(0) automaton as a set of mutually recursive functions rather than as a huge table gives more flexibility (use lookahead only where it is needed).

- Frown :-( generates good code (huge but fast). The output is fairly readable.

- With some additional work the generated parsers can produce good error messages (monadic lexer that keeps track of line numbers).