



Parse Tools

Copyright © 1997-2020 Ericsson AB. All Rights Reserved.
Parse Tools 2.2
June 21, 2020

Copyright © 1997-2020 Ericsson AB. All Rights Reserved.

Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file except in compliance with the License. You may obtain a copy of the License at <http://www.apache.org/licenses/LICENSE-2.0> Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License. Ericsson AB. All Rights Reserved..

June 21, 2020

1 Reference Manual

The **Parsetools** application contains utilities for parsing and scanning. Yecc is an LALR-1 parser generator for Erlang, similar to yacc. Yecc takes a BNF grammar definition as input, and produces Erlang code for a parser as output. Leex is a regular expression based lexical analyzer generator for Erlang, similar to lex or flex.


```
myparser:parse_and_scan({Function, Args})
myparser:parse_and_scan({Mod, Tokenizer, Args})
```

The tokenizer Function is either a fun or a tuple {Mod, Tokenizer}. The call apply(Function, Args) or apply({Mod, Tokenizer}, Args) is executed whenever a new token is needed. This, for example, makes it possible to parse from a file, token by token.

The tokenizer used above has to be implemented so as to return one of the following:

```
{ok, Tokens, Endline}
{eof, Endline}
{error, Error_description, Endline}
```

This conforms to the format used by the scanner in the Erlang io library module.

If {eof, Endline} is returned immediately, the call to parse_and_scan/1 returns {ok, eof}. If {eof, Endline} is returned before the parser expects end of input, parse_and_scan/1 will, of course, return an error message (see above). Otherwise {ok, Result} is returned.

More Examples

1. A grammar for parsing infix arithmetic expressions into prefix notation, without operator precedence:

```
Nonterminals E T F.
Terminals '+' '*' '(' ')' number.
Rootsymbol E.
E -> E '+' T: {'$2', '$1', '$3'}.
E -> T: '$1'.
T -> T '*' F: {'$2', '$1', '$3'}.
T -> F: '$1'.
F -> '(' E ')': '$2'.
F -> number: '$1'.
```

2. The same with operator precedence becomes simpler:

```
Nonterminals E.
Terminals '+' '*' '(' ')' number.
Rootsymbol E.
Left 100 '+'.
Left 200 '*'.
E -> E '+' E: {'$2', '$1', '$3'}.
E -> E '*' E: {'$2', '$1', '$3'}.
E -> '(' E ')': '$2'.
E -> number: '$1'.
```

3. An overloaded minus operator:

```
Nonterminals E uminus.
Terminals '*' '-' number.
Rootsymbol E.

Left 100 '-'.
Left 200 '*'.
Unary 300 uminus.

E -> E '-' E.
E -> E '*' E.
E -> uminus.
E -> number.

uminus -> '-' E.
```

4. The Yecc grammar that is used for parsing grammar files, including itself:

```
Nonterminals
grammar declaration rule head symbol symbols attached_code
token tokens.
Terminals
atom float integer reserved_symbol reserved_word string char var
'-'> ':' dot.
Rootsymbol grammar.
Endsymbol 'end'.
grammar -> declaration : '$1'.
grammar -> rule : '$1'.
declaration -> symbol symbols dot: {'$1', '$2'}.
rule -> head '-'> symbols attached_code dot: {rule, ['$1' | '$3'],
'$4'}.
head -> symbol : '$1'.
symbols -> symbol : ['$1'].
symbols -> symbol symbols : ['$1' | '$2'].
attached_code -> ':' tokens : {erlang_code, '$2'}.
attached_code -> '$empty' : {erlang_code,
[{atom, 0, '$undefined'}]}.
tokens -> token : ['$1'].
tokens -> token tokens : ['$1' | '$2'].
symbol -> var : value_of('$1').
symbol -> atom : value_of('$1').
symbol -> integer : value_of('$1').
symbol -> reserved_word : value_of('$1').
token -> var : '$1'.
token -> atom : '$1'.
token -> float : '$1'.
token -> integer : '$1'.
token -> string : '$1'.
token -> char : '$1'.
token -> reserved_symbol : {value_of('$1'), line_of('$1')}.
token -> reserved_word : {value_of('$1'), line_of('$1')}.
token -> '-'> : {'-'>', line_of('$1')}.
token -> ':' : {':', line_of('$1')}.
Erlang code.
value_of(Token) ->
    element(3, Token).
line_of(Token) ->
    element(2, Token).
```

Note:

The symbols '-'>', and ':' have to be treated in a special way, as they are meta symbols of the grammar notation, as well as terminal symbols of the Yecc grammar.

5. The file `erl_parse.yrl` in the `lib/stdlib/src` directory contains the grammar for Erlang.

Note:

Syntactic tests are used in the code associated with some rules, and an error is thrown (and caught by the generated parser to produce an error message) when a test fails. The same effect can be achieved with a call to `return_error(Error_line, Message_string)`, which is defined in the `yeccpre.hrl` default header file.

Files

```
lib/parsetools/include/yeccpre.hrl
```

See Also

Aho & Johnson: 'LR Parsing', ACM Computing Surveys, vol. 6:2, 1974.

leex

Erlang module

A regular expression based lexical analyzer generator for Erlang, similar to lex or flex.

Note:

The Leex module should be considered experimental as it will be subject to changes in future releases.

DATA TYPES

```
ErrorInfo = {ErrorLine,module(),error_descriptor()}
ErrorLine = integer()
Token = tuple()
```

Exports

`file(FileName) -> LeexRet`

`file(FileName, Options) -> LeexRet`

Types:

```
FileName = filename()
Options = Option | [Option]
Option = - see below -
LeexRet = {ok, Scannerfile} | {ok, Scannerfile, Warnings} | error |
{error, Errors, Warnings}
Scannerfile = filename()
Warnings = Errors = [{filename(), [ErrorInfo]}]
ErrorInfo = {ErrorLine, module(), Reason}
ErrorLine = integer()
Reason = - formatable by format_error/1 -
```

Generates a lexical analyzer from the definition in the input file. The input file has the extension `.xrl`. This is added to the filename if it is not given. The resulting module is the Xrl filename without the `.xrl` extension.

The current options are:

`dfa_graph`

Generates a `.dot` file which contains a description of the DFA in a format which can be viewed with Graphviz, www.graphviz.com.

`{includefile, Includefile}`

Uses a specific or customised prologue file instead of default `lib/parsetools/include/leexinc.hrl` which is otherwise included.

`{report_errors, bool()}`

Causes errors to be printed as they occur. Default is `true`.

`{report_warnings, bool()}`

Causes warnings to be printed as they occur. Default is `true`.

`\c`

Any other character literally, for example `\\` for backslash, `\ "` for `"`.

The following examples define simplified versions of a few Erlang data types:

```
Atoms [a-z][0-9a-zA-Z_]*
```

```
Variables [A-Z_][0-9a-zA-Z_]*
```

```
Floats (\+|-)?[0-9]+\.[0-9]+((E|e)(\+|-)?[0-9]+)?
```

Note:

Anchoring a regular expression with `^` and `$` is not implemented in the current version of Leex and just generates a parse error.