

Parsing

Overview

- **parser:** program extracting structure from linear sequence of elements
 - e.g. transforming string "3+4*5" to a tree representing the expression
- **domain specific language (DSL):** small programming language for a narrow domain
 - often embedded in existing languages, adding new features particular to the domain, while otherwise using existing functionality
- if DSL can be parsed by extending host language parse, it's much more convenient to use
- Prolog handles this well:
 - `read/1` reads a term
 - `op/3` allows you to extend the language by defining new operators

Operator Precedence

- **operator precedence:** simple parsing technique based on operator's:
 - **precedence:** which operator binds tightest
 - **associativity:** if repeated infix operators associate to left/right/both
* i.e. $a - b - c$ is $(a - b) - c$ or $a - (b - c)$ or an error
 - **fixity:** infix, prefix, postfix

Prolog Operators

- Prolog's `op/3` predicate declares an operator
 - precedence:
 - * larger numbers are lower precedence
 - * 1000: goal precedence
 - fixity: 2/3 letter symbol giving fixity and associativity
 - * `f`: operator
 - * `x`: subterm at lower precedence
 - * `y`: subterm at higher precedence
 - operator: operator to declare

```
1 :- op(precedence, fixity, operator).
```

e.g. Prolog imperative for loop

```

1 :- op(950, fx, for).
2 :- op(940, xfx, in).
3 :- op(600, xfx, '..').
4 :- op(1050, xfy, do).

5
6 for Generator do Body :- 
7     (   call(Generator),
8         call(Body),
9         fail
10    ;   true
11    ).
12
13 Var in Low .. High :- 
14     between(Low, High, Var).
15
16 Var in [H|T] :- 
17     member(Var, [H|T]).
```

Haskell Operators

- simpler and more limited than Prolog
- only supports infix operators
- declare as `associativity` `precedence` `operator`
- associativity can be:
 - `infixl`: left associative infix operator
 - `infixr`: right associative infix operator
 - `infix`: non-associative infix operator
- precedence: integer 1-9
 - lower numbers are lower precedence (looser)

e.g. define `%` as synonym for `mod`

```

1 infix 7 %
2
3 (%) :: Integral a => a -> a -> a
4 a % b = a `mod` b
```

Grammars

- parsing is based on a **grammar** which specifies the language to be parsed
- **terminals:** symbols of the language
- **non-terminals:** specify a linguistic category
- grammar comprised of set of rules

$$(\text{non-terminal} \cup \text{terminal})^* \rightarrow (\text{non-terminal} \cup \text{terminal})^*$$

- most commonly, LHS of arrow is a single non-terminal:

expression → expression' + expression
 expression → expression' - expression
 expression → expression'* expression

Definite Clause Grammars

- Prolog directly supports **definite clause grammars**, which adhere to the following rules:
 - Non-terminals are written using goal-like syntax
 - Terminals are written between single quotes
 - LHS and RHS separated with -->
 - parts on RHS separated with ,
 - empty terminal: [] or ''
- e.g. expression grammar as Prolog DCG:

```

1 expr --> expr, '+', expr.
2 expr --> expr, '*', expr.
3 expr --> expr, '-', expr.
4 expr --> expr, '/', expr.
5 expr --> number.
```

- note this can only test whether a given string is an element of the language
- to produce a **parse tree**, i.e. a data structure representing the input, add arguments to the non-terminals

```

1 expr(E1+E2) --> expr(E1), '+', expr(E2).
2 expr(E1*E2) --> expr(E1), '*', expr(E2).
3 expr(E1-E2) --> expr(E1), '-', expr(E2).
4 expr(E1/E2) --> expr(E1), '/', expr(E2).
5 expr(N) --> number(N).
```

Recursive Descent Parsing

- **recursive descent parsing:** DCGs map each non-terminal to a predicate that nondeterministically parses one instance of that non-terminal
- to use a grammar, you use the `phrase/2` predicate: `phrase(nonterminal, string)`.
- recursive descent parsing cannot handle left recursion

Left Recursion

- `expr(E1+E2) --> expr(E1), '+', expr(E2)`. is left recursive
 - before parsing any terminals, it calls itself recursively
 - as DCGs are transformed to ordinary Prolog code, this becomes a clause that calls itself recursively consuming no input: infinite recursion
- DCGs can be transformed to remove left recursion:
 - rename left recursive rules to `A_rest` and remove the first non-terminal
 - add a rule for `A_rest` matching empty input
 - add `A_rest` to the end of the non-left recursive rules
- DCGs with arguments: non-left recursive rules
 - replace argument of non-left recursive rules with a fresh variable
 - use original argument of `_rest` added non-terminal
 - add fresh variable as second argument of `_rest` added non-terminal e.g.

```

1 expr(N) --> number(N).
2 % becomes
3 expr(E) --> number(N), expr_rest(N, E).
```

- DCGs with arguments: left recursive rules
 - use argument of left-recursive non-terminal as first head argument, and fresh variable as second
 - use original argument of head as first argument of `_tail` call, and fresh variable as second argument of head and `_tail` call

```

1 expr(E1+E2) --> expr(E1), '+', expr(E2).
2 % becomes
3 expr_rest(E1,R) --> '+', expr(E2), expr_rest(E1+E2, R).
```

Disambiguating a grammar

- original grammar is ambiguous: `expr(E1-E2) --> expr(E1), '-' , expr(E2)`.
 - applied to “3-4-5” allows E1 to be “3-4” or “4-5”
- ensure only desired one is possible by splitting ambiguous non-terminal into separate non-terminals for each precedence level
- becomes (before elimination of left recursion)

```
1 expr(E-F) --> expr(E), '-' factor(F)
```

Final Grammar

```
1 expr(E) --> factor(F), expr_rest(F, E).
2
3 expr_rest(F1, E) --> '+', factor(F2), expr_rest(F1+F2, E).
4 expr_rest(F1, E) --> '-', factor(F2), expr_rest(F1-F2, E).
5 expr_rest(F, F) --> [].
6
7 factor(F) --> number(N), factory_rest(N,F).
8
9 factor_rest(N1, F) --> '*', number(N2), factor_rest(N1*N2, F).
10 factor_rest(N1, F) --> '/', number(N2), factor_rest(N1/N2, F).
11 factor_rest(N, N) --> [].
```

Tokenisers

- syntax analysis = lexical analysis/tokenising + parsing**
- lexical analysis:** uses simpler class of grammar to group characters and tokens
 - eliminates meaningless text (whitespace, comments)
- you can use '`strings`' as terminals or lists if you need to
- you can also write normal Prolog code in a DCG wrapped in { }
- if it fails, the rule fails

```
1 number(N) -->
2   [C],
3   { '0' =< C, C =< '9' },
4   { N0 is C - '0' },
5   number_rest(N0, N).
6
7 number_rest(N0, N) -->
```

```
8      ( [C],
9       { '0' =< C, C =< '9' }
10      -> { N1 is N0 *10 + C - '0' },
11          number_rest(N1, N),
12          ; { N = N0}
13      ).
```

Working parser

```
1 ?- phrase(expr(E), '3+4*5'), Value is E.
2 E = 3+4*5,
3 Value = 23;
4 false.
```

Extras

- DCGs can run backwards to generate text from structure

```
1 flatten(empty) --> []
2 flatten(node(L, E, R)) -->
3   flatten(L),
4   [E],
5   flatten(R).
```

- parsing in Haskell
 - [ReadP](#), [Read](#), [Parsec](#)