

**Idee:** (Pseudocode)

$x :: t$  bedeutet  $x$  hat Typ  $t$ .

$x :: a \rightarrow b$  bedeutet  $x$  ist Funktion mit Eingabetyp  $a$  und Ausgabotyp  $b$ .

Wir definieren weiter den Typ  $\text{Parser} = \{\text{String} \rightarrow \{\text{String}\}\}$ .

Ein Parser ist also eine Funktion von einer Menge von Strings zu einer Menge von Strings. Später werden die folgenden Funktionen auf Indizes von Matches ( $\text{startIndex}$ ,  $\text{endIndex}$ ) laufen, jedoch ist zur Veranschaulichung dies besser.

**Definition:**  $\text{readChar}$

$\text{readChar} :: \text{Char} \rightarrow \text{Parser}$

$\text{readChar}(c) = S \rightarrow \{w_1, \dots, w_n \mid w \in S \wedge w_0 = c\}$

$\text{readChar}(c)$  gibt also einen Parser(Funktion) zurück, der den char  $c$  konsumiert.

**Beispiel:**  $\text{readChar}$

Mit  $f_a := \text{readChar}(a)$  haben wir

$$f_a(\{\text{ab}, \text{aac}, \text{xb}, \text{b}\}) = \{\text{b}, \text{ac}\}$$

**Definition:**  $\text{concatenate}$

$\text{concatenate} :: (\text{Parser}, \text{Parser}) \rightarrow \text{Parser}$

$\text{concatenate}(p_1, p_2) = p_2 \circ p_1$

**Beispiel:**  $\text{concatenate}$

Mit  $f_a$  wie oben und  $f_b$  analog definieren wir  $f_{ab} := \text{concatenate}(f_a, f_b) = f_b \circ f_a$ .

Dann haben wir beispielsweise

$$f_{ab}(\{\text{abc}, \text{aa}, \text{bb}, \text{xbz}\}) = f_b(f_a(\{\text{abc}, \text{aa}, \text{bb}, \text{xbz}\})) = f_b(\{\text{bc}, \text{a}\}) = \{\text{c}\}$$

**Definition:**  $\text{alternate}$

$\text{alternate} :: (\text{Parser}, \text{Parser}) \rightarrow \text{Parser}$

$\text{alternate}(p_1, p_2) = S \rightarrow p_1(S) \cup p_2(S)$

**Beispiel:**  $\text{alternate}$

Mit  $f_a, f_b$  wie oben definieren wir  $f_{a|b} := \text{alternate}(f_a, f_b)$ . Dann haben wir beispielsweise:

$$\begin{aligned} f_{a|b}(\{\text{abc}, \text{by}, \text{xbz}\}) &= f_a(\{\text{abc}, \text{by}, \text{xbz}\}) \cup f_b(\{\text{abc}, \text{by}, \text{xbz}\}) \\ &= \{\text{bc}\} \cup \{\text{y}\} = \{\text{bc}, \text{y}\} \end{aligned}$$

**Definition:**  $\text{iterate}$

$\text{iterate} :: \text{Parser} \rightarrow \text{Parser}$

$\text{iterate}(p) = S \rightarrow \bigcup_{i \in \mathbb{N}_0} p^i(S) = S \rightarrow S \cup p(S) \cup p(p(S)) \cup \dots$

**Beispiel:**  $\text{iterate}$

Mit  $f_a$  wie oben definieren wir  $f_{a*} := \text{iterate}(f_a)$ . Dann haben wir beispielsweise:

$$\begin{aligned} &f_{a*}(\{\text{aaaz}, \text{axa}, \text{edc}\}) \\ &= \{\text{aaaz}, \text{axa}, \text{edc}\} \cup f_a(\{\text{aaaz}, \text{axa}, \text{edc}\}) \cup f_a(f_a(\{\text{aaaz}, \text{axa}, \text{edc}\})) \cup \dots \\ &= \{\text{aaaz}, \text{axa}, \text{edc}\} \cup \{\text{aaz}, \text{xa}\} \cup \{\text{az}\} \cup \{\text{z}\} \cup \emptyset \cup \emptyset \cup \dots \\ &= \{\text{aaaz}, \text{axa}, \text{edc}, \text{aaz}, \text{xa}, \text{az}, \text{z}\} \end{aligned}$$

In dieser Weise lassen sich Induktiv Parser zusammenbasteln, z.B. mit dem Thompsonvisitor.

```
concatenate(
  iterate(
    alternate(
      readChar('a'),
      concatenate(
        readChar('b'),
        readChar('b')
      )
    )
  ),
  readChar('b')
)
```

Dies entspricht der Regex  $(a \mid bb)^*b$ .

Um jetzt aber vernünftig matchen zu können arbeiten wir nicht mehr mit Strings sondern Indizes auf einem globalen gemeinsamen Eingabestring.

Wir haben also beispielsweise den Eingabestring `bbaaabbb` und Paare von `start` und `endIndex` als Matches der obigen Regex wie folgt:

(0, 1) für <code> b baaabbb</code> , wo also der Kleene-Stern 0 mal gematcht wird und dann das <code>b</code> .		
(0, 6) für <code> bbaaab bb</code>	(0, 8) für <code> bbaaabbb </code>	(1, 2) für <code>b b aaabbb</code>
(1, 6) für <code>b baaab bb</code>	(1, 8) für <code>b baaabbb </code>	(2, 6) für <code>bb aaab bb ...</code>

Die Definitionen sehen dann wie folgt aus:

Für ein gegebenes Eingabewort  $w = w_0w_1 \dots w_n$  haben wir

```
type Match = (Int, Int)      type Parser = {Match} -> {Match}
```

```
readChar(c) = S → {(i, j + 1) | (i, j) ∈ S ∧ wj = c}
```

Die restlichen Funktionen behalten ihre Definition.

Wenn wir also einen Parser `p` aus diesen Funktionen zusammengebaut haben, so können wir unsere Startmenge von Matches durch  $S := \{(i, i) \mid i \in [0, |w| - 1]\}$  definieren, wobei  $w$  das Eingabewort ist. Dies symbolisiert alle möglichen Startposition von Matches im gegebenen String. So können wir alle möglichen Matches parallel behandeln. Weiter haben wir also eine Funktion `runParser :: (Parser, String) -> {Match}` welche die Startmenge  $S$  erzeugt und diese dem Parser übergibt.

### Beispiel

Regex ist  $(a|b)c^*$ . Eingabewort  $w = xabccx$ , der Parser lässt sich bauen durch

```
p = concatenate(
  alternate(readChar('a'), readChar('b')),
  iterate(readChar('c'))
)
```

Wir setzen zur Veranschaulichung aber Zwischenfunktionen, also

$f_x := \text{readChar}(x)$  für  $x \in \{a, b, c\}$ ,  $f_{a|b} := \text{alternate}(f_a, f_b)$ ,  $f_{c^*} := \text{iterate}(f_c)$  und  $p := \text{concatenate}(f_{a|b}, f_{c^*})$ .

Wir nehmen also an das von überall Zugriff auf  $w$  ist, also  $w$  in Java z.B. ein `public static` der `Match`-Klasse ist. Der Aufruf `runParser(p, w)` läuft dann wie folgt ab:

$$\begin{aligned}
\text{runParser}(p, w) &= p(S) \\
&= f_{c*}(f_{a|b}(S)) \\
&= f_{c*}(f_a(S) \cup f_b(S)) \\
&= f_{c*}(\{(j, k+1) \mid (j, k) \in S \wedge w_k = a\} \cup f_b(S)) \\
&= f_{c*}(\{(1, 2)\} \cup f_b(S)) \\
&= f_{c*}(\{(1, 2)\} \cup \{(j, k+1) \mid (j, k) \in S \wedge w_k = b\}) \\
&= f_{c*}(\{(1, 2)\} \cup \{(2, 3)\}) \\
&= f_{c*}(\{(1, 2), (2, 3)\}) \\
&= \bigcup_{i \in \mathbb{N}_0} f_c^i(\{(1, 2), (2, 3)\}) \\
&= \{(1, 2), (2, 3)\} \cup \bigcup_{i \in \mathbb{N}} f_c^i(\{(1, 2), (2, 3)\}) \\
&= \{(1, 2), (2, 3)\} \cup \bigcup_{i \in \mathbb{N}} \{(j, k+i) \mid (j, k) \in \{(1, 2), (2, 3)\} \wedge \forall l \in [1, i] : w_{k+l} = c\} \\
&= \{(1, 2), (2, 3)\} \cup \{(2, 4)\} \cup \{(2, 5)\} \cup \emptyset \cup \dots \\
&= \{(1, 2), (2, 3), (2, 4), (2, 5)\}
\end{aligned}$$

Das Ergebnis entspricht den Matches  $\mathbf{x|a|bccx}$ ,  $\mathbf{xa|b|ccx}$ ,  $\mathbf{xa|bc|cx}$ ,  $\mathbf{xa|bcc|x}$ .

Grundsätzlich haben wir also eine Regex-Matching Funktion rekursiv über den Aufbau unserer Regexe definiert. In einem würde es sich auch wie folgt schreiben lassen:

Sei  $\Sigma$  ein Alphabet, und  $\mathcal{R}_\Sigma$  die Menge aller Regexe über  $\Sigma$ .

Sei ferner  $r \in \mathcal{R}_\Sigma$ ,  $w \in \Sigma^*$  und  $S := \{(i, i) \mid i \in [0, |w| - 1]\}$ . Dann definiere:

$$\begin{aligned}
\text{ext}(r, S) &:= \begin{cases} \{(i, j+1) \mid (i, j) \in S \wedge w_j = c\} & \text{falls } r = c \in \Sigma \\ \text{ext}(r_2, \text{ext}(r_1, S)) & \text{falls } r = (r_1 r_2) \text{ für } r_1, r_2 \in \mathcal{R}_\Sigma \\ \text{ext}(r_1, S) \cup \text{ext}(r_2, S) & \text{falls } r = (r_1 \mid r_2) \text{ für } r_1, r_2 \in \mathcal{R}_\Sigma \\ \text{iter}(t, S) & \text{falls } r = t^* \text{ für } t \in \mathcal{R}_\Sigma \end{cases} \\
\text{iter}(t, S) &:= \begin{cases} \emptyset & \text{falls } S = \emptyset \\ S \cup \text{iter}(t, \text{ext}(t, S)) & \text{sonst} \end{cases} \\
\text{runParser}(r, w) &:= \text{ext}(r, S)
\end{aligned}$$

Um Anforderungen wie z.B die Priorität in Alternationen zu erfüllen, müsste man der Klasse **Match** Informationen hinzufügen. Man kann bspw. also ein Tripel (startIndex, endIndex, priority) anstatt dem obigen Paar benutzen.

Dies erlaubt einem bei der Rekursion an jeder Stelle einer Alternation die Priorität der Matches entsprechend anzupassen.

Wenn wir 0 als höchste,  $\infty$  als geringste Priorität betrachten, definieren wir

$$\text{decrPrio}(S) := \{(i, j, p+1) \mid (i, j, p) \in S\}$$

Dann können wir die 3. Zeile der Definition von **ext** ersetzen durch:

$$\text{ext}(r, S) = \text{ext}(r_1, S) \cup \text{decrPrio}(\text{ext}(r_2, S)) \quad \text{falls } r = (r_1 \mid r_2) \text{ für } r_1, r_2 \in \mathcal{R}_\Sigma$$

Damit haben wir dann neben StartPosition und Länge ein neues Kriterium, nach welchem wir die Ergebnismenge von Matches des Parsers ordnen können. Die restlichen Funktionen entsprechend anzupassen sollte kein Problem darstellen.