



Programmierung 1 (WS 2020/21)

Zusatztutorialium 7 (Lösungsvorschläge)

Konkrete Syntax

Hinweis: Diese Aufgaben wurden von den Tutoren für das Zusatztutorialium erstellt. Sie sind für die Klausur weder relevant noch irrelevant. 🤔 markiert potentiell schwerere Aufgaben.

1 Lexer

Aufgabe Z7.1 (*Lexer*)

Schreiben Sie eine Prozedur `lex : string → token list`, die ganze Zahlen und Bezeichner einliest. In vielen Aspekten verhält sich die Sprache wie SML.

- Negative Zahlen werden mit `~` notiert.
- Bezeichner müssen mit einem Buchstaben anfangen, danach können Buchstaben oder Ziffern folgen, jedoch keine anderen Zeichen.
- Bei einem Zeichen, das nicht mehr zu einer Zahl bzw. einem Bezeichner gehört, soll ein neues Token begonnen werden.
- Als Trennzeichen wird ausschließlich das Leerzeichen betrachtet.

Gegeben ist folgender Datentyp:

```
1 datatype token = ICON of int | ID of string
```

Hinweis: Sie können `Char.isDigit` und `Char.isAlpha` verwenden. Es bietet sich an, die eigentliche Funktionalität in drei verschränkt rekursive Prozeduren `lex'`, `lexInt` und `lexId` auszugliedern.

Lösungsvorschlag Z7.1

```
1 exception Error of string
2
3 fun lex' [] = []
4   | lex' ("␣" :: cr) = lex' cr
5   | lex' ("~" :: c :: cr) = if Char.isDigit c then lexInt ~1 0 (c :: cr)
6                             else raise Error "lex"
7   | lex' (c :: cr) = if Char.isDigit c then lexInt 1 0 (c :: cr)
8                       else if Char.isAlpha c then lexId [] (c :: cr)
9                           else raise Error "lex"
10 and lexInt s a [] = [ICON (s * a)]
11   | lexInt s a (c :: cr) = if Char.isDigit c
12                             then lexInt s (a * 10 + ord c - ord #"0") cr
13                             else ICON (s * a) :: lex' (c :: cr)
14 and lexId a [] = [ID (implode (rev a))]
15   | lexId a (c :: cr) = if Char.isAlpha c orelse Char.isDigit c
16                           then lexId (c :: a) cr
17                           else ID (implode (rev a)) :: lex' (c :: cr)
18
19 fun lex s = lex' (explode s)
```

Aufgabe Z7.2 (*Strings in Strings*)

Wir möchten einen Lexer für Strings und Bezeichner schreiben, wie sie in SML vorkommen können. Zusätzlich soll eine Konkatenation für Strings existieren. Bedenken Sie, dass auch Zitate in einem String auftauchen können. Wir wollen uns dabei nach folgender Grammatik richten:



$$\begin{aligned} \text{exp} &::= \text{pexp} \, [\, \text{"^"} \, \text{exp}] \\ \text{pexp} &::= \text{id} \mid \text{" " str " " } \mid \text{"(" exp "}" \end{aligned}$$

- (a) Erstellen Sie den passenden Konstruktortypen `token` und eine Ausnahme `Invalid`
- (b) Schreiben Sie einen Lexer für Strings

Lösungsvorschlag Z7.2

(a)

```
1 datatype token = ID of string | CONCAT | LPAR | RPAR | LITERAL of string
2 exception Invalid
```

(b)

```
1 fun lex nil = nil
2   | lex (#" " :: xr) = lex xr
3   | lex (#"\t" :: xr) = lex xr
4   | lex (#"\n" :: xr) = lex xr
5   | lex (#"(" :: xr) = LPAR :: lex xr
6   | lex (#")" :: xr) = RPAR :: lex xr
7   | lex (#"^" :: xr) = CONCAT :: lex xr
8   | lex (#"\\" :: xr) = lexStr nil xr
9   | lex (c :: xr) = if Char.isAlpha c then lexId [c] xr
10                    else raise Invalid
11
12 and lexId c xr = if null xr orelse not (Char.isAlpha (hd xr))
13                 then ID(implode(rev c)) :: lex xr
14                 else lexId (hd xr :: c) (tl xr)
15
16 and lexStr c nil = raise Invalid
17   | lexStr c (#"\" :: #"\\" :: xr) = lexStr (#"\\" :: c) xr
18   | lexStr c (#"\" :: xr) = LITERAL (implode (rev c)) :: lex xr
19   | lexStr c (d::xr) = lexStr (d::c) xr
```

2 Parser & Grammatiken

Aufgabe Z7.3 (*Leist du noch oder parst du schon?*)

Wir betrachten die Sprache $\alpha\beta\gamma$. Zusätzlich zu Zahlen und expliziten Klammern seien folgende Operatoren gegeben: α , β und γ .

- α und β sind binäre Infix-Operatoren und γ ein unärer Operator.
 - α klammert rechts und β klammert links.
 - γ klammert stärker als β und α klammert am stärksten.
- (a) Erstellen Sie zunächst eine Grammatik, welche die obigen Bedingungen erfüllt.
- (b) Geben Sie einen Ausdruck nach dieser Grammatik an, in der jeder Operator einmal vorkommt.
- (c) Verändern Sie die Grammatik, sodass sie RA-tauglich wird.
- (d) Erstellen Sie einen Konstruktortyp `token` und einen Konstruktortyp `exp` für die Baumdarstellung der Ausdrücke.
- (e) Schreiben Sie nun einen Parser. Bei einer ungültigen Eingabe soll die Ausnahme `Parse` geworfen werden.

Lösungsvorschlag Z7.3

- (a)
- | | |
|---|--|
| $beta ::= gamma \mid beta \, \beta \, gamma$ | $beta ::= [beta \, \beta] \, gamma$ |
| $gamma ::= alpha \mid \gamma \, gamma$ | $gamma ::= \gamma \, gamma \mid alpha$ |
| $alpha ::= prim \mid prim \, \alpha \, alpha$ | $alpha ::= prim \, [\alpha \, alpha]$ |
| $prim ::= num \mid (\, beta \,)$ | $prim ::= num \mid (\, beta \,)$ |

- (b) $1 \, \beta \, \gamma \, 2 \, \alpha \, 3$

- (c) RA-taugliche Grammatik mit Hilfskategorien:

$$\begin{aligned} beta &::= gamma \, beta' \\ beta' &::= [\beta \, gamma \, beta'] \\ gamma &::= \gamma \, gamma \mid alpha \\ alpha &::= prim \, [\alpha \, alpha] \\ prim &::= num \mid (\, beta \,) \end{aligned}$$

- (d)
-
- ```
1 datatype token = ALPHA | BETA | GAMMA | ICON of int | LPAR | RPAR
2
3 datatype exp = Alpha of exp * exp | Beta of exp * exp | Gamma of exp | Icon of int
```
- 
- (e)
- 
- ```
1 exception Parse
2
3 fun beta ts = beta' (gamma ts)
4 and beta' (e, BETA :: tr) = let val (e', tr') = gamma tr in beta' (Beta (e, e'), tr') end
5   | beta' s = s
6
7 and gamma (GAMMA :: tr) = let val (e', tr') = gamma tr in (Gamma e', tr') end
8   | gamma s = alpha s
9
10 and alpha ts = (case prim ts of (e, ALPHA :: tr) => let val (e', tr') = alpha tr
11   | _ => s) in (Alpha (e, e'), tr') end
12
13 and prim (ICON x :: tr) = (Icon x, tr)
14   | prim (LPAR :: tr) = (case beta tr of (e', RPAR :: tr') => (e', tr')
15   | _ => raise Parse)
16   | prim _ = raise Parse
17
18 fun parse ts = case beta ts of (e, nil) => e
19   | _ => raise Parse
20
```
-

Aufgabe Z7.4 (XY)

Wir betrachten die Sprache XY , die aus den beiden binären Präfix-Operatoren X und Y und Zahlen besteht. Die zugehörige abstrakte Grammatik ist folgendermaßen definiert:

$e \in Exp = "X" \exp \exp \mid "Y" \exp \exp \exp \mid num$

- (a) Machen Sie sich klar, dass XY keine Klammern benötigt.
- (b) Stellen Sie eine konkrete Grammatik für XY auf.
- (c) Schreiben Sie einen Parser `parse: token list \rightarrow exp` für XY .

```
1 datatype token = BIGX | BIGY | ICON of int
2 datatype exp = X of exp * exp | Y of exp * exp * exp | Icon of int
```

Lösungsvorschlag Z7.4

- (a) Wenn man in der Eingabe jeden Operator durch seine Stelligkeit ersetzt, also X durch 2, Y durch 3 und jedes num -Token durch 0, dann handelt es sich um die Prälinearisierung des Syntaxbaums. Da wir aus der Prälinearisierung einen Baum eindeutig rekonstruieren können, können wir das auch bei einem Ausdruck aus der Sprache XY .
- (b) $exp ::= "X" \exp \exp \mid "Y" \exp \exp \exp \mid num$

- (c)
-
- ```
1 exception Parse
2 datatype token = BIGX | BIGY | ICON of int
3 datatype exp = X of exp * exp | Y of exp * exp * exp | Icon of int
4
5 fun exp (BIGX :: tr) = let val (e1, tr') = exp tr
6 val (e2, tr'') = exp tr' in (X (e1, e2), tr'') end
7 | exp (BIGY :: tr) = let val (e1, tr') = exp tr
8 val (e2, tr'') = exp tr'
9 val (e3, tr''') = exp tr'' in (Y (e1, e2, e3), tr''') end
10 | exp ((ICON i) :: tr) = (Icon i, tr)
11 | exp _ = raise Parse
12
13 fun parse ts = case exp ts of (e, nil) \Rightarrow e | _ \Rightarrow raise Parse
```
- 

### Aufgabe Z7.5 (JSON)

*JSON* (kurz für *JavaScript Object Notation*, <https://json.org>) ist, wie der Name schon vermuten lässt, ein kleiner Teil der Sprache *JavaScript*, mit dem sich Objekte repräsentieren lassen. Wegen der relativ einfachen Syntax wird *JSON* insbesondere bei Webanwendungen gerne als Format verwendet, um kleine strukturierte Datenmengen auszutauschen.

In dieser Aufgabe schreiben wir einen Parser für *JSON*. Die ganze *JSON*-Datei repräsentiert genau einen Wert. Ein Wert kann `null`, `true`, `false`, eine **Zahl**, ein **String**, ein **Array** oder ein **Objekt** sein. Ein Array besteht aus beliebig vielen Werten, die durch je ein Komma getrennt werden, und wird mit eckigen Klammern notiert, also zum Beispiel `[1, "abc", 1.2, true]` oder `[]`. Ein Objekt ist im Prinzip eine Liste von Schlüssel-Wert-Paaren, wobei ein Schlüssel immer ein String ist und mit einem Doppelpunkt vom Wert getrennt wird. Objekte werden mit geschweiften Klammern notiert, also zum Beispiel `{"abc": 123, "xyz": null}` oder `{}`.

Folgende Datentypen sind bereits gegeben:

---

```
1 datatype token = LBRACE (* { *) | RBRACE (* } *) | LBRACK (* [*) | RBRACK (*] *)
2 | COMMA (* , *) | COLON (* : *)
3 | STR of string | INT of int | REAL of real | TRUE | FALSE | NULL
4 datatype value = Null | True | False | Int of int | Real of real | Str of string
5 | Array of value list | Object of (string * value) list
```

---

- (a) Stellen Sie eine RA-taugliche Grammatik für *JSON* auf. Lassen Sie der Einfachheit halber zunächst das leere Array `[]` und das leere Objekt `{}` aus.

- (b) Implementieren Sie den zugehörigen Parser.
- (c) Es ist nicht schwer, den Parser auch [] und {} erkennen zu lassen. Erweitern Sie Ihren Parser entsprechend. Wenn Sie die Grammatik analog anpassen, ist diese nicht mehr RA-tauglich. Warum? Ist es möglich, dennoch RA-Tauglichkeit zu erlangen? 🤔

### Lösungsvorschlag Z7.5

(a)

$$\begin{aligned} \text{value} &::= \text{null} \mid \text{true} \mid \text{false} \mid \text{int} \mid \text{real} \mid \text{string} \mid \text{"[" list "]} \mid \text{"{" kvlist "}} \\ \text{list} &::= \text{value} [ \text{" , " list} ] \\ \text{kvlist} &::= \text{string} [ \text{" : " value [ \text{" , " kvlist} ]} \end{aligned}$$

(b)

---

```

1 exception Parse
2
3 fun value (NULL :: tr) = (Null, tr)
4 | value (TRUE :: tr) = (True, tr)
5 | value (FALSE :: tr) = (False, tr)
6 | value ((INT x) :: tr) = (Int x, tr)
7 | value ((REAL x) :: tr) = (Real x, tr)
8 | value ((STR x) :: tr) = (Str x, tr)
9 | value (LBRACK :: tr) = (case list tr of (vs, RBRACK :: tr') => (Array vs, tr')
10 | _ => raise Parse)
11 | value (LBRACE :: tr) = (case kvlist tr of (kvs, RBRACE :: tr') => (Object kvs, tr')
12 | _ => raise Parse)
13 | value _ = raise Parse
14 and list ts = (case value ts
15 of (v, COMMA :: tr') => (case list tr' of (vs, tr'') => (v :: vs, tr''))
16 | (v, tr') => ([v], tr'))
17 and kvlist ((STR k) :: COLON :: tr) = (case value tr
18 of (v, COMMA :: tr')
19 => (case kvlist tr'
20 of (kvs, tr'') => ((k, v) :: kvs, tr''))
21 | (v, tr') => [(k, v)], tr'))
22 | kvlist _ = raise Parse
23
24 fun parse ts = case value ts of (v, []) => v | _ => raise Parse

```

---

(c) Der Parser wird folgendermaßen erweitert:

---

```

1 fun value ...
2 | value (LBRACK :: RBRACK :: tr) = (Array [], tr)
3 | value (LBRACK :: tr) = (case list tr of (vs, RBRACK :: tr') => (Array vs, tr')
4 | _ => raise Parse)
5 | value (LBRACE :: RBRACE :: tr) = (Object [], tr)
6 | value (LBRACE :: tr) = (case kvlist tr of (kvs, RBRACE :: tr') => (Object kvs, tr')
7 | _ => raise Parse)
8 ...

```

---

Die zugehörige Grammatik ist folgende:

$$\begin{aligned} \text{value} &::= \text{null} \mid \text{true} \mid \text{false} \mid \text{int} \mid \text{real} \mid \text{string} \mid \text{"[" "]} \mid \text{"[" list "]} \mid \text{"{" "}} \mid \text{"{" kvlist "}} \\ \text{list} &::= \text{value} [ \text{" , " list} ] \\ \text{kvlist} &::= \text{string} [ \text{" : " value [ \text{" , " kvlist} ]} \end{aligned}$$

Diese Grammatik ist nicht RA-tauglich, da in der Kategorie *value* mindestens zwei Token benötigt werden, um die Fälle voneinander zu unterscheiden. Nach unserer Definition dürfen wir jedoch nur maximal ein Token verwenden. Ein Lösungsansatz ist nun, für "[" und "{" je eine Hilfskategorie einzuführen und dort die Fallunterscheidung vorzunehmen:

$$\begin{aligned}
value &::= null \mid true \mid false \mid int \mid real \mid string \mid "[" array' \mid "{" object' \\
array' &::= "]" \mid list "]" \\
object' &::= "}" \mid kvlst "{" \\
list &::= value [", " list] \\
kvlst &::= string ":" value [", " kvlst]
\end{aligned}$$

Doch auch diese Grammatik ist noch nicht RA-tauglich. Die Kategorien *list* und *kvlst* sind so zu behandeln, als könnten sie für etwas beliebiges stehen, also rein theoretisch auch ein "]" bzw. "}". Dementsprechend kann die Entscheidung noch immer nicht eindeutig getroffen werden.

Das Problem lässt sich lösen, indem man einmal *list* bzw. *kvlst* bis zu einem Terminal ersetzt, was allerdings für *list* sehr unschön ist:

$$\begin{aligned}
value &::= null \mid true \mid false \mid int \mid real \mid string \mid "[" array' \mid "{" object' \\
array' &::= "]" \mid null [", " list] "]" \mid true [", " list] "]" \mid false [", " list] "]" \mid int [", " list] "]" \\
&\quad \mid real [", " list] "]" \mid string [", " list] "]" \mid "[" array' [", " list] "]" \mid "{" object' [", " list] "]" \\
object' &::= "}" \mid string ":" value [", " kvlst] "}" \\
list &::= value [", " list] \\
kvlst &::= string ":" value [", " kvlst]
\end{aligned}$$

In *einzelnen* Fällen ist es also unter Verwendung gewisser Zusatzinformationen (in unserem Fall, dass eine schließende Klammer direkt nach einer öffnenden immer das leere Array bzw. das leere Objekt ist) einfacher, einen Parser zu einer nicht-RA-tauglichen Grammatik zu implementieren, als die Grammatik erst RA-tauglich zu machen.