

AIN - Sprachkonzepte

Übungsaufgaben

Bericht

WS2022/23

Inhaltsverzeichnis

Abbildungsverzeichnis

Quelltextverzeichnis

1	Aufgabe 1	1
2	Aufgabe 2	7
3	Aufgabe 3	13
4	Aufgabe 4	14

Abbildungsverzeichnis

1.1	Ausgabe Java Formatter	3
1.2	Ausgabe TimeLexer richtiger Input	6
1.3	Ausgabe TimeLexer falscher Input	6
2.1	Parse Tree Beispiele	8
2.2	UML-Diagramm	9

Quelltextverzeichnis

1.1	Funktion für Java Formatter	2
1.2	Input für Formatter	3
1.3	Lexer für Date/Time	4
1.4	alternativer Lexer	5
2.1	Lexer Grammar A2a	7
2.2	Parser Grammar A2a	8
2.3	Grundlegende AST-Klassen	10
2.4	Builder-Klasse für den AST	11

1. Aufgabe 1

a)

Schreiben Sie ein Java-Programm, das in einem String Formatspezifikationen gemäß *java.util.Formatter* findet.

Erstellen Sie dazu mit der Syntax von *java.util.regex.Pattern* einen regulären Ausdruck für eine solche Formatspezifikation.

Sie brauchen darin nicht zu berücksichtigen, dass bestimmte Angaben innerhalb einer Formatspezifikation nur bei bestimmten Konversionen erlaubt sind. Achten Sie aber bei `argument_index`, `width` und `precision` darauf, ob der Zahlbereich bei 0 oder 1 beginnt.

Beispieleingaben:

```
xxx %d yy%n
xxx%1$d yy
%1$-02.3dyyy
Wochentag: %tA Uhrzeit: %tT
```

Beispielausgaben:

```
TEXT("xxx ")FORMAT("%d")TEXT(" yy")FORMAT("%n")
TEXT("xxx")FORMAT("%1$d")TEXT(" yy")
FORMAT("%1$-02.3d")TEXT("yyy")
TEXT("Wochentag:")FORMAT("%tA")TEXT("Uhrzeit:")FORMAT("%tT")
```

a - Lösung

Code

```
1 private static String formatter(String input) {
2     Pattern patternGeneral =
3         Pattern.compile(
4
5         "(%([1-9]+\\$)?[--#0,(\\s]?\\d*(\\.\\d)?[bBhHsScCdoXxEfgGaA%n])"
6         );
7     Pattern patternDate =
8         Pattern.compile(
9
10        "(%([1-9]+\\$)?[--#0,(\\s]?\\d*[tT][HIklLMSpQZzsBbhAaCYyjmdrTrDFc])"
11        );
12    Pattern patternLeftover =
13        Pattern.compile(
14            "(%([1-9]+\\$)?[--#0,(\\s]?\\d*\\D)"
15            );
16    Pattern usePattern = Pattern.compile(
17        patternGeneral.pattern()
18        + "|" + patternDate.pattern()
19        + "|" + patternLeftover.pattern()
20    );
21
22    var builder = new StringBuilder();
23
24    Map<String, String> parts = new
25    TreeMap<>(Comparator.comparing(input::indexOf));
26
27    Arrays.stream(input.split(usePattern.toString()))
28        .forEach(x -> parts.put(x, "TEXT(\"" + x + "\")"));
29
30    usePattern.matcher(input).results()
31        .forEach(x -> parts.put(x.group(), "FORMAT(\"" + x.group() +
32        "\")"));
33
34    parts.forEach((x, y) -> builder.append(y));
35
36    return builder.toString();
37 }
```

Quelltext 1.1: Funktion für Java Formatter

Erklärung

Das Pattern ist realisiert anhand der [Java Formatter Docs](#).

Dazu unterteilen wir in 3 Unterscheidungen:

- generelle Formate

- Datum
- Rest/Unspezifiziert

Was sich alle Formate teilen, ist das beginnende '%'. Formate erlauben einen optionalen Format-Index anzugeben: '[1-9]+\\$', sowie verschiedene Formattierungs-Optionen: '[-+#0,(]?' und eine optionale Angabe für die minimale Anzahl an Charakteren: '\d*'. Zusammenhängend ergibt das: '%([1-9]+\\$)?[-+#0,(\s)?\d*'.

Bei den generellen Formaten gibt es die zusätzliche Option für die Anzahl Nachkommastellen: '(\.\d)?' gefolgt von den möglichen Konversionen: '[bBhHsScCdoXxEfgGaA%n]'. Für ein Datum fehlt noch '[tT]' und Konversionen: '[HIk1LMSpQZzsBbhAaCYyjmdeRTrDFc]'. Alle restlichen Konversionen sind illegal und reserviert für zukünftige Erweiterungen. Diese decken wir mit '\D' ab.

Anschließend wird der String in einzelne Teile zerlegt, die dann in einer Map gespeichert werden. Sortiert wird anhand der Position des Keys im Input-String. Die fertige Map wird dann in einen String umgewandelt.

Programmausgabe

```
1 public static void main(String[] args) {
2     System.out.println(formatter("xxx %d yyy%n"));
3     System.out.println(formatter("xxx%1$d yyy"));
4     System.out.println(formatter("%1$-02.3dyyy"));
5     System.out.println(formatter("Wochentag: %tA Uhrzeit: %tT"));
6 }
```

Quelltext 1.2: Input für Formatter

```
PS C:\Users\Admin\Desktop\Sprachkonzepte\src\u1> java .\Format.java
TEXT("xxx ")FORMAT("%d")TEXT(" yyy")FORMAT("%n")
TEXT("xxx")FORMAT("%1$d")TEXT(" yyy")
FORMAT("%1$-02.3d")TEXT("yyy")
TEXT("Wochentag: ")FORMAT("%tA")TEXT(" Uhrzeit: ")FORMAT("%tT")
```

Abbildung 1.1: Ausgabe Java Formatter

b)

Erkennen Sie mit ANTLR 4 Lexer-Regeln Zeitangaben im digitalen 12-Sunden-Format gemäß https://en.wikipedia.org/wiki/12-hour_clock. Beachten Sie auch die alternativen Schreibweisen 12 midnight und 12 noon. Testen Sie mit *org.antlr.v4.gui.TestRig*.

b - Lösung

Code

```
1  lexer grammar TimeLexer;  
2  
3  Time12H: Default|Noon|Midnight;  
4  
5  fragment Default: ('12:00'|(( [1-9]|'1'[01])':'[0-5][0-9]))WS[ap]'.m.';  
6  fragment Noon: 'Noon'|'12 noon';  
7  fragment Midnight: 'Midnight'|'12 midnight';  
8  
9  WS: [ \t\r\n]+ -> skip;
```

Quelltext 1.3: Lexer für Date/Time

Erklärung

Lexer Grammatiken beschreiben die Token, die vom Lexer erkannt werden sollen.

Fragmente sind Teile der Grammatik, die nicht direkt erkannt werden, sondern nur in anderen Regeln verwendet werden.

Der Ansatz hier war die Zeitangaben in drei Teile zu zerlegen:

- Default: Volle Uhrzeitangaben im klassischen 'HH:MM' Format mit AM/PM Angabe
- Noon: Zusätzlich die Mittagszeit '12 noon' und 'Noon'
- Midnight: Synchron dazu Mitternacht '12 midnight' und 'Midnight'

Noon und Midnight sind hierbei die Ausnahme, aber vorgegeben durch die Aufgabenstellung.

Alternativ könnte man mehr Token beschreiben:

```
1  lexer grammar TimeLexerV2;
2
3  TIME : HOUR SEPERATOR MINUTE (AM | PM)
4  | TWELVE SEPERATOR '00' (AM | PM)
5  | TWELVE 'noon'
6  | TWELVE 'midnight'
7  | 'Noon'
8  | 'Midnight';
9
10 TWELVE : '12';
11
12 HOUR : '1'[0-1]|[0-9];
13 MINUTE : [0-5][0-9];
14 SEPERATOR : ':';
15
16 AM : 'a.m.' ;
17 PM : 'p.m.' ;
18
19 WS: [ \t\r\n]+ -> skip;
```

Quelltext 1.4: alternativer Lexer

Allerdings wird die Lexer Grammatik hier etwas missbraucht, da die 'TIME' Regel eher als Parser Regel genutzt wird. Lexer Regeln sollten eigentlich nur die Token beschreiben, die vom Lexer erkannt werden sollen.

Programmausgabe

Korrekter Input:

12:00 AM

12:00 a.m.

12:00 am

3:00 am

1:12 am

Noon

12 noon

Midnight

11:59 PM

```
C:\PROG_uebungen\Sprachkonzepte\src\ui>java -cp ".;C:\Program Files\Java\lib\antlr-4.11.1-complete.jar;" org.antlr.v4.gui.TestRig TimeLexer tokens -tokens RichtigeBeispiele.txt
[0,0:7='12:00 AM',<Time12H>,1:0]
[01,10:19='12:00 a.m.',<Time12H>,2:0]
[02,22:29='12:00 am',<Time12H>,3:0]
[03,32:38='3:00 am',<Time12H>,4:0]
[04,41:47='1:12 am',<Time12H>,5:0]
[05,52:55='Noon',<Time12H>,7:0]
[06,58:64='12 noon',<Time12H>,8:0]
[07,67:74='Midnight',<Time12H>,9:0]
[08,77:84='11:59 PM',<Time12H>,10:0]
[09,85:84='<EOF>',<EOF>,10:8]
```

Abbildung 1.2: Ausgabe TimeLexer richtiger Input

Falscher Input:

```
0:00 am
0:00 pm
12:00
noon
12 Noon
13:00 PM
```

```
C:\PROG_uebungen\Sprachkonzepte\src\ui>java -cp ".;C:\Program Files\Java\lib\antlr-4.11.1-complete.jar;" org.antlr.v4.gui.TestRig TimeLexer tokens -tokens FalscheBeispiele.txt
line 1:0 token recognition error at: '0'
line 1:1 token recognition error at: ':'
line 1:2 token recognition error at: '0'
line 1:3 token recognition error at: '0'
line 1:5 token recognition error at: 'a'
line 1:6 token recognition error at: 'm'
line 2:0 token recognition error at: '0'
line 2:1 token recognition error at: ':'
line 2:2 token recognition error at: '0'
line 2:3 token recognition error at: '0'
line 2:5 token recognition error at: 'p'
line 2:6 token recognition error at: 'm'
line 3:0 token recognition error at: '12:00\r\nn'
line 4:1 token recognition error at: 'o'
line 4:2 token recognition error at: 'o'
line 4:3 token recognition error at: 'n'
line 5:0 token recognition error at: '12 N'
line 5:4 token recognition error at: 'o'
line 5:5 token recognition error at: 'o'
line 5:6 token recognition error at: 'n'
line 6:0 token recognition error at: '13'
line 6:2 token recognition error at: ':'
line 6:3 token recognition error at: '0'
line 6:4 token recognition error at: '0'
line 6:6 token recognition error at: 'P'
line 6:7 token recognition error at: 'M'
[0,48:47='<EOF>',<EOF>,6:8]
```

Abbildung 1.3: Ausgabe TimeLexer falscher Input

2. Aufgabe 2

a)

Denken Sie sich eine kleine Sprache aus. Definieren Sie deren Vokabular mit einer ANTLR4 lexer grammar und deren Grammatik mit einer ANTLR4 parser grammar. Erzeugen Sie für einige Beispieltexte mit Hilfe von *org.antlr.v4.gui.TestRig* den Ableitungsbaum (Parse Tree).

a - Lösung

Dargestellt ist der Lexer einer Sprache, welche die korrekte Kreation einer Java Klasse darstellt. Erlaubt sind Variablen - also einzelne Buchstaben, sowie Zahlen. Parameter werden durch Komma getrennt und dürfen auch eigene, neu erzeugte Klassen sein.

```
1 // CreationLexer.g4
2 lexer grammar CreationLexer;
3
4 KEYWORD : 'new' ;
5
6 NAME : [A-Za-z]+ ;
7 NUM : [0-9]+ ;
8
9 COMMA : ',' ;
10
11 PAR_OPEN : '(' ;
12 PAR_CLOSE : ')' ;
13
14 WS : [ \t\r\n]+ ;
15
16 InvalidChar: . ;
```

Quelltext 2.1: Lexer Grammar A2a

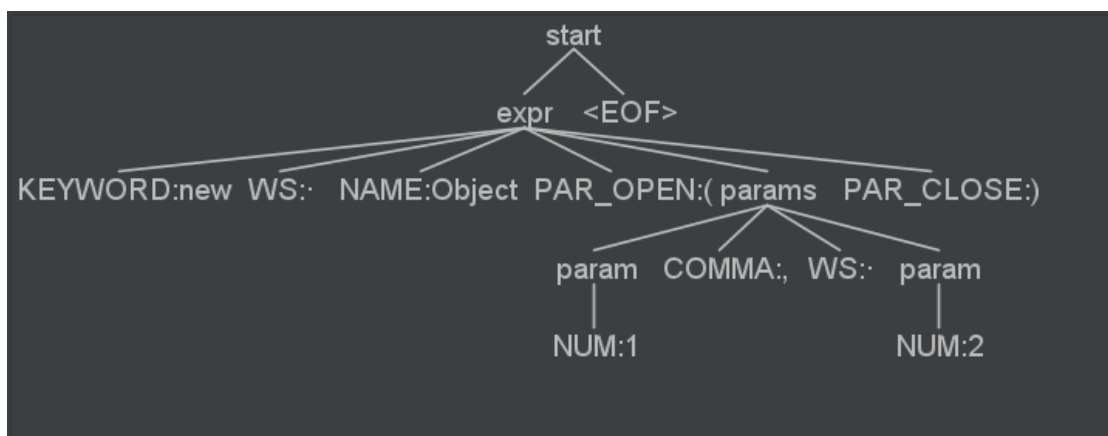
Der dazugehörige Parser:

```

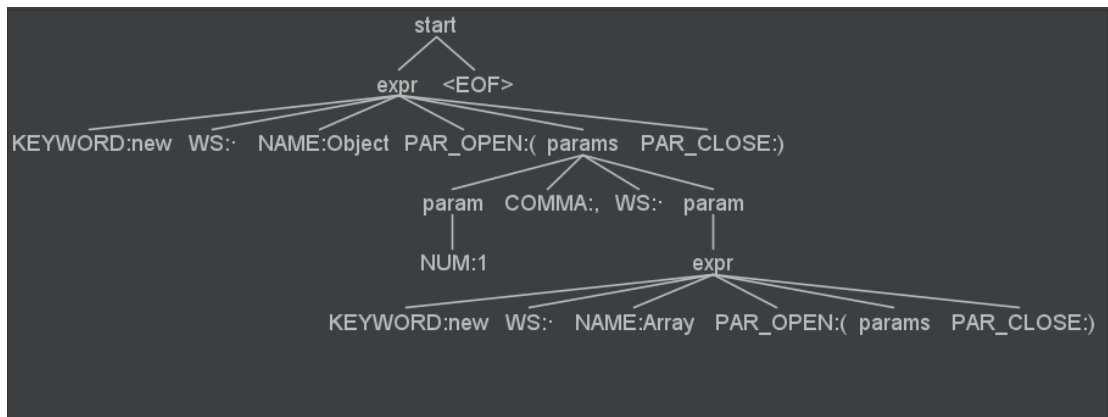
1  // CreationParser.g4
2  parser grammar CreationParser;
3  options { tokenVocab=CreationLexer; }
4
5  start : expr EOF;
6
7  expr : KEYWORD WS NAME PAR_OPEN params PAR_CLOSE ;
8
9  params : (param (COMMA WS? param)*)? ;
10
11 param : (expr | NAME | NUM) ;

```

Quelltext 2.2: Parser Grammar A2a



(a) Input: 'java org.antlr.v4.gui.TestRig Creation start -gui new Object(1, 2)'



(b) Input: 'java org.antlr.v4.gui.TestRig Creation start -gui new Object(1, new Array())'

Abbildung 2.1: Parse Tree Beispiele

Der Parser ist für die grammatikalische Anordnung der durch den Lexer vorgegebenen Token verantwortlich.

b)

Definieren Sie mit Java-Klassen die abstrakte Syntax Ihrer Sprache aus a) und schreiben Sie ein Java-Programm, das den ANTLR4 Parse Tree in einen AST überführt.

b - Lösung

Definition als UML-Diagramm:

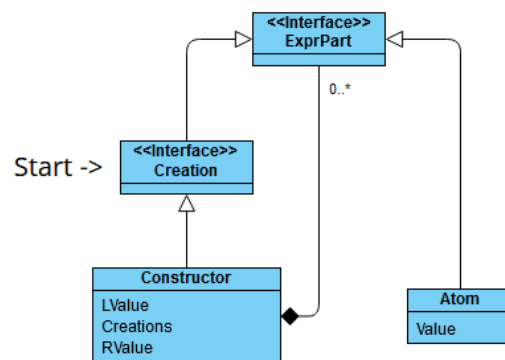


Abbildung 2.2: UML-Diagramm

Eine Baumstufe kann in drei Teile geteilt werden:

Links Der Keyword-, Namen- und Klammer-Teil vor der Mitte.

Mitte Rekursive Parameter-Definition - kann auch leer oder eine neue *Creation* sein.

Rechts In diesem Fall die Klammer zum Schließen der Parameter.

Dadurch ergibt sich die Vorgabe von mindestens einer *Constructor* Instanz. *L-* und *RValue* sind definierbare Strings und *Creations* ist die Liste der Parameter innerhalb der Klammern/Strings, dargestellt in der Relation. So können neue *Constructor* Objekte, sowie Werte definiert in der Klasse *Atom* als Parameter übergeben werden.

Um die Struktur beizubehalten, werden beispielhaft also folgende Klassen erstellt:

```
1 public interface Expr {
2 }
3
4 public interface Creation extends Expr {
5 }
6
7 public class Atom implements Expr {
8     private final String val;
9
10    public Atom(String val) {
11        this.val = val;
12    }
13    public String getVal() {
14        return val;
15    }
16    @Override
17    public String toString() {
18        return this.val;
19    }
20 }
21
22 public class Constructor implements Creation {
23     private final String leftVal;
24     private final List<Expr> params;
25     private final String rightVal;
26
27    public Constructor(String leftVal, List<Expr> params, String
rightVal) {
28        this.leftVal = leftVal;
29        this.params = params;
30        this.rightVal = rightVal;
31    }
32    public String getLeftVal() {
33        return leftVal;
34    }
35    public List<Expr> getParams() {
36        return params;
37    }
38    public String getRightVal() {
39        return rightVal;
40    }
41    @Override
42    public String toString() {
43        return this.leftVal + this.params + this.rightVal;
44    }
45 }
```

Quelltext 2.3: Grundlegende AST-Klassen

Um diese Struktur in einem Builder umzusetzen, muss ein eigener Stack für jedes neue Constructor Objekt angelegt werden.


```
51     }  
52   }  
53  
54 }
```

Quelltext 2.4: Builder-Klasse für den AST

Für jedes angefangene *Constructor* Objekt wird ein neuer Stack angelegt in der *enterExpr* Methode. Zusätzlich wird die Tiefe erhöht, um jeweils beim aktuell noch nicht abgeschlossenen *Constructor* Objekt zu bleiben. Dadurch bleibt die Reihenfolge erhalten und innerhalb eines noch nicht abgeschlossenen *Constructor* Objekts können *Atom* Objekte und weitere *Constructor* Objekte erstellt werden, ohne dass Konstruktoren die auf dem Stack liegenden Objekte in sich aufnehmen.

Die *exitExpr* Methode arbeitet den aktuellen Stack ab und setzt ihn als Parameter-Liste für das aktuelle *Constructor* Objekt. Danach wird der aktuelle Stack geleert und die Tiefe verringert. Abschließend wird der aktuelle Konstruktor auf den Stack gelegt für den nächsten Konstruktor oder den finalen Rückgabewert.

Die *enterParam* Methode ist der unterste Baustein des Baums, und dient zur Erstellung von *Atom* Objekten, die dann auf den Stack gelegt werden für Konstruktoren. Wenn die *Param* Regel wieder eine *Expr* Regel enthält, wird sie in der Methode ignoriert, da diese bereits in der *exitExpr* Methode abgearbeitet wird.

3. Aufgabe 3

4. Aufgabe 4