

SetlCup Tutorial

Jonas Eilers

April 25, 2016

Contents

1	Überblick	3
1.1	Motivation	3
2	Funktionalität	3
2.1	Benutzung	3
3	Aufrufsmethoden	3
3.1	Aufruf über Kommandozeile	4
3.2	Aufruf in SetlX	4
4	Dateistruktur	5
5	Aufbau der Definitionen	5
5.1	Kommentare	6
5.2	Scannerdefinition	6
5.3	Parserdefinition	7
6	Beispiele	9
6.1	Arithmetische Ausdrücke	9
6.2	Programmiersprachen Parser	10
6.3	Nutzung von Definitionen im Kommentarbereich	13

1 Überblick

1.1 Motivation

Zur Anwendung von formalen Sprachen in der Informatik gehören unter anderem Scanner und Parser[ASU86]. Programmgeneratoren, welche z.B. Parser generieren existieren bereits für Sprachen wie JAVA[SEHT16], C[DS] usw. In der Vorlesung "Formal Languages and their Applications"[Str16c] werden momentan verschiedene Tools genutzt um Parser zu generieren. Dazu gehören u.a. JAVACUP zum Erstellen von LALR-Parsern, sowie ANTLR4[Par12] für die Erstellung eines Top-Down-Parsers mit Nutzung von EBNF-Grammatiken. Diese Programme ermöglichen es beim Parsen eines Programmes Aktionen durchzuführen. Diese Aktionen werden in der Programmiersprache JAVA definiert. Ziel dieser Studienarbeit ist es eine Alternative zur Nutzung von JAVA(in der Parser-Generierung) zu ermöglichen. Dies soll geschehen, indem ein Parser-Generator in der Programmiersprache SETLX geschrieben wird. SETLX ist eine mengenbasierte Programmiersprache. Sie kann mit dem SETLX-Interpreter genutzt werden. Dieser ist unter der folgenden Adresse verfügbar:

<http://randoom.org/Software/SetlX>

Diese Programmiersprache wird in den Vorlesungen der theoretischen Informatik 1 und 2(Logik[Str16b], Algorithms[Str16a]) bereits ausgiebig genutzt. Eine Ergänzung für die Vorlesung über die formalen Sprachen kann somit den Studenten helfen die Konzepte des Parsens zu verstehen, ohne in eine andere Programmiersprache "umzudenken". Da der Parser-Generator sehr stark an JAVACUP orientiert ist, entstand der Name SETLCUP. Ein großer Unterschied ist jedoch, dass SETLCUP ein **LR**-Parser-Generator ist. Ein **LALR**-Parser-Generator kann weniger Grammatiken verarbeiten als ein **LR**-Parser. Zwar hat der **LALR**-Parser den Vorteil, dass die Parse-Tabellen kleiner sind als bei **LR**-Parsern, jedoch fällt dies bei den Heute verfügbaren Rechnern aufgrund der Größe moderner Hauptspeicher kaum noch ins Gewicht.

2 Funktionalität

2.1 Benutzung

Der Benutzer erstellt eine Grammatik, die mit Aktionen attribuiert ist. Daraus erstellt SETLCUP dann eine SETLX-Datei, die Parser und Scanner erhält. Diese beinhaltet einen kanonischen LR-Parser, welcher auf den gegebenen Definitionen beruht. Anschließend kann eine Eingabedatei von dem generierten Parser überprüft werden. Dabei wird ggf. angegebener Code bei der Reduzierung der entsprechenden Regeln durchgeführt. Zunächst wird erklärt, wie die Komponente aufgerufen werden kann.

3 Aufrufsmethoden

Es gibt mehrere Möglichkeiten SETLCUP aufzurufen:

3.1 Aufruf über Kommandozeile

1. `setlx setlcup.stlx -p parser_scanner_datei`
`setlx setlcup.stlx -p examples/math_simple_expression.g`

Mit diesem Aufruf wird ein Parser gemäß der in der Eingabedatei gegebenen Definitionen erstellt.

2. `setlx setlcup.stlx -p parser_scanner_datei -d`
`setlx setlcup.stlx -p examples/math_simple_expression.g -d`

Um die Art und Weise, wie der Parser generiert wird nachvollziehen zu können, kann mit der Option "-d" das Debugging eingeschaltet werden. Dabei wird empfohlen die Ausgabe in eine Datei umzuleiten.

3. `setlx parser_datei -p eingabe_datei`
`setlx math_simple_expressionGrammar.stlx -p simple_statement.txt`

Der erstellte Parser kann mit dem o.g. Befehl aufgerufen werden und versucht die Eingabedatei nach den angegebenen Regeln zu überprüfen.

4. `setlx parser_datei -p eingabe_datei -d`
`setlx math_simple_expressionGrammar.stlx -p simple_statement.txt -d`

Analog zur Parsererstellung wird durch die Option "-d" das Debugging eingeschaltet.

5. `setlx setlcup.stlx -p -help`

Dieser Aufruf zeigt die Hilfe an, wie SETLCUP aufgerufen werden kann.

6. `setlx test_setlcup.stlx`

Dieser Aufruf testet die Funktionen des Parsergenerators mit verschiedenen Beispielen. Dabei sollte das folgende Ergebnis erzeugt werden:

```
Testing math_expression_grammar..  
parser generated succesfully for math_expression_grammar_ast.g  
grammar: math_expression_grammar_ast.g input: math_expression_input.txt  
successfull: yes  
Testing interpreter_grammar..  
parser generated succesfully for interpreter_grammar_ast.g  
grammar: interpreter_grammar_ast.g input: factorial.sl successfull: yes  
grammar: interpreter_grammar_ast.g input: solve.sl successfull: yes  
grammar: interpreter_grammar_ast.g input: sum.sl successfull: yes  
grammar: interpreter_grammar_ast.g input: sum-for.sl successfull: yes
```

3.2 Aufruf in SetlX

Ein Aufruf von SETLCUP ist auch in SETLX selbst möglich. Dazu ist es notwendig die Funktionalität im eigenen Quellcode zu laden:

```
load("setlcup_load.stlx");  
generate_parser('examples\math_expression_grammar_ast.g', true);  
load('examples\math_expression_grammar_astGrammar.stlx');  
result := test_parser_from_file('examples\math_expression_input.txt', true);
```

4 Dateistruktur

Der Parser-Generator SETLCUP ist in mehrere Dateien aufgeteilt.

setlcup.stlx Diese Datei wird aufgerufen, falls man über die Kommandozeile mit SETLCUP arbeiten möchte. Sie nimmt die Definition der Grammatik als Parameter an und gibt sie an die jeweiligen Programme weiter, damit die Generierung des Parsers gestartet werden kann.

setlcup_load.stlx Diese Datei lässt zunächst den Scanner und daraufhin den Parser generieren. Falls über SETLX mit SETLCUP gearbeitet werden soll, ist es möglich diese Datei als Grundlage zu nutzen. Sie wird u.a. für den Test der Funktionalität benutzt (siehe "test_setlcup").

scanner_generator.stlx In dieser Datei wird die Eingabe analysiert. Es wird der Scanner aus den angegebenen Tokens generiert. Dieser wird daraufhin in einer Ausgabedatei abgelegt.

Auch der Parser aus der Datei "_sr_parser_part.stlx" wird der Datei angehängen. Zusätzlich wird die Grammatikdefinition analysiert und in einer Token-Liste zurückgegeben.

parser_generator.stlx Der Parser Generator erstellt den kanonischen LR-Parser. Er erzeugt die notwendigen Action-, State- und Gototabellen. Diese werden zusammen mit den Regeln und dem vom User gewünschten Code auch in der o.g. Datei abgespeichert.

_sr_parser_part.stlx Der Shift-Reduce-Parser ist in dieser Datei umgesetzt. Er wird zusammen mit den benötigten Tabellen, sowie dem Scanner in eine neue Datei kopiert.

_Grammar.stlx In dieser Datei wird der erstellte Scanner und Parser abgelegt. Beim Aufruf wird die übergebene Datei geparsed. Dabei werden die o.g. Tabellen, sowie der vorher erstellte Scanner genutzt. Der Parser ist ein Shift-Reduce Parser. Beim Reducen wird der vom User angegebene Code ausgeführt. Der Wert der Variable "result" wird am Ende zurückgegeben.

test_setlcup.stlx Für einen Regressionstest kann mit dieser Datei überprüft werden ob alle Beispiele momentan geparsed werden können. Dabei wird für mehrere Grammatiken der Parser erstellt. Diese werden anschließend mit einer Eingabedatei aufgerufen und die Ergebnisse mit den vorher festgelegten verglichen.

5 Aufbau der Definitionen

Der Aufbau der Scanner- und Parserdefinitionen ist in drei, durch "%%%" getrennte, Abschnitte zu unterteilen:

1. Kommentare

2. Scannerdefinition
3. Parserdefinition

5.1 Kommentare

Die Kommentar-sektion wird 1:1 in den späteren Parser reinkopiert. Das heißt, dass Kommentare über die Grammatik mit `"/"` auskommentiert werden müssen. Jedoch bietet dies zusätzlich die Möglichkeit Funktionsdefinitionen abzuspeichern, welche im Action-Code genutzt werden können. Ein Beispiel folgt im Abschnitt 6.3. Dieser Abschnitt endet mit dem Symbol `%%`.

5.2 Scannerdefinition

Die Aufgabe des Scanners ist es, die Eingabe in eine Liste von Tokens zu zerlegen. Die Syntax (siehe Figur 1, Seite 6) wird im Folgenden erklärt.

1	INTEGER	:= 0 [1-9][0-9]*	;
2	ASTERIKS	:= *	;
3	WHITESPACE	:= [\t\v\r\s]	;
4	SKIP	:= {WHITESPACE} \n	;

Figure 1: Scanner Definition

1. In Zeile 1 wird der Token "INTEGER" definiert. Tokens werden auf die folgende Weise deklariert:
`token_name := regex ;`
 Dabei ist folgendes bei der Syntax des regulären Ausdrucks zu beachten:
 - (a) Die Rückgabe von Capture-Gruppen¹ z.B. `"ab(.*?)ab(.*?)ab"` wird konkateniert zurückgegeben (es werden alle einzelnen Unterergebnisse aneinander gehangen).
 - (b) Die Nutzung der geschweiften Klammer wurde überlagert, sodass die regulären Ausdrücke der großgeschriebene Wörter innerhalb der Klammern im Nachhinein ersetzt werden. Siehe Zeile 4 - `"{WHITESPACE}"`.
 - (c) Ansonsten sind die bereits in SETLX bzw. JAVA vorhandenen Regex-Ausdrücke benutzbar.
2. Wie in SETLX (siehe Zeile 2 und 3) müssen auch in SETLCUP vordefinierte Symbole wie `"*"`, `"+"`, `"?"`, `"|"`, `"{"`, `"}"`, `"("`, `")"`, `"."`, `"["`, `"]"`, `"\"` escaped werden, wenn diese wörtlich gematcht werden sollen.
3. In Zeile 4 wird das "SKIP"-Token genutzt. In manchen Fällen werden gewisse Tokens nicht benötigt. Diese können mithilfe des "SKIP"-Tokens ignoriert werden. Dabei ist die Eingabe mit der o.g. Ersetzstrategie

¹Gruppierungen in regulären Ausdrücken - https://de.wikipedia.org/wiki/Regulärer_Ausdruck#Gruppierungen_und_R.C3.BCckw.C3.A4rtsreferenzen

("{TOKENNAME}") oder die Nutzung eines regulären Ausdrucks möglich. Verschiedene Tokens müssen mit der Pipe "|" separiert werden.

5.3 Parserdefinition

Die Definition der Grammatik für den Parser benutzt Konzepte aus JAVACUP und ANTLR (siehe Figur 2, Seite 7).

```

1 grammar      ::= definition_list;
2 definition_list ::= rule_definition definition_list
3               |
4               ;
5 rule_definition ::= VARIABLE '::=' body_list ';' ;
6 body_list      ::= body '|' neBody_list
7               | body
8               |
9               ;
10 neBody_List   ::= body body_list;
11 body          ::= element_list action_code;
12 action_code   ::= '{:' CODE ':}'
13               |
14               ;
15 element_list  ::= element element_list
16               |
17               ;
18 element       ::= token;
19 token         ::= LITERAL id
20               | TOKEN_NAME id
21               | VARIABLE id ;
22 id            ::= ':' VARIABLE
23               |
24               ;
25 VARIABLE      ::= [a-z] [a-zA-Z_0-9]*;
26 TOKEN_NAME    ::= [A-Z] [A-Z_0-9]*;
27 LITERAL       ::= '[' '^' ']' '*' ' ';
28 CODE          ::= [.\n]*;

```

Figure 2: Grammatik zum Aufbau der Parserdefinition

VARIABLE beschreibt den Namen einer Regel z.B. "body". Somit können Regeln referenziert werden. Außerdem haben auch die IDs einer Variable den gleichen Aufbau. Dies ermöglicht es den Rückgabewert einer Regel bzw. eines regulären Ausdrucks zu nutzen.

TOKEN_NAME bezieht sich auf einen Token aus der Scanner Definition.

LITERAL beschreibt wortwörtliche Ausdrücke, welche genutzt werden können (z.B. '+', ',', '*', 'print').

CODE ist ein optionaler Teil einer Regel. Er kann am Ende einer Regel hinzugefügt werden. Der CODE innerhalb der Klammern wird beim Reduzieren der Regel ausgeführt. Durch die Nutzung der Variable "result" können Ergebnisse zwischen den Regeln transferiert werden. Die IDs der Elemente der jeweiligen Regel können im Code benutzt werden. Der Code selber darf keine Anführungszeichen enthalten. Eine Alternative bietet der Literalstring z.B. 'a+b = 15'. Außerdem sollte **generell** auf die Nutzung von Dollarzeichen verzichtet werden. Diese führen zu Komplikationen bei der Serialisierung des Parsers, da sie in SETLX eine Sonderfunktion einnehmen.

6 Beispiele

Das erste Beispiel beschreibt eine Grammatik für einfache arithmetische Ausdrücke. Das zweite Beispiel beschreibt eine simple Programmiersprache.

6.1 Arithmetische Ausdrücke

Der arithmetische Parser (siehe Figur 3, Seite 9) kann mit der oben beschriebenen Syntax definiert werden. Die Generierung des Parsers kann nun über

```
1  %%%
2  INTEGER      := 0|[1-9][0-9]* ;
3  WHITESPACE   := [\t\v\r\s] ;
4  SKIP         := {WHITESPACE} | \n ;
5  %%%
6  term
7  ::= expr:e    {: print(e);      :}
8  ;
9  expr
10 ::= expr:e '+' prod:p  {: result := e+p;    :}
11   | prod:p      {: result := p;          :}
12 ;
13 prod
14 ::= prod:p '*' fact:f  {: result := p*f;    :}
15   | fact:f       {: result := f;          :}
16 ;
17 fact
18 ::= INTEGER:n      {: result := eval(n); :}
19 ;
```

Figure 3: Parserdefinition für einen simplen arithmetischen Ausdruck

```
setlx setlcup.stlx -p examples\math_simple_expression.g
```

gestartet werden. Mit einem Treiber Programm(siehe Figur 4, Seite 10) können eingegebene Ausdrücke geparsed werden. Der Treiber kann mit dem folgenden Kommando aufgerufen werden:

```
setlx math_simple_expression_driver.stlx
```

Beispielausgabe:

```
Give an arithmetic Expression: (q for end)2 + 3 * 4 * 5 * 6 * 7 + 8
2530
Give an arithmetic Expression: (q for end)7 * 8 + 3 * 4 + 5 + 6 + 9 * 2
97
```

```

1 load("math_simple_expressionGrammar.stlx");
2 query := "Give an arithmetic Expression: (q for end)";
3 while(true){
4   arith_expr := read(query);
5   if(!("q" in arith_expr)){
6     test_parser_from_string(arith_expr, true);
7   }
8   else{
9     break;
10  }
11 }

```

Figure 4: Beispiel für Treiber

6.2 Programmiersprachen Parser

Eine Erweiterung der arithmetischen Ausdrücken um Ausgabe-, sowie Zuweisungsfunktionen spiegelt diese simple Programmiersprache wieder. Mit dieser simplen Programmiersprache ist es möglich einer Variablen einen Term zuzuweisen ("Assignemnt"). Auch eine Variable selbst kann Teil des Terms sein. Außerdem können Terme mit dem "print"-Befehl ausgegeben werden. Der Parser gibt für ein Programm, welches aus den o.g. Elementen besteht, den AST² zurück. Der Scanner (siehe Figur 5, Seite 10) besteht aus den Tokens für Dezimalzahlen, Variablen, Kommentaren und Leerzeichen o.ä. Die Anweisungen und Defini-

```

1 WHITESPACE := [ \t\v\r\s] ;
2 INTEGER    := 0|[1-9][0-9]* ;
3 ZID        := [a-zA-Z_][a-zA-Z0-9_]* ;
4 SKIP       := {WHITESPACE}|\n|//[^\n]* ;

```

Figure 5: Scannerdefinition für eine sehr einfache Programmiersprache

tionen (siehe Figur 6, Seite 11) beschreiben den Aufbau der Eingabedatei. Sie besteht aus Zuweisungen oder Ausgaben. Diese können wieder aus arithmetischen Ausdrücken bestehen. Ein Beispiel Programm (siehe Figur 7, Seite 11) wird durch den Syntaxbaum (siehe Figur 8, Seite 12) abgebildet. Die Ausgabe des Programms könnte dabei für die Erstellung eines Interpreters genutzt werden.

²AST - abstract syntax tree, siehe https://de.wikipedia.org/wiki/Abstrakter_Syntaxbaum

```

1 program
2   ::= stmtList:d                {: result := Program(d); :}
3   ;
4 stmtList
5   ::= statement:s stmtList:s1 {: result := [s] + s1 ; :}
6   |
7   ;
8 statement
9   ::= ZID:id '=' expr:e ';'     {: result := Assignment(id, e); :}
10  | 'print' '(' expr:e ')' ';'  {: result := PrintExpr(e); :}
11  ;
12 expr
13  ::= expr:e '+' prod:p         {: result := Sum(e,p); :}
14  | expr:e '-' prod:p          {: result := Difference(e,p); :}
15  | prod:p                     {: result := p; :}
16  ;
17 prod
18  ::= prod:p '*' fact:f         {: result := Product(p,f); :}
19  | prod:p '/' fact:f          {: result := Quotient(p,f); :}
20  | fact:f                     {: result := f; :}
21  ;
22 fact
23  ::= '(' expr:e ')'           {: result := e; :}
24  | INTEGER:n                 {: result := Integer(eval(n)); :}
25  | ZID:id_1                  {: result := Variable(id_1); :}
26  ;

```

Figure 6: Grammatik für eine sehr einfache Programmiersprache

```

1 a := 5;
2 b := 7;
3 c := b*a +a/b*(7+6);
4 print(c/2);

```

Figure 7: Beispiel-Eingabe

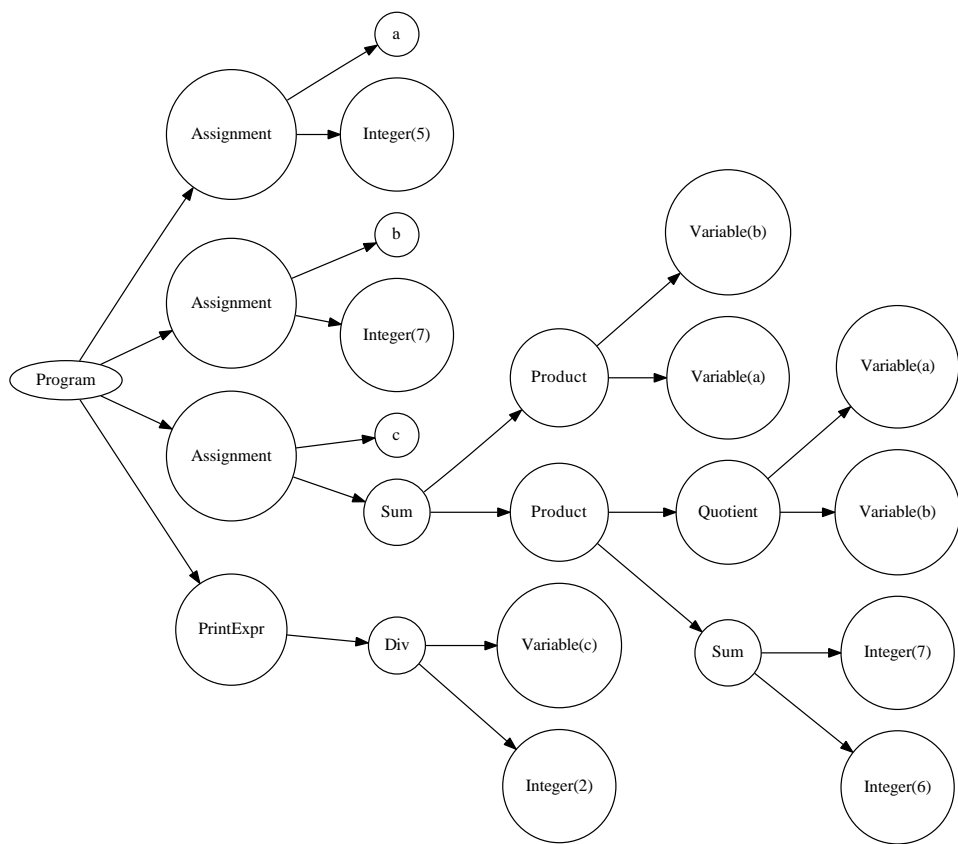


Figure 8: Interpreter AST

6.3 Nutzung von Definitionen im Kommentarbereich

Eine Ergänzung ist durch die Nutzung der Kommentarsektion(siehe Figur 9, Seite 13) möglich. Zu finden unter:

```
examples/interpreter_with_comment_action.g
```

Dabei muss die Grammatik(siehe Figur 6, Seite 11) ein wenig verändert werden(siehe Figur 10, Seite 14). Wenn für diese angepasste Definition ein Parser erzeugt wird, kann dieser die Beispieleingabe(siehe Figur 7, Seite 11) interpretieren.

```
setlx setlcup.stlx -p examples/interpreter_with_comment_action.g
cd examples
setlx interpreter_with_comment_actionGrammar.stlx -p simple_interpreter_statements.txt
```

Dabei wird der korrekte Rückgabewert geliefert:

155/7

```
1 parseElement := closure(element, rw values){
2   match(element){
3     case Program(d):
4       parseElement(d, values);
5     case [h|t]:
6       parseElement(h, values);
7       parseElement(t, values);
8     case Assignment(id, e_1):
9       values["$id$"] := parseElement(e_1, values);
10    case PrintExpr(e_2):
11      print(parseElement(e_2, values));
12    case Sum(e,p):
13      return parseElement(e, values) + parseElement(p, values);
14    case Difference(e,p):
15      return parseElement(e, values) - parseElement(p, values);
16    case Product(p,f):
17      return parseElement(p, values) * parseElement(f, values);
18    case Quotient(p,f):
19      return parseElement(p, values) / parseElement(f, values);
20    case Integer(n):
21      return n;
22    case Variable(id_1):
23      return values[id_1];
24  }
25 };
```

Figure 9: Beispiel-Kommentar

```

1 file
2   ::= program:p           {: parseElement(p, {}); :}
3   ;
4 program
5   ::= stmtList:d          {: result := Program(d); :}
6   ;

```

Figure 10: Abgeänderte Grammatikzeilen für Interpreter

References

- [ASU86] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers, Principles, Techniques, and Tools*. Addison-Wesley, 1986.
- [DS] Charles Donnelly and Richard Stallman. *Bison: the Yacc-compatible parser generator*.
- [Par12] Terence Parr. *The Definitive ANTLR 4 Reference*. Pragmatic Bookshelf, 2012.
- [SEHT16] Georgia Tech Scott E. Hudson, GVU Center and TUM. CUP - LALR parser generator for Java, 2016. Available at <http://www2.cs.tum.edu/projects/cup/>.
- [Str16a] Karl Stroetmann. Algorithms, 2016. Available at <https://github.com/karlstroetmann/Algorithms/>.
- [Str16b] Karl Stroetmann. Computer science I: Logic and set theory, 2016. Available at <https://github.com/karlstroetmann/Logik/>.
- [Str16c] Karl Stroetmann. Formal languages and their applications, 2016. Available at <https://github.com/karlstroetmann/Formal-Languages/>.