## Compilerbau

Martin Plümicke

WS 2021/22

## Agenda

Überblick Vorlesung Literatur

Compiler Überblick

Compiler (I)

Scanner

Parser

#### Literatur

Bauer and Höllerer.

Übersetzung objektorientierter Programmiersprachen. Springer-Verlag, 1998, (in german).

Alfred V. Aho, Ravi Lam, Monica S.and Sethi, and Jeffrey D. Ullman. Compiler: Prinzipien, Techniken und Werkzeuge.

Pearson Studium Informatik. Pearson Education Deutschland, 2. edition, 2008.

(in german).

Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. Compilers Principles, Techniques and Tools. Addison Wesley, 1986.

Reinhard Wilhelm and Dieter Maurer.

Übersetzerbau.

Springer-Verlag, 2. edition, 1992. (in german).

#### Literatur II

James Gosling, Bill Joy, Guy Steele, Gilad Bracha, and Alex Buckley. The Java® Language Specification.

The Java series. Addison-Wesley, Java SE 8 edition, 2014.

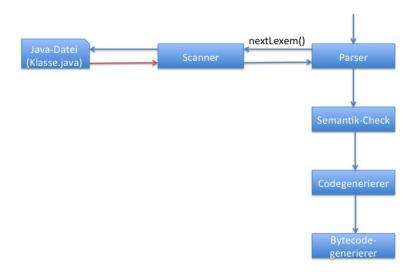
Tim Lindholm, Frank Yellin, Gilad Bracha, and Alex Buckley. The Java® Virtual Machine Specification. The Java series. Addison-Wesley, Java SE 8 edition, 2014.

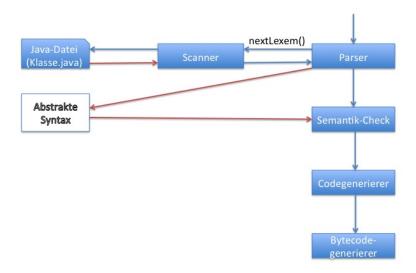
Bryan O'Sullivan, Donald Bruce Stewart, and John Goerzen. Real World Haskell.

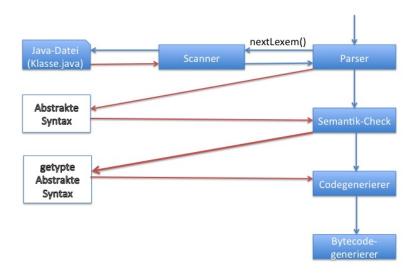
O'Reilly, 2009.

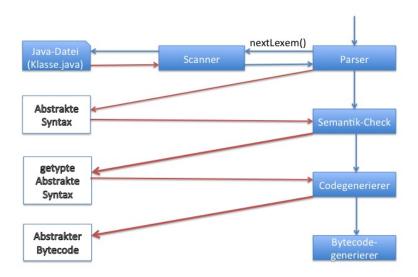
Peter Thiemann. Grundlagen der funktionalen Programmierung. Teubner, 1994.

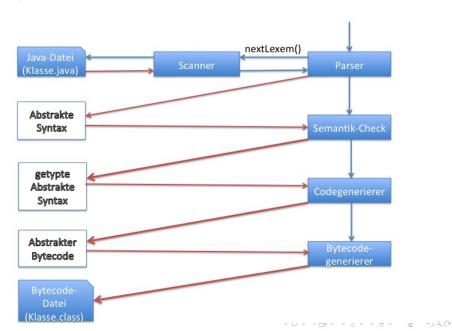




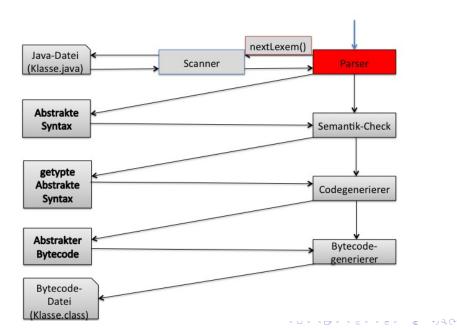








#### Parser



## Programmiersprachen

Programmiersprachen werden als formale Sprachen über einem Alphabet von Tokens definiert.

## Spezifikation eines Parser

Eingabe: Grammatik  $G = (N, \Sigma, \Pi, S), w \in \Sigma^*$ 

Ausgabe:  $erg \in \{True, False\}$ 

Nachbedingung:  $erg = (w \in \mathcal{L}(G))$ 

Mit anderen Worten: Es muss eine Ableitung  $S \stackrel{*}{\to} w$  gefunden werden.

$$G = (N, T, \Pi, S)$$
 mit  
 $N = \{S, A\}$   
 $T = \{a, b, c\}$   
 $\Pi = \{S \rightarrow cAb$   
 $A \rightarrow ab \mid a\}$ 

$$G = (N, T, \Pi, S)$$
 mit  $N = \{S, A\}$   $T = \{a, b, c\}$   $\Pi = \{S \rightarrow cAb \ A \rightarrow ab \mid a\}$   $W = cab$ 

$$G = (N, T, \Pi, S)$$
 mit

 $N = \{S, A\}$ 
 $T = \{a, b, c\}$ 
 $\Pi = \{S \rightarrow cAb$ 
 $A \rightarrow ab \mid a\}$ 
 $w = cab$ 
 $\Rightarrow Ableitung : S \rightarrow cAb \rightarrow cab$ 

$$G = (N, T, \Pi, S)$$
 mit
 $N = \{S, A\}$ 
 $T = \{a, b, c\}$ 
 $\Pi = \{S \rightarrow cAb$ 
 $A \rightarrow ab \mid a\}$ 
 $w = cab$ 
 $\Rightarrow Ableitung : S \rightarrow cAb \rightarrow cab$ 
Ergebnis:  $erg = True$ 

#### Ansatz für einen Parser

Programmiersprachen werden als kontextfreie Grammatiken mit kontextsensitiven Nebenbedingungen beschrieben

#### Ansatz für einen Parser

Programmiersprachen werden als kontextfreie Grammatiken mit kontextsensitiven Nebenbedingungen beschrieben

- Cocke-Younger-Kasami–Algorithmus Nachteil:
  - ► Voraussetzung: Grammatik Chomsky-Normalform
  - Aufwand  $O(n^3)$

#### Ansatz für einen Parser

Programmiersprachen werden als kontextfreie Grammatiken mit kontextsensitiven Nebenbedingungen beschrieben

- Cocke-Younger-Kasami–Algorithmus Nachteil:
  - ► Voraussetzung: Grammatik Chomsky-Normalform
  - Aufwand  $O(n^3)$
- (Nicht deterministische) Push–Down–Automaten Nachteil:
  - Polynomialer Aufwand
- ⇒ Betrachtung einer Teilmenge der Chomsky-2–Sprachen (LR–Sprachen) (sind äquivalent zu den Sprachen, die durch deterministische Push–Down–Automaten erkannt werden.
- ⇒ Effiziente Implementierung möglich.

## Ableitungsbaum

Man kann die Ableitung eines Wortes als Baum betrachten.

## Ableitungsbaum

Man kann die Ableitung eines Wortes als Baum betrachten.

#### Aufbau:

Top-down: Linksableitungen (man erhält eine Ableitung, bei der immer das am weitesten links stehende Nichtterminal abgeleitet wird)

## Ableitungsbaum

Man kann die Ableitung eines Wortes als Baum betrachten.

#### Aufbau:

Top-down: Linksableitungen (man erhält eine Ableitung, bei der immer das am weitesten links stehende Nichtterminal abgeleitet wird)

Bottom-Up: Rechtsableitungen (man erhält (rückwärts) eine Ableitung bei der immer das am weitesten rechts stehende Nichtterminal abgeleitet wird)

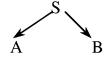
$$G = (N, T, \Pi, S)$$
 mit  $N = \{S, A, B, C\}$   $T = \{a, b, c\}$  und  $\Pi = \{S \rightarrow AB, A \rightarrow a, B \rightarrow bC, C \rightarrow c\}$ 

Eingabewort: w = abc

$$G = (N, T, \Pi, S)$$
 mit  $N = \{S, A, B, C\}$   $T = \{a, b, c\}$  und  $\Pi = \{S \rightarrow AB, A \rightarrow a, B \rightarrow bC, C \rightarrow c\}$ 

Eingabewort: w = abc

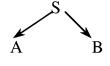
nexttoken() = a

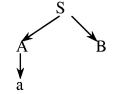


$$G = (N, T, \Pi, S)$$
 mit  $N = \{S, A, B, C\}$   $T = \{a, b, c\}$  und  $\Pi = \{S \rightarrow AB, A \rightarrow a, B \rightarrow bC, C \rightarrow c\}$ 

Eingabewort: w = abc

$$nexttoken() = a$$

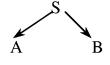


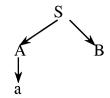


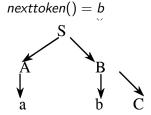
$$G = (N, T, \Pi, S)$$
 mit  $N = \{S, A, B, C\}$   $T = \{a, b, c\}$  und  $\Pi = \{S \rightarrow AB, A \rightarrow a, B \rightarrow bC, C \rightarrow c\}$ 

Eingabewort: w = abc

$$nexttoken() = a$$



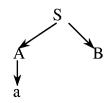


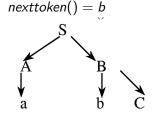


$$G = (N, T, \Pi, S)$$
 mit  $N = \{S, A, B, C\}$   $T = \{a, b, c\}$  und  $\Pi = \{S \rightarrow AB, A \rightarrow a, B \rightarrow bC, C \rightarrow c\}$ 

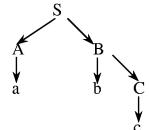
Eingabewort: w = abc

nexttoken() = a









## Recursive Decent-Syntaxanalyse

- ► Eingabe wird durch eine Menge rekursiver Funktionen abgebarbeitet.
- ▶ Jedem Nichtterminal der Grammatik entspricht eine Funktion.
- Die Folge der Funktionsaufrufe bestimmt implizit den Ableitungsbaum.

$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ Exp \rightarrow let \ var = Exp \ in \ Exp + Exp \ | \ var \ | \ digits \}$$

Problem: Linksrekursion

$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ Exp \rightarrow let \ var = Exp \ in \ Exp \\ | Exp + Exp \\ | var \\ | digits \}$$

Problem: Linksrekursion Eingabewort 1 + 1 + 1

$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ Exp \rightarrow let \ var = Exp \ in Exp \ | Exp + Exp \ | var \ | digits \}$$

# Problem: Linksrekursion Eingabewort 1 + 1 + 1



$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ \begin{array}{l} \textit{Exp} \rightarrow \textit{let var} = \textit{Exp in Exp} \\ | \textit{Exp} + \textit{Exp} \\ | \textit{var} \\ | \textit{digits} \} \end{array}$$

# Problem: Linksrekursion Eingabewort 1 + 1 + 1



$$G = (N, T, \Pi, S)$$
 mit

$$\Pi = \{ \begin{array}{ll} \textit{Exp} \rightarrow \textit{ let var } = \textit{Exp in Exp} \\ \mid \textit{Exp} + \textit{Exp} \\ \mid \textit{var} \\ \mid \textit{digits} \} \end{array}$$

# Problem: Linksrekursion

Eingabewort 1 + 1 + 1





$$G = (N, T, \Pi, S)$$
 mit

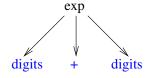
$$\Pi = \{ \begin{array}{l} \textit{Exp} \rightarrow \textit{ let var} = \textit{Exp in Exp} \\ \mid \textit{Exp} + \textit{Exp} \\ \mid \textit{var} \\ \mid \textit{digits} \} \end{array}$$

# Problem: Linksrekursion

Eingabewort 1 + 1 + 1





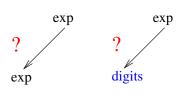


$$G = (N, T, \Pi, S)$$
 mit

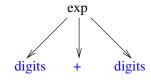
$$\Pi = \{ \begin{array}{l} \textit{Exp} \rightarrow \textit{ let var} = \textit{Exp in Exp} \\ \mid \textit{Exp} + \textit{Exp} \\ \mid \textit{var} \\ \mid \textit{digits} \} \end{array}$$

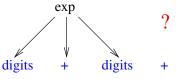
# Problem: Linksrekursion

Eingabewort 1 + 1 + 1









digits

$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ Exp \rightarrow let \ var = Exp \ in Exp \ | Exp + Exp \ | var \ | digits \}$$

Problem: Linksrekursion Eingabewort 1 + 1 + 1

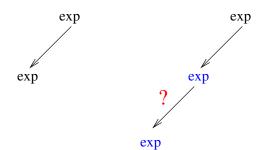
$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ Exp \rightarrow let \ var = Exp \ in Exp \ | Exp + Exp \ | var \ | digits \}$$

# Problem: Linksrekursion Eingabewort 1 + 1 + 1



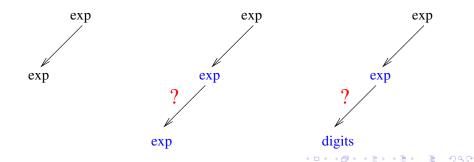
$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ Exp \rightarrow let \ var = Exp \ in \ Exp + Exp \ | \ var \ | \ digits \}$$

# Problem: Linksrekursion Eingabewort 1 + 1 + 1



$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ Exp \rightarrow let \ var = Exp \ in \ Exp \ | Exp + Exp \ | var$$

Problem: Linksrekursion Eingabewort 1 + 1 + 1



digits}

#### Elimination der Linksrekursion

$$G = (N, T, \Pi, S)$$
 mit 
$$\Pi = \{ Exp \rightarrow let \ var = Exp \ in \ Exp \ | Exp + Exp \ | var \ | var \ | digits \}$$

#### Elimination der Linksrekursion

```
G = (N, T, \Pi, S) mit
                         \Pi = \{ Exp \rightarrow let \ var = Exp \ in \ Exp \}
                               | Exp + Exp
| var
| digits}
                        \Pi = \{ Exp \rightarrow TExp Exp' \}
                                 Exp' \rightarrow + TExp Exp'
                                 TExp \rightarrow let \ var = Exp \ in \ Exp
```

 $\Pi = \{S \rightarrow AB, A \rightarrow a, B \rightarrow bC, C \rightarrow c\}$ Eingabewort: w = abc

 $\Pi = \{S \to AB, A \to a, B \to bC, C \to c\}$ 

Eingabewort: w = abc

nexttoken() = a wende an:  $A \rightarrow a$ 



$$\Pi = \{S \rightarrow AB, A \rightarrow a, B \rightarrow bC, C \rightarrow c\}$$

Eingabewort: w = abc

$$nexttoken() = a$$
  $nexttoken() = b$   
wende  $an: A \rightarrow a$   $nexttoken() = c$ 

nexttoken() = a nexttoken() = bwende an:  $C \rightarrow c$ 

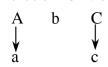


$$\Pi = \{S \rightarrow AB, A \rightarrow a, B \rightarrow bC, C \rightarrow c\}$$

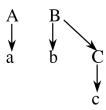
Eingabewort: w = abc

$$nexttoken() = a$$
  $nexttoken() = b$   
wende an:  $A \rightarrow a$   $nexttoken() = c$ 

wende an:  $C \rightarrow c$ 



nexttoken() = a nexttoken() = b wende  $an: B \rightarrow bC$ 



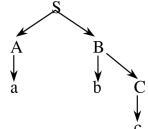
$$\Pi = \{S \to AB, A \to a, B \to bC, C \to c\}$$

Eingabewort: w = abc

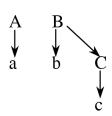
wende an:  $A \rightarrow a$  nexttoken() = c

wende an:  $C \rightarrow c$ 

wende an:  $S \rightarrow AB$ 



nexttoken() = a nexttoken() = b wende  $an: B \rightarrow bC$ 



$$\Pi = \{S \to AB, A \to a, B \to bC, C \to c\}$$

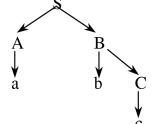
Eingabewort: w = abc

wende an:  $A \rightarrow a$  nexttoken() = c

wende an:  $C \rightarrow c$ 

nexttoken() = a nexttoken() = b wende  $an: B \rightarrow bC$ 

wende an:  $S \rightarrow AB$ 



Betrachtet man die Konstruktion rückwärts:

$$\underline{S} \rightarrow A\underline{B} \rightarrow Ab\underline{C} \rightarrow \underline{A}bc \rightarrow abc$$

```
G = (N, T, \Pi, S) mit N = \{S, A, B\}, T = \{a, b, c, d, e\} und \Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\} Eingabestring: w = abbcde
```

```
G = (N, T, \Pi, S) mit N = \{S, A, B\}, T = \{a, b, c, d, e\} und \Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\} Eingabestring: w = abbcde
```

#### Bottom-Up Syntaxanalyse:

a.bbcde

```
G = (N, T, \Pi, S) mit N = \{S, A, B\},\ T = \{a, b, c, d, e\} und \Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\}
```

Eingabestring: w = abbcde

#### Bottom-Up Syntaxanalyse:

 $a.bbcde \leftarrow ab.bcde$ 

```
G = (N, T, \Pi, S) mit N = \{S, A, B\}, T = \{a, b, c, d, e\} und \Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\} Eingabestring: w = abbcde
```

#### Bottom-Up Syntaxanalyse:

 $a.bbcde \leftarrow ab.bcde \leftarrow aA.bcde$ 

```
G = (N, T, \Pi, S) mit N = \{S, A, B\}, T = \{a, b, c, d, e\} und \Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\} Eingabestring: w = abbcde
```

#### Bottom-Up Syntaxanalyse:

 $a.bbcde \leftarrow ab.bcde \leftarrow aA.bcde \leftarrow aAb.cde$ 

```
G = (N, T, \Pi, S) mit

N = \{S, A, B\},

T = \{a, b, c, d, e\} und

\Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\}
```

Eingabestring: w = abbcde

#### Bottom-Up Syntaxanalyse:

 $a.bbcde \longleftarrow ab.bcde \longleftarrow aA.bcde \longleftarrow aAb.cde \longleftarrow aAbc.de$ 

```
G = (N, T, \Pi, S) mit N = \{S, A, B\}, T = \{a, b, c, d, e\} und \Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\} Eingabestring: w = abbcde

Bottom-Up Syntaxanalyse: a.bbcde \longleftarrow aA.bcde \longleftarrow aAb.cde \longleftarrow aAbc.de
```

aAA cde

```
G = (N, T, \Pi, S) mit

N = \{S, A, B\},

T = \{a, b, c, d, e\} und

\Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\}
```

Eingabestring: w = abbcde

#### Bottom-Up Syntaxanalyse:

$$a.bbcde \leftarrow ab.bcde \leftarrow aA.bcde \leftarrow aAb.cde \leftarrow aAbc.de \leftarrow Shift-Reduce-Konflikt$$

$$aAA cde$$

 $aA.de \leftarrow aAd.e$ 

```
G = (N, T, \Pi, S) mit

N = \{S, A, B\},

T = \{a, b, c, d, e\} und

\Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\}
```

Eingabestring: w = abbcde

#### Bottom-Up Syntaxanalyse:

$$a.bbcde \leftarrow ab.bcde \leftarrow aA.bcde \leftarrow aAb.cde \leftarrow aAbc.de \leftarrow Shift-Reduce-Konflikt$$

$$aAA.cde$$

 $aA.de \leftarrow aAd.e \leftarrow aAB.e$ 

```
G = (N, T, \Pi, S) mit

N = \{S, A, B\},

T = \{a, b, c, d, e\} und

\Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\}
```

Eingabestring: w = abbcde

#### Bottom-Up Syntaxanalyse:

$$a.bbcde \leftarrow ab.bcde \leftarrow aA.bcde \leftarrow aAb.cde \leftarrow aAbc.de \leftarrow Shift-Reduce-Konflikt$$

$$aAA.cde$$

$$aA.de \leftarrow aAd.e \leftarrow aAB.e$$

Reduce-
Reduce-
Konflikt

 $aAA = e$ 

```
G = (N, T, \Pi, S) mit
N = \{S, A, B\},\
T = \{a, b, c, d, e\} und
\Pi = \{S \rightarrow aABe, A \rightarrow Abc \mid b \mid d, B \rightarrow d\}
Eingabestring: w = abbcde
Bottom-Up Syntaxanalyse:
a.bbcde \leftarrow ab.bcde \leftarrow aA.bcde \leftarrow aAb.cde \leftarrow aAbc.de \leftarrow
                                                               aAA cde
aA.de \leftarrow aAd.e \leftarrow aAB.e \leftarrow aABe \leftarrow S
                        aAA.e
```

## Konfliktlösungsansätze

- LR(0): Es muss ohne Vorausschau möglich sein zu entscheiden, ob geshiftet oder reduziert wird.
- SLR(1): An Hand der Bildung der Menge aller möglichen folgenden Terminalsymbole auf ein Nichtterminal, wird entschieden ob reduziert wird.
- LR(1): Für jeden möglichen Ableitungsschritt in einer Produktion wird die Menge der darauffolgenden Terminalsymbolde zur Unterscheidung betrachtet.
- LALR(1): Es werden alle Mengen von LR(1)–Elementen zusammengefasst, die den gleichen Ableitungsschritt vollziehen.

$$LR(0) \subset SLR(1) \subset LALR(1) \subset LR(1) (= \mathcal{L}(\mathcal{DPDA}))$$



# Parsergenerator—Tool jaooy

► Lexer/Scanner fasst die Strings zusammen, die an einer Stelle als Terminalsymbol in einer Grammatik stehen dürfen:

Lexeme mit gleicher Bedeutung: z.B.: [<html>, <HTML>]
Lexeme unterschiedlicher Bedeutung: z.B.: Text zwischen Tags

Ziel: Definition einer Datenstruktur, die diesem Rechnung trägt.

## Klasse yyTokenclass

```
class yyTokenclass {
  public int tokennr;
  public Object value;
  yyTokenclass () {
      this.tokennr=-1;
  yyTokenclass (Object o) {
      this.value = o;
```

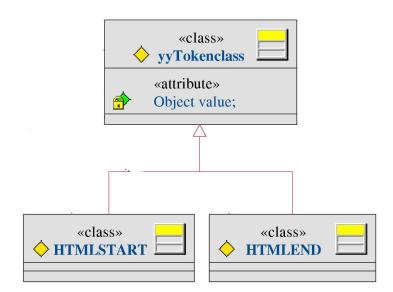
Alle Tokens erben von yyTokenclass.

## JLex-Spezifikation mit Tokens

```
%%
%class browserlexer
%eofval{
  return new EOF();
%eofval}
ws = [ \t r n] +
%%
"<"(h|H)(t|T)(m|M)(1|L)">" { return new HTMLSTART(yytext());}
"</"(h|H)(t|T)(m|M)(1|L)">" { return new HTMLEND(yytext());}
[^\<]+
              { System.out.println(yytext());
                       return new WORT(yytext());}
```

Wo kommen die Tokens (Java-Klassen) her?

#### Token-Klasse



## jay-Tool (jaooy)

### Eingabe (Grammatik)

- Namen der Tokenklassen (Terminale)
- Produktionen

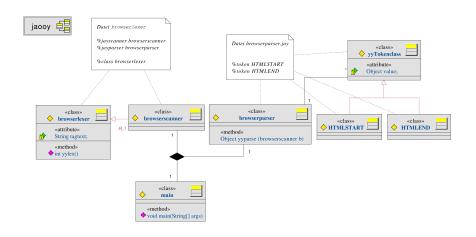
#### Ausgabe

 Kellerautomat, der die Sprache der eingegebenen Grammatik erkennt als Java-Programm

## jay-Spezifikation

```
%{
class browserparser {
%}
%token HTMLSTART
%token HTMLEND
%%
S : HTMLSTART Head Body HTMLEND { }
Head : HEADSTART Title    HEADEND { }
. . .
%%
```

# Klassendiagramm JLex und jay



## JLex-Spezifikation mit Klassendeklarationen

```
%%
%jayscanner browserscanner
%jayparser browserparser
%class browserlexer
%eofval{
  return new EOF();
%eofval}
ws = [ \ \ \ \ \ \ \ \ \ \ \ ] +
%%
"<"(h|H)(t|T)(m|M)(1|L)">" { return new HTMLSTART(yytext());}
"</"(h|H)(t|T)(m|M)(1|L)">" { return new HTMLEND(yytext());}
[^\<]+
              { System.out.println(yytext());
                        return new WORT(yytext());}
```

### Java-Klasse für den Aufruf

```
public class main {
    public static void main (String [] args) {
        browserscanner scanner =
            new browserscanner
                  (new java.io.InputStreamReader (System.in));
        browserparser parser = new browserparser() ;
        try {
            parser.yyparse(scanner);
        catch (Exception e) {
            e.printStackTrace();
```

#### Makefile

```
Main.class: Main.java yyTokenclass.class browserscanner.class
        javac Main.java
yyTokenclass.class: browserparser.java
        javac browserparser.java
browserscanner.class: browserlexer.java
        javac browserlexer.java
browserlexer.java: browserlexer
        java -cp JLex2.jar JLex2.Main browserlexer
browserparser.java: browserparser.jay skeleton.jaooy
        jaooy -v browserparser.jay < skeleton.jaooy > browserparser.java
clean:
        rm -f *.class browserlexer.java browserparser.java
```

## Aufgabe

Entwickeln Sie eine Mini-html-Parser als Grundlage für einen Browser.

- ▶ Definieren Sie eine Grammatik für Mini-html.
- Ergänzen Sie das JLex- und das Jay-File dementsprechend.