

LUCRAREA 6

Analiza sintactică folosind CUP

Cuprinsul lucrării:

- Generatorul de parseri sintactici CUP
- Sintaxa unui fișier
- Definirea gramaticii
- Conflicte de tip shift/reduce

Unele exemple și figuri din acest capitol au fost preluate și adaptate din [5].

6.1 Generatorul de parseri sintactici CUP

CUP (Constructor of Useful Parsers) este un program care generează parseri sintactici transformând o gramatică independentă de context într-un program Java care parsează secvența de simboluri de intrare în funcție de regulile gramaticii. Gramatica este specificată într-un fișier cu extensia .cup, iar pentru fiecare regulă a gramaticii este menționată și o acțiune care se va executa când regula respectivă ajunge să fie aplicată.

CUP crează un parser de tipul LALR (Look-Ahead Left-to-Right), care reprezintă o formă particulară de parser LR(1). Parserul sintactic astfel generat trebuie integrat cu un analizor lexical. În acest sens, în prezenta lucrare se va utiliza un lexer generat folosind JLex, și se va arăta cum CUP

și JLex pot fi utilizate împreună pentru a scrie un program complet care implementează analiza lexicală și sintactică.

Pentru acest laborator se va folosi un mediu de lucru gata configurat care conține atât codul sursă pentru CUP și JLex cât și un exemplu minimal de compilator descris prin cele două fișiere de specificații `example.lex` și `example.cup`[4]. Fișierele pot fi descărcate de la adresa:

<https://github.com/gmh33/Java-CUP-jLex-Example>

6.2 Integrarea JLex și CUP. Clasa Symbol

Pentru a integra JLex și CUP este nevoie de a specifica în cadrul fișierului JLex faptul că ieșirea lexer-ului va fi interfațată cu parserul generat de CUP. Acest lucru înseamnă în primul rând că atomii lexicali generați de lexer nu vor mai fi de tipul specificat de utilizator, ci vor fi obligatoriu de tipul `Symbol`, o clasă predefinită din CUP. Pentru aceasta, în prima secțiune a fișierului JLex (cea dedicată codului de utilizator) se adaugă linia: `import java_cup.runtime.Symbol;`

În continuare, în a doua secțiune a fișierului se introduce directiva `%cup` care va specifica JLex faptul că atomii lexicali vor fi trimiși mai departe către un program generat de CUP. Orice altă directivă care specifică tipul de atomi, de exemplu `%integer` trebuie îndepărtată, întrucât tipul de atomi care va fi returnat este `Symbol`.

Un obiect de tip `Symbol` are mai multe date membre, ce pot fi văzute la sfârșitul fișierului `CUP/java_cup/runtime/Symbol.java`. Dintre acestea, de interes sunt `sym` și `value`. Data membră `sym` va primi o valoare din clasa `sym`, ce va fi creată de CUP, iar `value` va primi o valoare dată de utilizator în fișierul JLex și este de tipul generic `Object`, ceea ce înseamnă că poate primi orice tip de dată.

De interes sunt și constructorii clasei `Symbol`. Aceștia sunt apelați în fișierul JLex atunci când se specifică acțiunea ce urmează a fi executată la găsirea unui atom de un anumit tip. Dacă atomul nu are o valoare asociată (spre exemplu, un șir de caractere sau o valoare numerică), se va apela constructorul cu prototipul `public Symbol(int sym_num)`, spre exemplu la întâlnirea unui operator: `return new Symbol(sym.PLUS)`. În cazul în care atomului trebuie să i se asocieze o valoare, se apelează constructorul `public Symbol(int id, Object o)`, ca de exemplu la întâlnirea unui identificator, situație în care trebuie returnat pe lângă tipul de atom și denumirea identificatorului sub forma unui `String`: `return new Symbol(sym.ID, new String(yytext()))`.

6.3 Fișierul CUP

Fișierul de specificații CUP are de regulă extensia `.cup` și conține descrierea gramaticii independente de context pe baza căreia CUP va crea clasele `parser` și `sym` (care stau la baza parser-ului propriu-zis). Acest fișier începe de regulă cu liniile de cod care vor fi copiate direct în parser. În exemplul dat, primele linii se referă la includerea unor anumite clase Java:

```
import java.cup.runtime.*;
import java.util.*;
```

În Codul sursă 6.1, se utilizează directiva `parser code` pentru a crea funcția `main` direct în clasa `parser`.

Codul sursă 6.1: Crearea funcției `main` direct în clasa `parser`.

```
parser code {:
    /* this part is included in the parser class */
    public static void main (String args[]) throws
        Exception{
        new parser(new Yylex(System.in)).parse();
    }
:}
```

Constructorul apelat primește ca parametru un nou obiect de tip `Yylex`, prin aceasta specificându-se faptul că parserul integrează un lexer definit prin clasa `Yylex` ce trebuie în prealabil generată folosind `JLex`.

Urmează enumerarea simbolurilor gramaticii, mai întâi cele terminale (vezi Codul sursă 6.2). În general, acestea sunt aceleași cu atomii lexicali returnați de lexer, iar CUP copiază denumirile specificate în clasa `sym` și atribuie fiecăreia o valoare numerică. Orice simbol terminal menționat aici dar care nu e prezent în nicio regulă a gramaticii va determina afișarea unui mesaj de avertizare în etapa generării parserului.

Codul sursă 6.2: Enumerarea simbolurilor terminale.

```
terminal VAR,BEGIN,END,IF,THEN,WHILE,DO,ODD;
terminal EQ,NEQ,LT,LTE,GT,GTE,LPAREN,RPAREN,PERIOD;
terminal COMMA,SEMI,ASSIGN,PLUS,MINUS,MULT,DIV,PRINT;
terminal Integer INT;
terminal String ID;
```

Specificațiile CUP continuă cu enumerarea simbolurilor neterminale (vezi Codul sursă 6.3), fiecare dintre ele trebuind să se regăsească ulterior în membrul stâng al unei producții. La fel, orice simbol neterminal menționat aici dar nefolosit sau inaccesibil va genera un mesaj de avertizare la momentul generării parserului. Un lucru important de reținut este că simbolurile

terminale sunt enumerate în ordinea crescătoare a precedenței (priorității) acestora.

Codul sursă 6.3: Enumerarea simbolurilor neterminale.

```
non terminal Program program;
non terminal Block block;
non terminal VarLst var_lst;
non terminal StmtLst stmt_lst;
non terminal Statement statement;
non terminal Condition condition;
non terminal Expr expr, term, factor;
```

După enumerarea simbolurilor gramaticii urmează explicitarea producțiilor acestora (vezi Codul sursă 6.4). Fiecare simbol neterminal se regăsește în membrul stâng al cel puțin unei producții, iar fiecare producție este însoțită și de o acțiune semantică. Aceste acțiuni sunt folosite pentru a produce anumite efecte la aplicarea unei anumite reguli, mai exact pentru a genera codul de asamblare corespunzător. Regulile gramaticii sunt aplicate după potrivirea membrului drept, datorită faptului că parserul este unul de tip LR.

Codul sursă 6.4: Explicitarea producțiilor gramaticii.

```
expr ::=
    term:t
        {: RESULT = t; :}
    | PLUS term:t
        {: RESULT = t; :}
    | MINUS term:t
        {: RESULT = new InvExpr(t); :}
    | expr:e PLUS term:t
        {: RESULT = new PlusExpr(e,t); :}
    | expr:e MINUS term:t
        {: RESULT = new MinusExpr(e,t); :}
    ;
```

În cadrul fiecărei producții a gramaticii, fiecare terminal are asociată o variabilă care preia valoarea terminalului respectiv. Această variabilă este definită prin simbolul : urmat de numele variabilei. Spre exemplu, `expr ::= term:t`, unde `t` este variabila care are valoarea simbolului terminal `term`. Folosind aceste variabile se pot specifica anumite acțiuni care implică simbolurile terminale. Pentru același exemplu, acțiunea este `{: RESULT = t; :}`, ceea ce înseamnă că valoarea atribuită lui `RESULT` este transmisă automat către neterminalul din partea stângă a producției, în cazul de față `expr`.

În cazul unei alte producții a simbolului neterminal `expr`, și anume `expr ::= MINUS term:t`, se observă faptul că acțiunea semantică utilizează variabila `t` dar în același timp crează un nou obiect care este trimis prin `RESULT` către `expr`, și anume: `{: RESULT = new InvExpr(t); :}`. Crearea de astfel de obiecte semantice, sau în alte situații funcții de acțiune semantică se impune în situațiile în care acțiunile care trebuie întreprinse pentru o anumită producție au un nivel de complexitate ridicat. În acest fel, fișierul de specificații CUP își păstrează o structură lizibilă.

Definirea și explicitarea acestor obiecte și funcții de acțiune semantică este făcută într-un fișier java separat, în cazul exemplului din acest laborator, este vorba de fișierul `Program.java`.

6.4 Generarea parserului

Mediul de lucru pentru prezentul laborator are automatizat procesul de compilare și generare a programelor aferente folosind un fișier de tip `makefile` (denumit `makefile`). Astfel, în linia de comandă se introduc următoarele comenzi:

```
make init
make all
```

care vor compila sursele JLex și Cup (`init`) precum și sursele compilatorului (`all`).

Fără utilizarea unui fișier de tip `makefile`, pașii necesari sunt descriși mai jos (ei se regăsesc în `makefile`, lucru verificabil prin deschiderea acestui fișier cu un editor de text):

```
javac JLex/*.java
javac java_cup/runtime/*.java
javac java_cup/*.java
javac -classpath . java_cup.Main example.cup
javac JLex.Main example.lex
mv example.lex.java Yylex.java
javac Program.java
javac sym.java parser.java Yylex.java
```

După compilarea surselor, sunt generate două clase Java care descriu implementarea parserului: `sym` și `parser`. Clasa `sym`, descrisă în fișierul `sym.java` cuprinde definirea de valori numerice întregi pentru fiecare din simbolurile terminale enumerate în fișierul de specificații CUP. Clasa `parser`, cuprinsă în fișierul `parser.java`, cuprinde descrierea tabelară a producțiilor, acțiunilor și regulilor de reducere corespunzătoare.

De asemenea, fișierul include și codul scris de utilizator la începutul fișierului de specificații CUP, în cazul nostru funcția `main`, care crează un obiect de tipul clasei (`parser`). Tot în fișierul `parser.java` se regăsește un constructor al clasei care primește ca parametru un obiect de tip `Scanner`. Acest `Scanner` e de fapt lexer-ul generat de JLex, fapt observabil în definiția clasei `Yylex` din fișierul `Yylex.java`. Prin aceste specificații, parserul știe că va folosi clasa `Yylex` ca lexer, cu alte cuvinte ca sursă de atomi lexicali.

După generare, parserul poate fi lansat în execuție furnizându-i-se un cod de intrare. În exemplul ilustrat de Codul sursă 6.5, fișierul `test.pl` conține un program exemplificativ pentru limbajul PL/0, care implică declararea de variabile, operații matematice cu acestea, afișări și o structură repetitivă (bucă de tip `while`).

Limbajul PL/0 este un limbaj minimalist (având o sintaxă mult simplificată¹), cu aplicații în domeniul educațional, mai exact în studiul și proiectarea compilatoarelor. A fost dezvoltat de informaticianul elvețian Niklaus Wirth și prezentat în cartea acestuia, *Algorithms + Data Structures = Programs* [14].

Codul sursă 6.5: Cod sursă PL/0.

```
var a,b,c;
begin
    a := 1;
    b := 2;
    print c;
    c := a + b;
    print c;
    if odd a then
        a := 0 - a;
    print a;
    b := 10;
    while b # 0 do
    begin
        print b;
        b := b - 1
    end
end.
```

Compilarea și lansarea în execuție a acestui program se realizează prin comanda:

```
make test
```

sau se poate explicita numele fișierului prin comanda:

```
make run < test.pl
```

¹<https://en.wikipedia.org/wiki/PL/0>

Din nou, este folosit fișierul `makefile` de automatizare a întregului flux. Comenzile de mai sus în fapt substituie comanda:

```
java classpath . parser < test.pl
```

După executarea parserului pe codul sursă din `test.pl`, pe ecran ar trebui să fie afișate valorile numerice din program în urma execuției comenzilor `print`. Urmăriți codul sursă și valorile afișate.

6.5 Temă de laborator

La specificațiile compilatorului de PL/0 prezentate în prezentul laborator, adăugați codul necesar, atât în fișierul `JLex` cât și în specificațiile `CUP` pentru a implementa instrucțiunea repetitivă `for` și lucrul cu metode (subprograme).