

LUCRAREA 5

Analiza lexicală folosind JLex

Cuprinsul lucrării:

- Generatori de analizori lexicali. JLex¹.
- Interfața JLex.
- Metodologia de lucru.
- Generarea de analizori lexicali folosind JLex.

Unele exemple și figuri din acest capitol au fost preluate și adaptate din [5] și [11].

5.1 Generatori de analizori lexicali. JLex

Un analizor lexical (*lexer* în limba engleză) împarte un șir de caractere de intrare în atomi lexicali. Procesul de elaborare a unui astfel de analizor poate fi extrem de dificil dacă este realizat manual, astfel încât au fost concepute instrumente software dedicate pentru acest lucru, care poartă numele de generatoare de analizori lexicali.

Cel mai cunoscut astfel de analizor este Lex [6]. Acesta este un instrument software ce rulează sub Linux și este destinat a lucra în limbajul de programare C. Primește ca intrare un fișier care conține specificațiile analizorului lexical și crează un fișier sursă C asociat cu acel analizor.

¹<https://www.cs.princeton.edu/~appel/modern/java/JLex/>

JLex este un program care se bazează pe Lex și care generează codul Java corespunzător unui lexer.

Metodologia de lucru este prezentată în Figura 5-1. Se proiectează specificațiile lexer-ului conform cerințelor într-un fișier cu extensia `.jlex`. Acest fișier este procesat de generatorul JLex (clasa `JLex.Main`), care pe baza specificațiilor va genera codul sursă (Java) al lexer-ului propriu zis. Acest cod est mai apoi compilat folosind `javac` pentru a obține fișierul de tip `.class` al lexerului, care este ulterior instanțiat și folosit în programul propriu zis. În final, lexer-ul astfel generat va produce, pe baza șirului de caractere primit la intrare, un flux de atomi lexicali recunoscuți.

5.2 Configurarea mediului de lucru

- (a) Verificați că Java este instalat (comenzile `java` și `javac` funcționează într-un terminal/comand prompt)
- (b) Creați în spațiul de lucru un director `lab4` și în acesta încă un director `lab4/JLex`
- (c) Descărcați clasa `Main.java` de la adresa: <https://www.cs.princeton.edu/~appel/modern/java/JLex/current/Main.java> <https://www.cs.princeton.edu/~appel/modern/java/JLex/current/Main.java> în directorul `lab4/JLex`
- (d) Compilați `Main.java` ca o clasă Java obișnuită:

```
javac Main.java
```

În urma acestei comenzi ar trebui să fie generate mai multe fișiere de tip clase Java, printre care și `Main.class`
- (e) Urmăriți în continuare laboratorul. Fișierele sursă le veți scrie în directorul `lab4/`

5.3 Interfața JLex

Un fișier sursă JLex are de regulă extensia `.jlex`, iar codul este structurat în trei părți despărțite prin caracterele `%%`. Un exemplu simplist de fișier sursă JLex este prezentat în Codul sursă 5.1 .

Codul sursă 5.1: Exemplu de fișier sursă JLex.

```
%%  
%integer
```

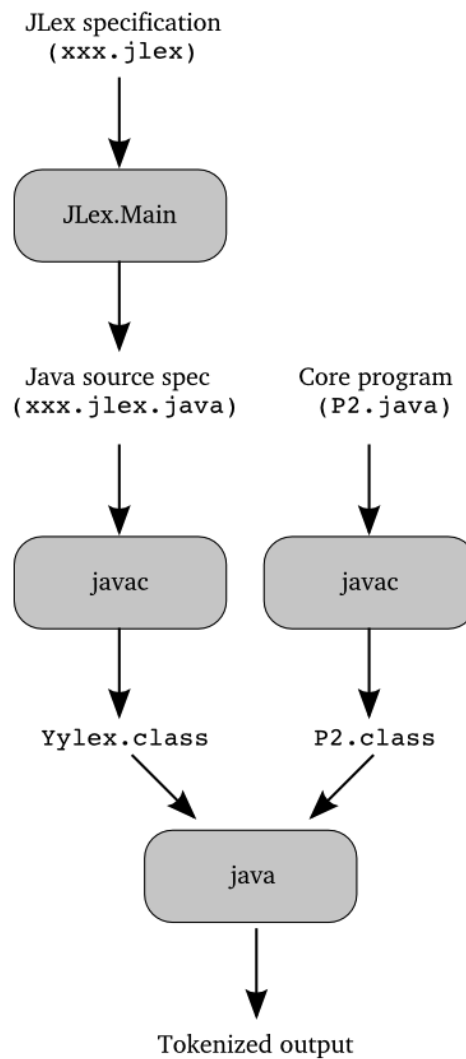


Figura 5-1: Metodologia de lucru JLex.

```

numbers = 0|[1-9][0-9]*
whitespace = [\ \t\n\r]+
%%
{whitespace} { System.out.println("S-a gasit un spatiu:");
}
{digits} { return new Token(INT, Integer.parseInt(yytext()))
; }
. { System.out.println("ERROR!"); }
  
```

Prima secțiune conține cod Java care este copiat direct, fără a fi executat în prealabil, în clasa Java Yylex care este numele lexer-ului. În cazul exemplului de mai sus, această secțiune este goală. A doua secțiune cuprinde diferite definiții și directive JLex. Directivele sunt precedate întotdeauna de caracterul %. În exemplul dat, secțiunea conține o directivă și trei definiții.

Directiva `%integer` specifică faptul că atomii care sunt returnați în urma analizei lexicale vor fi de tip întreg(`int`). Specificarea acestei directive mai conduce la definirea de către JLex a constantei `YYEOF = -1` în clasa Yylex, constantă ce va fi returnată în mod automat de Yylex la întâlnirea sfârșitului de fișier (end of file).

Cele două definiții (`numbers` și `whitespace`) definesc valori numerice și respectiv spații folosind expresii regulate. Astfel, `numbers = 0|[1-9][0-9]*` definește un număr ca fiind fie cifra 0, fie orice succesiune de oricâte cifre, dar măcar una, din care prima este obligatoriu diferită de 0. A doua definiție, `whitespace = [\\t\\n\\r]+`, include caracterele de tip spațiu, tab și sfârșit de linie.

A treia secțiune din fișierul `.lex` specifică expresiile regulate pentru care se caută potriviri și acțiunile specifice în fiecare caz. În exemplul dat, expresiile regulate vizează spațiile (`whitespace`) și numerele (`numbers`). În cazul identificării de spații, se afișează un mesaj corespunzător care anunță acest lucru dar nu se execută alte acțiuni. La identificarea de valori numerice, nu se afișează nimic, în schimb se returnează un atom (Token) cu valoarea numerică identificată pe baza expresiei regulate. Se observă caracterul punct (.) aflat pe ultima linie din fișierul `.lex`. Acest caracter cuprinde orice caracter cu excepția unei noi linii din fișier și este utilizat ultimul pentru a semnala printr-un mesaj de eroare faptul că s-a găsit un caracter ce nu a putut fi potrivit cu expresiile regulate specificate mai sus.

Trebuie reținut faptul că în ultima secțiune din fișierul `.lex` contează ordinea în care sunt specificate expresiile regulate, această ordine marcând prioritatea cu care sunt potrivite caracterele din fișier. Dacă un șir de caractere se potrivește cu mai multe expresii regulate, JLex alege expresia aflată „mai sus” în listă, deci cu prioritate mai mare.

5.4 Generarea de analizori lexicali folosind JLex

5.4.1 Generarea unui scanner de spații

Pentru început se va implementa un lexer care va scana spațiile dintr-un fișier text (Codul sursă 5.2).

Codul sursă 5.2: Codul sursă al unui scanner de spații.

```
class IntLit{
}
%%
WHITESPACE = [\n\t\040]
%type IntLit
%function next_token
%%
{WHITESPACE}+ { System.out.println("Am gasit un spatiu!");
}
```

În a doua secțiune a fișierului .jlex sunt incluse o definiție și două directive. Definiția `WHITESPACE = [\n\t\040]` cuprinde toate caracterele de tip linie nouă, tab, sau spațiu (definit prin caracterul `\040`). Directiva următoare (`%IntLit`) specifică faptul că lexer-ul va returna atomi lexicali de tip `IntLit`. Acest tip de atom trebuie definit în prima secțiune a fișierului jlex (lucru realizat prin declarația `class IntLit {}`).

Prin directiva `%function next_token()` este specificat numele metodei care va fi generată și folosită pentru a returna atomi lexicali.

A treia secțiune specifică, folosind expresii regulate și definițiile anterioare, ce șiruri de caractere vor fi găsite de lexer. Astfel, expresia `{WHITESPACE}+` va găsi toate șirurile de cel puțin un spațiu, iar pentru fiecare șir găsit se va executa acțiunea definită în continuare `{ System.out.println("Am gasit spatii!"); }`.

JLex procesează specificațiile din fișierul .jlex și le transformă în codul Java care implementează propriu-zis scannerul definit în mod abstract prin aceste specificații.

Pentru aceasta, la momentul compilării, compilatorului de Java trebuie să i se specifice locația pachetului `Jlex Main.class`. Spre exemplu, dacă directorul de lucru în care se găsesc fișierele .lex este `lab-jlex`, iar sursele JLex se găsesc în `../lab-jlex/JLex`, pentru a compila sursa comanda este următoarea:

```
java JLex.Main whitespace.jlex
```

Dacă nu au fost erori, rezultatul ar trebui să fie similar cu cel de mai jos:

```
Processing first section -- user code.
Processing second section -- JLex declarations.
Processing third section -- lexical rules.
Creating NFA machine representation.
NFA comprised of 9 states.
```

Working on character classes...:
NFA has 3 distinct character classes.
Creating DFA transition table.
Working on DFA states....
Minimizing DFA transition table.
3 states after removal of redundant states.
Outputting lexical analyzer code.

Se pot observa din mesajele returnate de compilator următoarele:

1. Procesarea secvențială a celor trei secțiuni ale codului sursă .jlex descrise anterior în laborator: cod utilizator, declarații JLex și reguli lexicale.
2. Metodologia de lucru este următoarea: pe baza specificațiilor introduse, JLex construiește un AFN, care este mai apoi convertit la un AFD și în final automatul rezultat este minimizat.

Puteți vizualiza codul Java produs de JLex pe baza specificațiilor, deschizând într-un editor de text fișierul `whitespace.jlex.java`. Studiați fișierul:

1. Observați diferența între specificațiile minimale din fișierul .jlex și dimensiunea și complexitatea codului rezultat
2. Fișierul definește două clase: `IntLit` (tipul de atom definit de către noi) și `Ylex` („motorul” propriu-zis al lexer-ului)

Acest fișier Java generat trebuie compilat:

```
javac whitespace.jlex.java
```

Compilarea va produce două fișiere de tip clasă Java, corespunzătoare claselor identificate mai sus: `IntLit.class` și `Ylex.class`. Aceste clase nu pot fi utilizate pur și simplu. Este nevoie de scrierea unui program Java în care acestea sunt instanțiate și folosite corespunzător. Acest lucru este realizat în fișierul `Driver.java`, conform Codului sursă 5.3.

Codul sursă 5.3: Driver.java

```
import java.io.*;
public class Driver{
    public static void main(String[] args) throws Exception
    {
        File f = new File("in.txt");
        FileInputStream fis = new FileInputStream (f);

        Ylex scanner = new Ylex(fis);
```

```

        do{
            scanner.next_token();
        } while(true);
    }
}

```

Programul deschide fișierul `in.txt` ca obiect de tip `InputStream` deoarece acesta este tipul care trebuie furnizat lexerului. În continuare, sunt citați atomi lexicali din fișier într-o buclă infinită. Citirea următorului atom se realizează prin apelul funcției `scanner.next_token()`.

Pasul următor îl reprezintă compilarea fișierului sursă `Driver.java`:

```
javac Driver.java
```

Compilarea a generat fișierul `Driver.class` care poate fi lansat în execuție. Înainte de a face acest lucru, trebuie creat fișierul `in.txt`. Creați acest fișier în care introduceți spații de toate tipurile (linie nouă, tab, spațiu). Apoi se poate executa clasa `Driver` prin comanda:

```
java Driver
```

Programul ar trebui să afișeze un mesaj de tipul:

```
Am gasit un spatiu!
```

Exerciții:

1. De ce programul rămâne blocat în execuție?
2. De ce este identificat un singur atom lexical și în consecință afișat un singur mesaj deși fișierul conține mai multe spații?
3. Modificați fișierul `in.txt` introducând și alte caractere în afară de spații. Rulați din nou `Driver`. Observați rezultatul.
4. Modificați fișierul `whitespace.jlex` astfel încât, în condițiile de la punctul 3, programul să nu își mai termine execuția printr-o eroare ci să afișeze un mesaj corespunzător.

5.4.2 Generarea unui scanner care însumează valorile numerice

În cele ce urmează vom modifica fișierul `whitespace.jlex` pentru a recunoaște un flux de valori numerice întregi, conform Codului sursă 5.4. Vom redenumi noul fișier ca `sum.jlex`.

Codul sursă 5.4: `sum.jlex`

```
class IntLit{
```

```

        private int value;
        public IntLit(String textIn) {
            value = Integer.parseInt(textIn);
        }
        public String toString() {
            return "Value: " + value;
        }
        public int valueGet(){
            return value;
        }
    }
%%
WHITESPACE = [\n\t\040]
DIGIT = [0-9]
not_star = [^*];
%type IntLit
%function next_token

%%
{DIGIT}+ { System.out.println("Am gasit o valoare numerica"
    );
        return new IntLit(yytext()); }
{WHITESPACE}+ { System.out.println("Am gasit un spatiu!");
    }
{not_star} { System.out.println("not star"); }
. { System.out.println("ERROR!"); }

```

Față de codul din `whitespace.jlex`, observăm o serie de modificări. În a doua secțiune a fost introdusă o nouă definiție care cuprinde cifrele `DIGIT = [0-9]`. Regula lexicală pe baza acestei definiții, ce se regăsește în secțiunea a treia a fișierului, specifică acțiunea ce se execută în momentul întâlnirii unui șir de unul sau mai multe cifre (`{DIGIT}+`). În acest caz, lexer-ul va afișa un mesaj corespunzător și va returna valoarea găsită (prin intermediul funcției `yytext()`, o metodă din clasa `Yylex`). De asemenea, a mai fost introdusă o definiție `not_star`, care pe baza expresiei regulate `[^*]` va găsi toate caracterele diferite de caracterul `*`. Pentru această definiție, a fost adăugată și o regulă lexicală care afișează un mesaj corespunzător la găsirea unui caracter care se potrivește cu expresia regulată.

Pentru a putea returna valoarea numerelor găsite, trebuie modificat corespunzător și constructorul clasei atomului `IntLit`. Astfel, au fost implementate:

- data membră `value` (de tip `private int`)
- constructorul `public IntLit(String textIn)` care inițializează data `value` cu valoarea primită ca parametru

- metoda `public String toString()` care permite afișarea datei membre `value` la afișarea unui atom
- o metodă `public int valueGet()` care permite citirea datei membre private `value`

În conformitate cu modificările aduse fișierului sursă `sum.jlex`, trebuie modificat și programul care instanțiază și utilizează lexer-ul generat pe baza specificațiilor. Astfel, se va scrie un nou fișier `Driver2.java` conform Codului sursă 5.5.

Codul sursă 5.5: `Driver2.java`

```
import java.io.*;

public class Driver2{
    public static void main(String[] args) throws Exception
    {
        File f = new File("in2.txt");
        FileInputStream fis = new FileInputStream (f);

        Ylex scanner = new Ylex(fis);
        int sum = 0;
        do{
            IntLit token = scanner.next_token();
            System.out.println(token);
            if (token == null) {
                break;
            }
            sum = sum + token.valueGet();
        } while(true);

        System.out.println("Suma este " + sum);
    }
}
```

Modificările aduse fișierului `Driver2.java` față de versiunea corespunzătoare parserului de spații implementează următoarele funcționalități:

- fiecare atom lexical recunoscut de analizor este scris într-un obiect de tip `IntLit` (`IntLit token = scanner.next_token()`)
- valoarea atomului este afișată (`System.out.println(token)`)
- se verifică dacă atomul lexical are valoarea `null`, caz în care se oprește execuția programului (cu instrucțiunea `break` se iese din bucla `while`)

- valorile atomilor lexicali sunt adunați la variabila `sum`, care este afișată pe ecran înainte de terminarea execuției programului

Exerciții:

1. Executați fișierul `sum.jlex` folosind JLex și compilați fișierul java rezultat.
2. Scrieți un fișier `in2.txt` care să conțină spații, valori numerice, caractere oarecare precum și caractere `*`.
3. Compilați și executați `Driver2.java`. Observați ce afișează programul în funcție de fișierul `in2.txt` introdus.
4. Ce se întâmplă dacă nu includem în `Driver2.java` porțiunea de cod cu ajutorul căreia se iese din buclă la întâlnirea unui atom nul?
5. Ce se întâmplă dacă se inversează ultimele două linii din fișierul `sum.jlex`. Argumentați răspunsul.

5.4.3 Specificarea de stări și tranziții în fișierul JLex

În fișierul sursă `.jlex` se pot specifica stări precum și tranzițiile dintre ele, ceea ce reflectă faptul că lexer-ul rulează practic un AFD. Vom modifica fișierul `sum.jlex` de mai sus în felul următor:

1. Se adaugă o nouă directivă `%state negative`, care va defini o nouă stare cu denumirea „negative“
2. Se adaugă o nouă regulă lexicală (antepenultima, înainte de regulile `{not_star}` și `.`) care va conduce în starea `negative` definită la întâlnirea șirului de caractere „NEG“. Tranziția într-o altă stare se realizează cu metoda `yybegin()` a clasei `Yylex`, ce primește ca parametru numele stării definite ca directivă `NEG { yybegin(negative); }`
3. Regulile lexicale se vor modifica astfel încât, odată intrat în starea `negative`, toate valorile numerice întâlnite vor fi considerate a fi negative, iar la întâlnirea unui nou șir de caractere „NEG“ se iese din starea `negative` și se intră în starea implicită `YYINITIAL`

În consecință, secțiunea a treia a noului fișier `.jlex` va arăta conform Codului sursă 5.6.

Codul sursă 5.6: Secțiunea a treia din noul fișier `sum.jlex`.

```

<negative>{DIGIT}+ { System.out.println("Am gasit o valoare
    numerica negativa!");
    yybegin(YYINITIAL);
    return new IntLit("-" + yytext()); }
{DIGIT}+ { System.out.println("Am gasit o valoare numerica
    pozitiva!");
    return new IntLit(yytext()); }
{WHITESPACE}+ { System.out.println("Am gasit un spatiu!");
}
<negative>"NEG" { yybegin(YYINITIAL); }
"NEG" { yybegin(negative); }
{not_star} { System.out.println("not star"); }
. { System.out.println("ERROR!"); }

```

Modificați fișierul `in2.txt` pentru a arăta ca în exemplul de mai jos:

```

NEGNEG
1
NEG
NEG3
42 NEG 20
a b c *

```

Exerciții:

1. Executați fișierul `sum.jlex` folosind `JLex` și compilați fișierul `java` rezultat.
2. Executați `Driver2.java`. Observați ce afișează în funcție de fișierul `in2.txt` introdus. Modificați fișierul `in2.txt` și analizați ce efecte au loc asupra mesajelor afișate de lexer.

5.5 Temă de laborator

Presupunem că, într-un cod Java, există o funcție denumită `myFunction` a cărei execuție consumă foarte mult timp, deci într-un program există limitări a numărului de apeluri datorită cerințelor de performanță.

Proiectați specificațiile `JLex` care generează un scanner ce numără apelurile acestei funcții dintr-un cod Java, urmând cerințele de mai jos:

- Scannerul trebuie să afișeze un mesaj de fiecare dată când întâlnește un apel la funcția `myFunction` însoțit de numărul liniei din cod (folosind `yyline`).

- După terminarea de parcurs a codului, se va afișa un mesaj cu numărul total de apeluri ale funcției `myFunction`.
- Codul sursă de intrare este scris în Java.
- În cod pot exista și alte variabile denumite `myFunction` sau alte funcții care au în denumirea lor `myFunction`.
- Cu toate acestea, definirea (declararea, antetul) funcției `myFunction` nu este prezentă în codul de intrare, doar apeluri la această funcție
- Nu există constrângeri cu privire la tipul returnat de `myFunction`, deci scannerul implementat va fi unul generic, acceptând orice tip returnat de această funcție (pentru a simplifica implementarea)