Programming Languages & Compilers

Phase1: Lexical Analyzer Generator

Mostafa Nabil Mohamed (54)
Tarek Mohamed Kamal (23)
Mohamed Mostafa Alnashar (46)
Merna Mohamed Mostafa (56)

Objective

This phase of the assignment aims to practice techniques for building automatic lexical analyzer generator tools.

Problem Statement

The lexical analyzer generator is required to automatically construct a lexical analyzer from a regular expression description of a set of tokens. The tool is required to construct a nondeterministic finite automata (NFA) for the given regular expressions, combine these NFAs together with a new starting state, convert the resulting NFA to a DFA, minimize it and emit the transition table for the reduced DFA together with a lexical analyzer program that simulates the resulting DFA machine.

The generated lexical analyzer has to read its input one character at a time, until it finds the longest prefix of the input, which matches one of the given regular expressions. It should create a symbol table and insert each identifier in the table. If more than one regular expression matches some longest prefix of the input, the lexical analyzer should break the tie in favor of the regular expression listed first in the regular specifications. If a match exists, the lexical analyzer should produce the token class and the attribute value. If none of the regular expressions matches any input prefix, an error recovery routine is Page ¼ to be called to print an error message and to continue looking for tokens.

The lexical analyzer generator is required to be tested using the given lexical rules of tokens of a small subset of Java. Use the given simple program to test the generated Lexical analyzer.

Keep in mind that the generated lexical analyzer will integrate with a generated parser which you should implement in phase 2 of the assignment such that the lexical analyzer is to be called by the parser to find the next token.

Lexical Rules

- The tokens of the given language are: identifiers, numbers, keywords, operators and punctuation symbols.
- The token id matches a letter followed by zero or more letters or digits.
- The token num matches an unsigned integer or a floating-point number. The
- consists of one or more decimal digits, an optional decimal point followed by one or
- more digits and an optional exponent consisting of an E followed by one or more
- keywords are reserved. The given keywords are: int, float, boolean, if, else, while.
- Operators are: +, -, *, /, =, <=, <, >, >=, !=, ==
- Punctuation symbols are parentheses, curly brackets, commas and semicolons.
- Blanks between tokens are optional.

Lexical Rules Input File Format

- Regular definitions are lines in the form LHS = RHS
- Regular expressions are lines in the form LHS: RHS
- Keywords are enclosed by { } in separate lines.
 Punctuations are enclosed by [] in separate lines
- \L represents the Lambda symbol.
- The following symbols are used in regular definitions and regular expressions with the meaning discussed in class: -1 + *()
- Any reserved symbol needed to be used within the language, is preceded by an
- escape backslash character.

Used Data Structure

Graph: It's graph contains the nodes and edges of the automaton. Each graph contains some attributes like:

- Node start → it holds the first node of the graph.
- Accept_state → it holds the last state of the graph that accepts definition or expression in that automaton the graph represents.

Node: It represents one of the states that an automaton consists of. Each node contains some attributes like:

- Int id → it holds the number of this node with respect to all automaton's nodes.
- String accepted_input → it holds the pattern that was accepted if it was an accept
 node and if not then the value would be non to represent that it's not accept node.
- Int priority → it holds the priority of the accepted node depending on the order of the lines entered in the input file.
- Vector<Edges> entering, leaving → representing the edges between nodes.

Edge: It represents one of the edges that an automaton has between states. Each edge contains some attributes like:

- Node src → it holds the source node.
- Node dest → it holds the destination node.
- RegularDef weight → it holds the weight of each edge.

RegularDef: It's mainly a graph. it holds the weight of each edge. It can be a one character weight or a whole graph representing a regular definition. Each regularDef contains some attributes like:

Graph regular_def → graph that holds the definition sequence.

NFA: It represents the Non Deterministic Automaton that will represent the rules of the language. Mainly it's a graph. Also, It's a singleton as there's only one NFA that represents all the rules and many classes will need to use it. Each NFA contains some attributes like:

Graph automaton → it holds the Non Deterministic Automaton graph.

DFA: It represents the Deterministic Finite Automaton that will represent the rules of the language after converting the NFA to DFA. Mainly it's a table that represents the transitions from one set of nodes to another set under some input. Also, It's a singleton as there's only one DFA that represents all the rules and many classes will need to use it. Each DFA contains some attributes like:

- vector<pair<Node*, map<RegularDef*,Node*>>> transition_table → it holds the
 transition table of the NFA after converting it. Each pair in the vector represents a
 row in the table. Each row has Node as the state and map that holds the inputs to
 that state as RegularDef and a Node that represents the state to reach under that
 input.
- vector<set<Node*>> stateMappingTable → it's a vector to keep track of the states
 we have while building the final transition table.

This was the main classes structure we have in this project that will be used as data structures in the whole project along with the others like: vector, map, set, and array.

Algorithms And Techniques

Overview of the program

The lexical analyzer generator starts by reading the input file which represents the rules of the language and after that it starts to parse it line by line and represent it as tokens after that we send the tokens of each line to be represented as an Non Deterministic Automaton by building each token as nodes in a graph then we merge the automatons of all lines together by ORing them and send the NFA to be converted into DFA and then minimize it. Finally we read the program that is needed to be tested character by character and test it through our minimized DFA and check if it works matches any of the patterns we created from the rules file and print each token's type after that.

Tokenizing the rules

After reading the input file each line is tokenized into a vector<string>.

We use **ReadRules** class to read the file line by line and then send it to **ExpNfa** class to function **construct_automata** which send it at the beginning to **Tokenizing** class to split the line after each space and check if it contains special cases that need further splits like having operation like * + - | = at the end or in between of a word and that is achieved by parsing each word character by character to check if there's any operation so we will split the word before it. Finally we return a vector of tokens back to the constuct_automata to continue building it.

Creating NFA

We used Thompson's construction algorithm to create the NFA by representing each char as automata then merging the automatons together as the pattern state.

After the line is tokenized we have a vector of tokens that we start working recursively by splitting each in its constituent sub expressions.

First, we check the type of the line if it's a definition, expression, keyword, or punctuation. then depending on the type we build the automata as follows:

Definition: we split the definition into the simple forms as if it's A-Z we create automata for each character and then merge them together to create the automata. After that we add it to a map<string, Definition*> to keep track of these definitions and to use them again if needed.

Expression: we create the automata of that expression using the definitions created or creating new ones and then store it in a vector of graphs to be ORed later to create the last NFA.

Keyword/Punctuation: create automata by oring the characters together.

We create the graphs of these automatons by splitting the tokens over *, +, |, and each character is created as a single automata. If it's 'a' then we have 0---a--->0 and then merge graphs together depending on the operation using the rules of thompson.

Finally after creating all the NFAs of all the lines we build the final NFA by adding another node at the beginning and edges to each NFA's starting node using LAMBDA.

Deriving DFA

After creating the NFA we send it to the **DFA** class to start building the transition table of the DFA after converting the NFA. We do that by using the subset construction algorithm.

For each state in the DFA transition table not visited yet, we get its destination set for each transition, get their epsilon closure and convert it to a new state in the DFA transition table. When no state for a certain transition is available, we direct the transition to a dummy state representing no accepted state.

Minimizing DFA

For each regular language, there also exists a minimal automaton that accepts it, that is, a DFA with a minimum number of states and this DFA is unique (except that states can be given different names). The minimal DFA ensures minimal computational cost for tasks such as pattern matching. There are two classes of states that can be removed or merged from the original DFA without affecting the language it accepts to minimize it.

- Unreachable states are the states that are not reachable from the initial state of the DFA, for any input string.
- Non Distinguishable states are those that cannot be distinguished from one another for any input string.

First, The main partition, which has all the states in the beginning, is divided into accepted and non-accepted partitions. Then, take each partition of them and check each two states with the same input definition to divide that partition into other partitions in condition that having the output states in the same partition.

Then check every time the previous number of partitions equals the next number of partitions if yes so we finish the minimization and begin to build the Minimized transition state table with the new partitions.

With conditions that:

- The starting group is the partition that has the starting state in it.
- The accepted input for each group is the accepted input for the first state in that partition.

Parse Program

We take the program and split into tokens after each space then parse it character by character and send it to the minimized DFA to check if it matches any pattern by following the states under the inputs given until it reaches accept state if it matches. If it does we will print the matched pattern. If not we will print not matched.

Sample Runs

1. Test 1

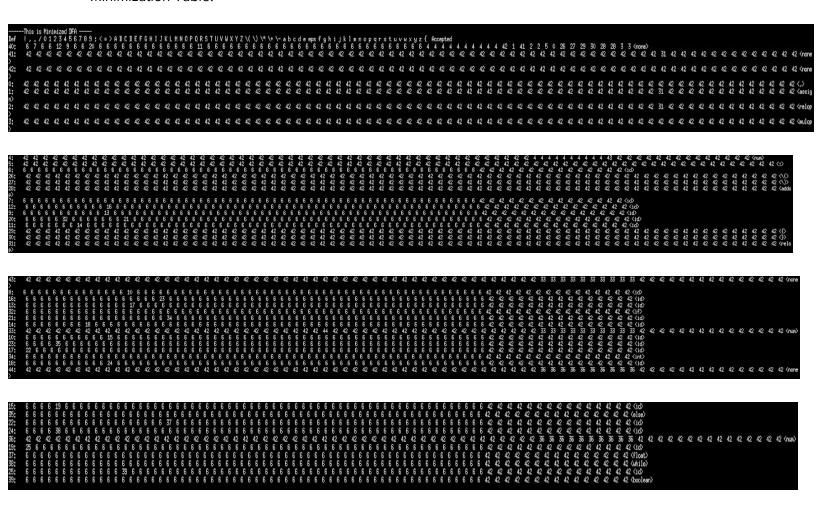
Rules:

```
letter = a-z | A-Z
digit = 0 - 9
id: letter (letter|digit)*
digits = digit+
{boolean int float}
num: digit+ | digit+ . digits ( \L | E digits)
relop: \=\= | !\= | > \= | < | <\=
assign: =
{ if else while }
[; , \( \) { }]
addop: \+ | -
mulop: \* | /|</pre>
```

Program:

```
int sum , count , pass ,
mnt; while (pass != 10)
{
pass = pass + 1;
}
```

Minimization Table:



Output:

```
int
id
id
id
id
while
id
relop
num
(
id
assign
id
addop
num
;
}
```

2. Test 2

Rules:

```
letter = a-z | A-Z
digit = 0 - 9
id: letter (letter|digit)*
digits: digit+
```

Program:

mostafa tarek 1245 nashar 1212 merna ada21 as2d1 2assad 5421

Minimized Table:

:



Output:

id
id
digits
id
digits
id
id
digits
id
digits
id
digits