

CS143 Compilers

Lecture 2
July 2, 2001

Recognition of Tokens

- Now that we know how to specify tokens, we can figure out how to recognize them.
- The basic idea is as follows:
 - First, we construct a regular expression for each token.
 - Then, we scan through the input, one character at a time, until the input is recognized as a string that matches one of the patterns described by the regular expressions.
 - Return the corresponding token, and repeat this process for the remaining input.

7/1/2001

2

Transition Diagrams

- In attempting to recognize a token, we could try to match the input to each regular expression in turn until we are successful.
- This is woefully inefficient, as this forces us to pass through the input several times.
- We can instead use a *transition diagram* to effectively combine the regular expressions into a single regular expression that can recognize any valid token.
- A well-designed transition diagram minimizes the need to backtrack through the input.

7/1/2001

3

Anatomy of Transition Diagrams

- A transition diagram consists of a set of *states* and a set of *transitions*.
- Each diagram contains a single state that is called the *start state*. There can be no transitions to the start state.
- Each diagram contains at least one state that is called a *final or accepting state*. There can be no transitions from an accepting state.

7/1/2001

4

Using Transition Diagrams

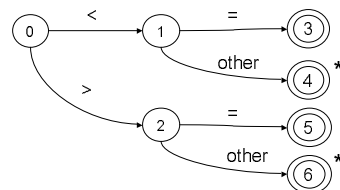
- Transitions between states are guided by input symbols from a given alphabet.
- Initially, the scanner is in the start state. It then reads the first input symbol, and determines which transition to use.
- Input is read and transitions followed until an accepting state is reached. This signifies that a token has been recognized, and the string of input symbols read since the start state is the corresponding lexeme.
- Control returns to the start state to read the next token.

7/1/2001

5

The following transition diagram can be used to recognize the relational operators <, <=, > and >=.

State 0 is the start state, and double circles denotes accepting states. The * next to a state indicates that the scanner should move back by one input symbol.



7/1/2001

6

Implementing a Transition Diagram

- Once a transition diagram has been constructed, it is easy to write a program to implement it.
- Each state corresponds to a segment of code.
- If there are edges leaving the state, then this code must read the next input character and cause a branch to the code for the appropriate state.
- If a state is a final state, a token has been recognized.
- If a state is not a final state, and no transition to another state is available, then a lexical error has been detected.

7/1/2001

7

Ad-hoc Scanners

- A lexical analyzer implemented from a transition diagram is called an *ad-hoc scanner*. This terminology stems from the manner in which the transition diagram is built from patterns for each token.
- It would be desirable to have a systematic method of creating this transition diagram, given a set of patterns that specify the lexemes for the tokens of the source language.
- Such a method would allow lexical analyzers to be created automatically.

7/1/2001

8

Recognizers

- A *recognizer* for a language is a program that can determine whether a given input string belongs to the language.
- Our goal is to be able to systematically build a recognizer from any given regular expression.
- We will then show how such a recognizer can be used to create a lexical analyzer for a source language, provided that each token can be described using a regular expression.
- We will see that given regular expressions for each token, the recognizer can be generated automatically.

7/1/2001

9

Finite Automata

- To convert a regular expression into a recognizer, we can construct a generalized transition diagram called a *finite automaton*.
- A *deterministic finite automaton* (DFA) is a transition diagram in which only one transition may be made from a state on a given input symbol.
- A *nondeterministic finite automaton* (NFA) allows more than one transition from a state on any given input symbol.
- Generally, using DFA's yields faster recognizers, but they are more difficult to design and are more complex.

7/1/2001

10

Nondeterministic Finite Automata

- A nondeterministic finite automaton is a mathematical model consisting of:
 - A set of *states* S
 - A set of input symbols Σ
 - A transition function *move* that maps state-symbol pairs to sets of states
 - A state s_0 that is the *start* or *initial state*
 - A set of states F called *accepting* or *final states*

7/1/2001

11

Transition Graphs

- An NFA can easily be represented by its *transition graph*. The transition graph is a generalized transition diagram that is built from the NFA as follows:
 - The states in the transition graph correspond to the states of the NFA, including the designation of the start state and accepting states.
 - The transitions in the transition graph represent the NFA's transition function, *move*. If *move* maps a state s_1 and input symbol a to state s_2 , then the transition graph includes a transition from state s_1 to state s_2 on symbol a .

7/1/2001

12

Acceptance of Strings

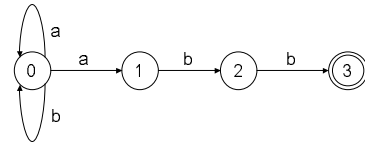
- An NFA *accepts* a string x if there exists a path in the NFA's transition graph from the start state to an accepting state, such that the concatenation of the input symbols along this path yields the string x .
- The *language defined by an NFA* is the set of input strings that it accepts.
- If the NFA is built from a regular expression, then the language that it defines is the language denoted by that regular expression.

7/1/2001

13

Example

This NFA recognizes the language denoted by the regular expression $(a|b)^*abb$.



7/1/2001

14

Deterministic Finite Automata

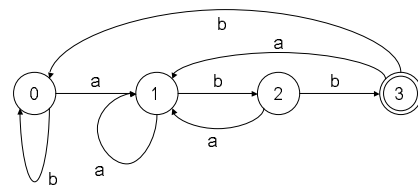
- A deterministic finite automaton is an NFA with the following restrictions:
 - No state has a transition on the empty string ϵ (known as an ϵ -transition)
 - For each state s and input symbol a , there is at most one transition from s on a .
- As a result of these restrictions, it is very easy to determine whether a DFA accepts a given string, i.e. whether a string matches the regular expression corresponding to the DFA.

7/1/2001

15

Example

This DFA also recognizes the language denoted by the regular expression $(a|b)^*abb$.



7/1/2001

16

From a Regular Expression to an NFA

Our primary goal, building an efficient recognizer for a regular expression, can be accomplished by building an NFA that accepts exactly those strings in the language denoted by the regular expression, and then using one of two approaches:

- converting the NFA to an equivalent DFA, which can easily be simulated
- simulating the NFA directly.

We will examine the trade-offs involved in both approaches.

7/1/2001

17

Construction of an NFA from a Regular Expression

An NFA can be built from a regular expression over an alphabet Σ by using the expression's syntactic structure as a guide for the construction:

- First, we show how to build an NFA that accepts the empty string ϵ , and each of the input symbols in Σ .
- We then use induction to handle all other regular expressions over Σ . Given NFA built from two regular expressions r and s , we show how to build NFA representing $r|s$, rs , and r^* .

7/1/2001

18

Thompson's Construction

- Let r and s be regular expressions over Σ . For any regular expression r over Σ , let $M(r)$ be an NFA that accepts the strings in $L(r)$.
 - If r is the empty string ϵ , then $M(r)$ consists of an initial state i , a final state f , and a single transition from i to f on ϵ .
 - If r is a symbol a from Σ , then $M(r)$ contains the states i and f , and a transition from i to f on a .
 - $M(r|s)$ contains the states i and f , along with ϵ -transitions from i to the initial states of $M(r)$ and $M(s)$, and ϵ -transitions from the final states of $M(r)$ and $M(s)$ to f .

7/1/2001

19

Thompson's Construction, cont'd

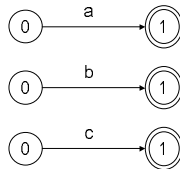
- $M(rs)$ is built as follows: The initial state of $M(rs)$ is the initial state of $M(r)$, and the final state is the final state of $M(s)$. The final state of $M(r)$ is merged with the initial state of $M(s)$ to form a new state, whose transitions are the union of the transitions to the final state of $M(r)$ and the transitions from the initial state of $M(s)$.
- $M(r^*)$ is built using $M(r)$, an initial state i , and a final state f . Four ϵ -transitions are added: from i to f from the final state of $M(r)$ to its initial state, from i to the initial state of $M(r)$, and from the final state of $M(r)$ to f .

7/1/2001

20

Example

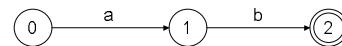
We now show how an NFA for the regular expression $(ab|b)^*c$ can be built using Thompson's construction. We begin with the following basic NFA's that recognize each of our alphabet symbols, a , b , and c :



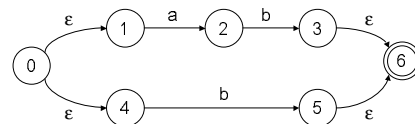
7/1/2001

21

Using these basic NFA's, we build an NFA for the regular expression ab .



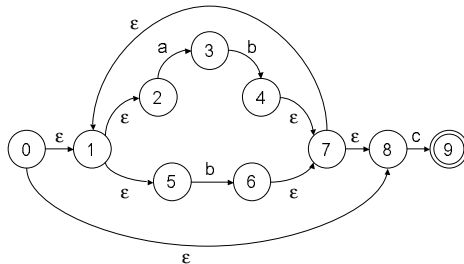
Using the above NFA, we then build an NFA for $(ab|b)$.



7/1/2001

22

We complete the construction by implementing the Kleene closure operator $*$, and then concatenating c .



7/1/2001

23

Conversion of an NFA into a DFA

- While NFA's are easier to design for a given language than DFA's, they require more effort to simulate.
- Using an algorithm called *subset construction*, a given NFA can be converted into an equivalent DFA.
- Each state in the DFA represents the set of states in the NFA that can be reached using a given input string.
- The number of states in the DFA can be as large as the number of subsets of the set of the NFA's states, which is 2^n , where n is the number of states in the NFA.

7/1/2001

24

Constructing a DFA from an NFA

- To convert an NFA N to an equivalent DFA D , we proceed by constructing a transition function $Dtran$ for D so that D will simulate all possible moves N can make on a given input string.
- Each state in D represents a set of states from N . To keep track of sets of NFA states, we use the ϵ -closure operation, defined as follows:
 - Given a set of states T from N , ϵ -closure(T) is the set of states of N that can be reached from some state in T on ϵ -transitions.

7/1/2001

25

E-closure

ϵ -closure can be implemented as follows:

```

push all states in  $T$  onto  $stack$ 
 $\epsilon$ -closure( $T$ ) = { $T$ }
while  $stack$  is not empty
  pop  $t$ , the top element, off of  $stack$ 
  for each state  $u$  such that  $move(t, \epsilon) = u$ 
    if  $u$  is not in  $\epsilon$ -closure( $T$ )
      add  $u$  to  $\epsilon$ -closure( $T$ )
      push  $u$  onto  $stack$ 
end
end
    
```

7/1/2001

26

Subset Construction

ϵ -closure(s_0) is the only state in $Dstates$, and is unmarked

while there is an unmarked state T in $Dstates$

 mark T

 for each input symbol a

$U = \epsilon$ -closure($move(T, a)$)

 if U is not in $Dstates$

 add U as an unmarked state to $Dstates$

$Dtran[T, a] = U$

 end

end

7/1/2001

27

Notes on Subset Construction

In the preceding pseudo-code:

- T is a set of states from the NFA N . It corresponds to a state in the newly constructed DFA D for which outgoing transitions have not yet been assigned.
- Given the input symbol a , U is the set of all states from N that can be reached from some state in T .
- ϵ -closure(T) is a depth-first search through the graph whose vertices are the states of N and whose edges are the ϵ -transitions.

7/1/2001

28

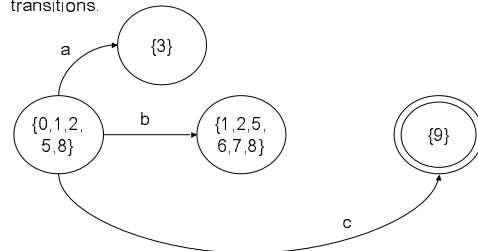
First, we create the initial state of our DFA, by computing the ϵ -closure of state 0.



7/1/2001

29

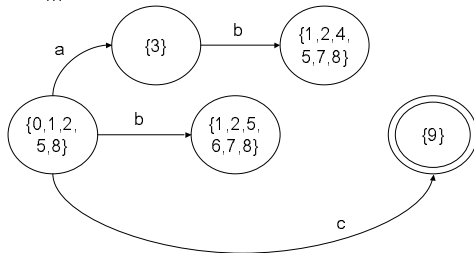
Then, we add the ϵ -closure of any states that can be reached from ϵ -closure(0), as well as the appropriate transitions.



7/1/2001

30

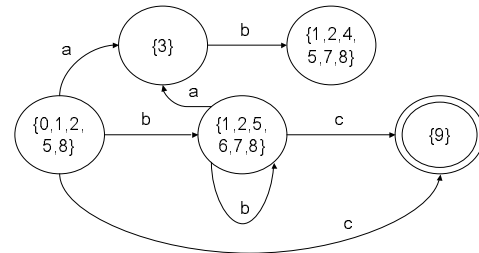
Next, we mark {3} and compute its going transitions, causing us to add the ϵ -closure of state 4 as a new unmarked state.



7/1/2001

31

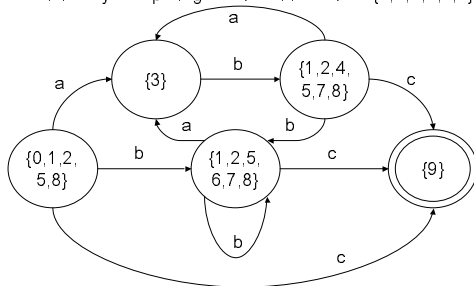
Next, we mark {1,2,5,6,7,8} and compute its outgoing transitions.



7/1/2001

32

We finish by computing the transitions from {1,2,4,5,7,8}.



7/1/2001

33

Simulation of an NFA

To determine whether an NFA N , with initial state s_0 and final states F , accepts an input string x , we proceed as follows:

```

 $S = \epsilon\text{-closure}(s_0)$ 
 $a = \text{first character of } x$ 
while characters remain in  $x$ 
     $S = \epsilon\text{-closure}(\text{move}(S, a))$ 
     $a = \text{next character in } x$ 
end
if  $S$  and  $F$  are not disjoint
     $N$  accepts  $x$ 

```

7/1/2001

34

Two-Stack Simulation of an NFA

- The preceding algorithm can be implemented efficiently using two stacks and a bit vector indexed by the states of the NFA.
- The bit vector is used to remember whether a state has been added to a stack, so we do not add it twice.
- One stack represents the set S , and the other stores $\epsilon\text{-closure}(S, a)$ while it is being computed. $\epsilon\text{-closure}(S, a)$ is completed, the roles of the stacks are interchanged.

7/1/2001

35

Efficiency Analysis

- Each state has at most two outgoing transitions, and the stack can hold at most $|N|$ elements, where $|N|$ is the number of states in N . Therefore, computing the $\epsilon\text{-closure}(S, a)$ requires time proportional to $|N|$.
- The $\epsilon\text{-closure}$ is computed $|x|$ times, where $|x|$ is the length of x .
- Therefore, this implementation runs in time proportional to $|N||x|$.

7/1/2001

36

Simulating a DFA

To determine whether a DFA D accepts a string x , we proceed as follows.

```

 $s$  = the start state of  $D$ 
 $c$  = first character in  $x$ 
while characters remain in  $x$ 
     $s = \text{move}(s, c)$ 
     $c$  = next character in  $x$ 
end
if  $s$  is an accepting state of  $D$ ,
     $x$  is accepted
    
```

7/1/2001

37

Time-space Tradeoffs

- Let r be a regular expression, and x be an input string. Let $|r|$ be the length of r , and $|x|$ be the length of x .
- To determine whether x belongs to $L(r)$ using an NFA, $O(|r|)$ space and $O(|r||x|)$ time is required.
- To perform the same task using a DFA, $O(2^{|r|})$ space is required, but only $O(|x|)$ time is needed.
- For this reason, when a DFA is used, transitions are usually computed only when they are needed for the particular input string. This approach, called "lazy transition evaluation," requires only $O(|r|)$ space, the same as for an NFA.

7/1/2001

38

Design of a Lexical Analyzer Generator

- We are now ready to show how a lexical analyzer can be generated for a given source language.
- We assume that the tokens for the source language can be recognized using a given set of patterns, described by regular expressions.
- We will create a recognizer that will try to match some prefix of the text of a source program against the patterns. If more than one pattern matches a prefix, the recognizer must choose the longest lexeme matched. If two or more patterns match the longest lexeme, the first pattern listed is chosen. The recognizer then continues with the remainder of the input.

7/1/2001

39

Pattern Matching based on NFA's

A lexical analyzer recognizing tokens using patterns p_1, p_2, \dots, p_n may be generated using an NFA as follows:

- For each pattern p_i construct an NFA $M(p_i)$.
- Construct an NFA N from the regular expression $p_1 | p_2 | \dots | p_n$. Add a start state s_0 to N , and ϵ -transitions from s_0 to the initial states of each $M(p_i)$.
- Simulate N , but instead of accepting a lexeme as soon as an accepting state is reached, continue until there are no transitions on the current input symbol.
- Retract the input pointer to the most recent position at which an accepting state was reached, and recognize the corresponding token.

7/1/2001

40

DFA for Lexical Analyzers

- A similar approach may be used with a DFA, built from the NFA described in the previous algorithm.
- As with an NFA, the DFA must be simulated until we reach a state from which there is no transition on the current symbol.
- Upon termination, we backtrack through the input until we reach a position at which the DFA entered an accepting state, and recognize the corresponding token.

7/1/2001

41

The Lookahead Operator

- Often, a lexical analyzer must look beyond the end of a lexeme in order to ensure that the proper token is recognized.
- For this purpose, the lookahead operator $/$ can be used to describe the pattern for such a token.
- The pattern r/s is interpreted as follows: a string in $L(r)$ is a lexeme for the given token only if it is followed by a string in $L(s)$.

7/1/2001

42

Implementing the Lookahead Operator

- When converting a pattern that includes the lookahead operator into an NFA, we treat the / as if it were ϵ , so as not to look for / in the input.
- If we recognize a token corresponding to a pattern that includes /, the end of the lexeme is not the most recent position at which the NFA was in an accepting state.
- In this case, we must find the most recent position in the input at which the NFA had a transition on /.

7/1/2001

43

Optimization of DFA-based Pattern Matchers

- Various algorithms have been developed to optimize pattern matchers using DFA. Some techniques are:
 - Creating a DFA directly from a regular expression, without building an NFA.
 - Reducing the number of states of a DFA
 - Building a more compact representation of the transition table of a DFA to save space.
- For more information about these methods, consult the text by Aho, Sethi and Ullman (the "Dragon book")

7/1/2001

44