CMPT 379 Compilers

#### Anoop Sarkar http://www.cs.sfu.ca/~anoop

1

# Lexical Analysis

• Also called *scanning*, take input program *string* and convert into tokens



# Token Attributes

1

- Some tokens have attributes
  - T\_IDENT "sqrt"
  - T\_INTCONSTANT
- Other tokens do not
  - T\_WHILE
- *Token*=T\_IDENT, *Lexeme*="sqrt", *Pattern*
- Source code location for error reports

# Lexical errors

- What if user omits the space in "doublef"?
  - No lexical error, single token
    T\_IDENT("doublef") is produced instead of sequence T\_DOUBLE, T\_IDENT("f")!
- Typically few lexical error types
  - E.g., illegal chars, opened string constants or comments that are not closed

# Implementing Lexers: Loop and switch scanners

- Ad hoc scanners
- Big nested switch/case statements
- Lots of getc()/ungetc() calls
  - Buffering
- Can be error-prone, use only if
  - Your language's lexical structure is simple
  - Tools don't do what you want
- Changing or adding a keyword is problematic
- Key idea: separate the defn from the implementation
- Problem: we need to reason about patterns and how they can be used to define tokens (recognize strings).

# Formal Languages: Recap

- Symbols: a, b, c
- Alphabet : finite set of symbols  $\Sigma = \{a, b\}$
- String: sequence of symbols bab
- Empty string:  $\varepsilon$  Define:  $\Sigma^{\varepsilon} = \Sigma \cup \{\varepsilon\}$
- Set of all strings:  $\Sigma^*$  cf. *The Library of Babel*, Jorge Luis Borges
- (Formal) Language: a set of strings
   { a<sup>n</sup> b<sup>n</sup> : n > 0 }

# Regular Languages

- The set of regular languages: each element is a regular language
- Each regular language is an example of a (formal) language, i.e. a set of strings
  e.g. { a<sup>m</sup> b<sup>n</sup> : m, n are +ve integers }

# Regular Languages

- Defining the set of all regular languages:
  - The empty set and  $\{a\}$  for all a in  $\Sigma^{\epsilon}$  are regular languages
  - If  $L_1$  and  $L_2$  and L are regular languages, then:  $L_1 \cdot L_2 = \{xy \mid x \in L_1 \text{ and } y \in L_2\}$  (concatenation)  $L_1 \cup L_2$  (union)  $L^* = \bigcup_{i=0}^{\infty} L^i$  (Kleene closure) are also regular languages
  - There are no other regular languages

# Formal Grammars

- A formal grammar is a concise description of a formal language
- A formal grammar uses a specialized syntax
- For example, a regular expression is a concise description of a regular language (*alb*)\**abb* : is the set of all strings over the alphabet {*a*, *b*} which end in *abb*

# Regular Expressions: Definition

- Every symbol of  $\Sigma \cup \{ \epsilon \}$  is a regular expression
- If  $r_1$  and  $r_2$  are regular expressions, so are
  - Concatenation:  $r_1 r_2$
  - Alternation:  $r_1 lr_2$
  - Repetition:  $r_1^*$
- Nothing else is.
  - Grouping re's: e.g. aalbc vs. ((aa)lb)c

# Regular Expressions: Examples

- Alphabet { 0, 1 }
- All strings that represent binary numbers divisible by 4 (but accept 0) ((0|1)\*00)|0
- All strings that do not contain "01" as a substring 1\*0\*

# **Regular Expressions**

- To describe all lexemes that form a token as a *pattern* 
  - (0|1|2|3|4|5|6|7|8|9)+
- Need decision procedure: to which token does a given sequence of characters belong (if any)?
  - Finite State Automata

# Finite Automata: Recap

• A set of states S

– One start state  $q_0$ , zero or more final states F

- An alphabet  $\sum$  of input symbols
- A transition function:

 $-\delta: S \ge \Sigma \Longrightarrow S$ 

• Example:  $\delta(1, a) = 2$ 

#### Finite Automata: Example

• What regular expression does this automaton accept?



Answer: (0|1)\*00

# FA: Pascal Example



### Building a Lexical Analyzer

- Token  $\Rightarrow$  Pattern
- Pattern  $\Rightarrow$  Regular Expression
- Regular Expression  $\Rightarrow$  NFA
- NFA  $\Rightarrow$  DFA
- DFA  $\Rightarrow$  Lexical Analyzer

# NFAs

- NFA: like a DFA, except
  - A transition can lead to more than one state, that is,  $\delta: S \times \Sigma \Rightarrow 2^S$
  - One state is chosen non-deterministically
  - Transitions can be labeled with ε, meaning states can be reached without reading any input, that is,

 $\delta: S \ge \Sigma \cup \{ \varepsilon \} \Longrightarrow 2^{S}$ 

# Thompson's construction

- Converts regexps to NFA
- Five simple rules
  - Symbols
  - Empty String
  - Alternation  $(r_1 \text{ or } r_2)$
  - Concatenation ( $r_1$  followed by  $r_2$ )
  - Repetition  $(r_1^*)$

• For each symbol *x* of the alphabet, there is a NFA that accepts it (include a *sinkhole* state)



• There is an NFA that accepts only ε



• Given two NFAs for  $r_1$ ,  $r_2$ , there is a NFA that accepts  $r_1 | r_2$ 



• Given two NFAs for  $r_1$ ,  $r_2$ , there is a NFA that accepts  $r_1r_2$ 



• Given a NFA for  $r_1$ , there is an NFA that accepts  $r_1^*$ 



# Example

- Set of all binary strings that are divisible by four (include 0 in this set)
- Defined by the regexp: ((0|1)\*00) | 0
- Apply Thompson's Rules to create an NFA

#### Basic Blocks 0 and 1



(this version does not report errors: no *sinkholes*)





 $(0|1)^{*}$ 



(0|1)\*00



# Simulating NFAs

- Similar to DFA simulation
- But have to deal with ε transitions and multiple transitions on the same input
- Instead of one state, we have to consider *sets* of states
- Simulating NFAs is a problem that is closely linked to converting a given NFA to a DFA

# NFA to DFA Conversion

- Subset construction
- Idea: subsets of set of all NFA states are *equivalent* and become one DFA state
- Algorithm simulates movement through NFA
- Key problem: how to treat ε-transitions?

#### ε-Closure

- Start state: q<sub>0</sub>
- $\varepsilon$ -closure(S): S is a set of states initialize:  $S \leftarrow \{q_0\}$   $T \leftarrow S$ repeat  $T' \leftarrow T$   $T \leftarrow T' \cup [\cup_{s \in T'} move(s, \epsilon)]$ until T = T'

# ε-Closure (T: set of states)

push all states in T onto *stack* initialize  $\varepsilon$ -*closure*(T) to T **while** *stack* is not empty **do begin** pop t off *stack* **for** each state u with  $u \in move(t, \varepsilon)$  **do if**  $u \notin \varepsilon$ -*closure*(T) **do begin** add u to  $\varepsilon$ -*closure*(T) push u onto *stack* **end** 

#### end

# NFA Simulation

- After computing the ε-*closure* move, we get a set of states
- On some input extend all these states to get a new set of states

 $\mathbf{DFAedge}(T, c) = \epsilon$ -closure  $(\cup_{q \in T} \mathbf{move}(q, c))$ 

# NFA Simulation

- Start state: q<sub>0</sub>
- Input:  $c_1, ..., c_k$

 $T \leftarrow \epsilon$ -closure( $\{q_0\}$ )

for  $i \leftarrow 1$  to k

 $T \leftarrow \mathbf{DFAedge}(T, c_i)$ 

# Conversion from NFA to DFA

- Conversion method closely follows the NFA simulation algorithm
- Instead of simulating, we can collect those NFA states that behave identically on the same input
- Group this set of states to form one state in the DFA
## Subset Construction

```
add \varepsilon-closure(q<sub>0</sub>) to Dstates unmarked

while \exists unmarked T \in Dstates do begin

mark T;

for each symbol c do begin

U := \varepsilon-closure(move(T, c));

if U \notin Dstates then

add U to Dstates unmarked

Dtrans[d, c] := U;

end

end
```

## Subset Construction

```
states[0] = \epsilon-closure({q<sub>0</sub>})
      p = j = 0
      while j \le p do begin
               for each symbol c do begin
                        e = DFAedge(states[j], c)
                        if e = states[i] for some i \le p
                        then Dtrans[j, c] = i
                        else p = p+1
                                 states[p] = e
                                 Dtrans[j, c] = p
               j = j + 1
               end
<sub>9/25/06</sub> end
```

#### Example: subset construction







### $move(\varepsilon$ -*closure*( $q_0$ ), 0)

()**E** ()

#### $\epsilon$ -*closure*(move( $\epsilon$ -*closure*( $q_0$ ), 0))



### move( $\epsilon$ -*closure*( $q_0$ ), 1)



#### $\epsilon$ -*closure*(move( $\epsilon$ -*closure*( $q_0$ ), 1))











- Algorithm for minimizing the number of states in a DFA
- Step 1: partition states into 2 groups: accepting and non-accepting



- Step 2: in each group, find a sub-group of states having property P
- P: The states have transitions on each symbol (in the alphabet) to the *same* group



9/25/06

- Step 3: if a sub-group does not obey P split up the group into a separate group
- Go back to step 2. If no further sub-groups emerge then continue to step 4



9/25/06

- Step 4: each group becomes a state in the minimized DFA
- Transitions to individual states are mapped to a single state representing the group of states



- Subset construction converts NFA to DFA
- Complexity:
  - in programs we measure time complexity in number of steps
  - For FSAs, we measure complexity in terms of the number of states

- Problem: An *n* state NFA can sometimes become a  $2^n$  state DFA, an exponential increase in complexity
  - Try the subset construction on NFA built for the regexp A\*aA<sup>n-1</sup> where A is the regexp (alb)
- Minimization can reduce the number of states
- But minimization requires determinization





9/25/06





 $2^5 = 32$  states 57

9/25/06

## NFA vs. DFA in the wild

<b>Engine Type</b>	Programs
DFA	<i>awk</i> (most versions), <i>egrep</i> (most versions), <i>flex</i> , <i>lex</i> , MySQL, Procmail
Traditional NFA	GNU <i>Emacs</i> , Java, <i>grep</i> (most versions), <i>less</i> , <i>more</i> , .NET languages, PCRE library, Perl, PHP (pcre routines), Python, Ruby, <i>sed</i> (most versions), vi
POSIX NFA	<i>mawk</i> , MKS utilities, GNU <i>Emacs</i> (when requested)
Hybrid NFA/DFA	GNU awk, GNU grep/egrep, Tcl



# Regexp to DFA: followpos

- *followpos(p)* tells us which positions can follow a position *p*
- There are two rules that use the *firstpos* {} and *lastpos* () information



9/25/06

*followpos*(j)+=k,l

$$\{k,l\} * (i,j)$$

. . .

*followpos*(i)+=k,l *followpos*(j)+=k,l





### Regexp to DFA: ((ab) | (ba))\*#

<i>root</i> ={1,3,5}	{1,3,5} A	<i>A: fp</i> (5),# {},# E,#
fp(1)=2 fn(3)=4	<i>A: fp</i> (1),a {2},a B,a	<i>B: fp</i> (2),b {1,3,5},b A,b
fp(3)=4 fp(2)=1,3,5	<i>A: fp</i> (3),b {4},b C,b	<i>C: fp</i> (4),a {1,3,5},a A,a
<i>fp</i> (4)=1,3,5		
1:a		a B
2:b		E) (A b
4:a	b	b b
5:#	a (C)	a C
9/25/06	<u> </u>	63

9/25/06

- (R|S)T = RT|ST
- $R^{**} == R^*$
- $R^{*}R^{*} == (R^{*})^{*} == R^{*}$  $== RR^{*} \epsilon$
- (R|S) == (S|R)
- (RS)T == R(ST)
- (R|S)|T == R|(S|T) == R(S|T) == RS | RTRISIT

- $R = R | R = R \epsilon$
- (RS)\*R == R(SR)\*
- $RR^* == R^*R$
- $(R|S)^* == (R^*S^*)^* ==$ (R\*S)\*R\* == $(R^*|S^*)^*$

## Equivalence of Regexps

- (01)\*
- $(01)(01)*|\epsilon$
- $(01)(01)*|(01)(01)*|\epsilon$   $R^* == RR^*|\epsilon$
- (01)(01)\*|(01)\*|
- 0(10)\*1|(01)\*

- $R^* == RR^* | \epsilon$
- R == R R
- RS == (RS)
- (RS)\*R == R(SR)\*

## Equivalence of Regexps

## Lexical Analyzer using DFAs

- Each token is defined using a regexp  $r_i$
- Merge all regexps into one big regexp  $-R = (r_1 | r_2 | ... | r_n)$
- Convert *R* to an NFA, then DFA, then minimize
  - remember orig NFA final states with each DFA state

## Lexical Analyzer using DFAs

• The DFA recognizer has to find the *longest match* for a token

- e.g. <print> and not <pr>>, <int>

- If two patterns match the same token, pick the one that was listed earlier in R
  - e.g. prefer final state (in the original NFA) of  $r_2$ over  $r_3$

# Lexical Analyzer using DFAs

- Alternative method:
  - Organize all the DFAs for each token in an ordered list
  - For input  $i_1, i_2, ..., i_n$  run all DFAs until some reach a final state (pick the longest match for each DFA)
  - Pick the token for which some DFA could read the longest match in the input,
    - e.g. prefer DFA #8 over all others because it read the input until  $i_{30}$  and none of the other DFAs reached  $i_{30}$
  - If two DFAs reach the same input character then pick the one that is listed first in the ordered list

- 2D array storing the transition table
- Adjacency list, more space efficient but slower
- Merge two ideas: array structures used for sparse tables like DFA transition tables
  - base & next arrays: Tarjan and Yao, 1979
  - Dragon book (default+base & next+check)



	a	b	c	d
0	-	1	-	2
1	1	-	1	-
2	1	2	1	-

	a	b	C	d
0	_	1	-	2
1	1	-	1	-
2	1	2	1	-





nextstate(s, x) :
L := base[s] + x
return next[L] if check[L] eq s









default

nextstate(s, x) :L := base[s] + x return next[L] if check[L] eq s else return nextstate(default[s], x)

9/25/06
## Summary

- Token  $\Rightarrow$  Pattern
- Pattern  $\Rightarrow$  Regular Expression
- Regular Expression  $\Rightarrow$  NFA
  - Thompson's Rules
- NFA  $\Rightarrow$  DFA
  - Subset construction
- DFA  $\Rightarrow$  minimal DFA
  - Minimization

## ⇒ Lexical Analyzer (multiple patterns)