

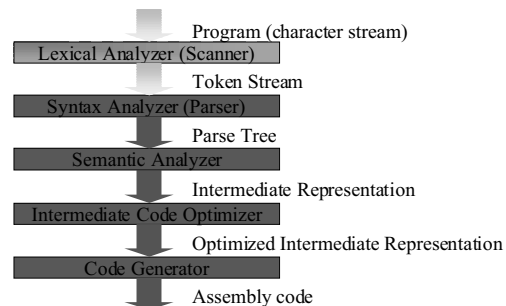
6.035

Fall 2000

Lecture 2: Lexical Analysis

Regular Expressions, NFA's, DFA's and Scanner Generation

Lexical Analysis



5

What is a Lexical Analyzer?

Source program text → Tokens

• Example of Tokens

- Operators = + - > ({ := == <
- Keywords if while for int double
- Numeric literals 43 6.035 -3.6e10 0x13F3A
- Character literals 'a' '~' '\"'
- String literals "6.891" "Fall 98" "\\\" = empty"

• Example of non-tokens

- White space space(' ') tab('\t') end-of-line('\n')
- Comments /*this is not a token*/

Lexical Analyzer needs to....

f o r v a r l = 1 0 v a r l < =

- for ID("var1") eq_op Num(10) ID("var1") leq_op
- Precisely separate the text stream into the correct stream of tokens
 - ID("var1") not ID("var") Num(1)
 - ID("var1") leq_op not ID("var1<=")
- Identify the type of token that matches the input string
 - 10 Num(10)
 - var1 ID("var1")

Lexical Analyzer needs to....

f o r v a r l = 1 0 v a r l < =

for ID("var1") eq_op Num(10) ID("var1") leq_op

- Describe different types of tokens in different languages
 - FORTRAN DO I=1,10
 - C++ for(int i=1; I<= 10; I++)
 - C-shell foreach i (1 2 3 4 5 6 7 8 9 10)

Lexical Analyzer needs to....

- Use *regular expressions* to precisely describe what strings each type of token can recognize

Examples of Regular Expressions

Regular Expression	Strings matched
<code>a</code>	<code>"a"</code>
<code>a · b</code>	<code>"ab"</code>
<code>a b</code>	<code>"a"</code> <code>"b"</code>
<code>—</code>	<code>""</code>
<code>a*</code>	<code>"a"</code> <code>"aa"</code> <code>"aaa"</code>

Examples of Regular Expressions

Regular Expression	Strings matched
<code>num = 0 1 2 3 4 5 6 7 8 9</code>	<code>"0"</code> <code>"1"</code> <code>"2"</code> <code>"3"</code> ...
<code>posint = num · num*</code>	<code>"8"</code> <code>"6035"</code> ...
<code>int = (_ -) · posint</code>	<code>"-42"</code> <code>"1024"</code> ...
<code>real = int · (_ (. · posint))</code>	<code>"-12.56"</code> <code>"12"</code> <code>"1.414"</code> ...

14

QUESTION

- Why `posint = num · num*` ?
Why not `posint = num*` ?

Few Additional Notations

- "one or more occurrence of" `r+ = r · r*`
- "zero or one occurrence of" `r? = r | _`

18

QUESTION

- What regular expression best identifies MIT course numbers?

`num = 0|1|2|3|4|5|6|7|8|9`

- 1) `class = num · num*`
- 2) `class = num · . · num*`
- 3) `class = num | . | num*`
- 4) `class = (num · . · num)*`

QUESTION

- What regular expression will match all numbers 0 to 256 and nothing else?

Definition: Formal Languages

- Alphabet Σ = finite set of symbols
 - $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$
- String s = finite sequence of symbols from the alphabet
 - $s = 6004$
- Empty string ϵ = special string of length zero
- Language L = set of strings over an alphabet
 - $L = \{6001, 6002, 6003, 6004, 6035, 6891, \dots\}$

Definition: Regular Expressions

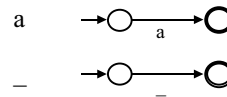
- For a regular expression r , the language $L(r) = \{\text{all the strings that match } r\}$
 - $L((a | \epsilon) \cdot b) = \{\text{"ab"}, \text{"b"}\}$
- Suppose r and s are regular expression denoting languages $L(r)$ and $L(s)$
 - $L(r | s) = L(r) \cup L(s)$
 - $L(r \cdot s) = \{xy \mid x \in L(r) \text{ and } y \in L(s)\}$
 - $L(r^*) = \{x_1 x_2 \dots x_k \mid x_i \in L(r) \text{ and } k \geq 0\}$
 - $L(\epsilon) = \{\epsilon\}$

More Regular Expressions

- We know:
 - $L(r | s)$ is the **union** of $L(r)$ and $L(s)$
 - $L(r \cdot s)$ is the **concatenation** of $L(r)$ and $L(s)$
 - $L(r^*)$ is the **Kleene closure** of $L(r)$
 - "zero or more occurrence of"

22

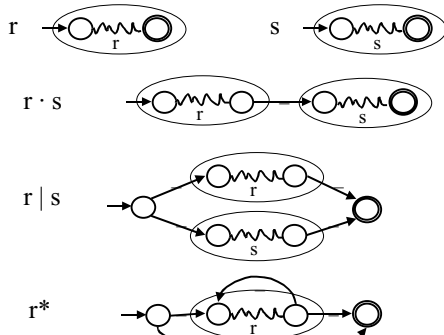
Constructing a NFA for a regular expression



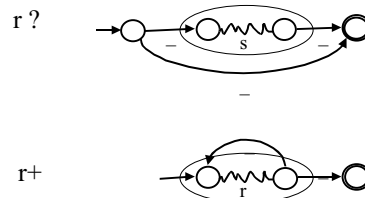
If r and s are regular expressions with the NFA's

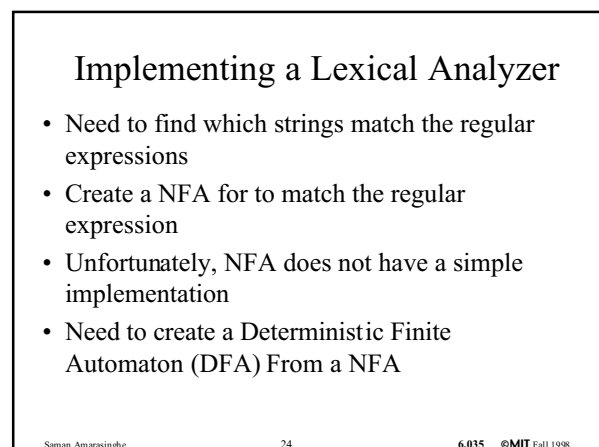
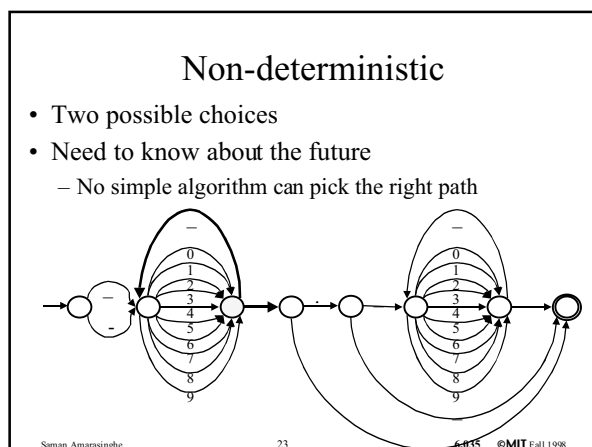
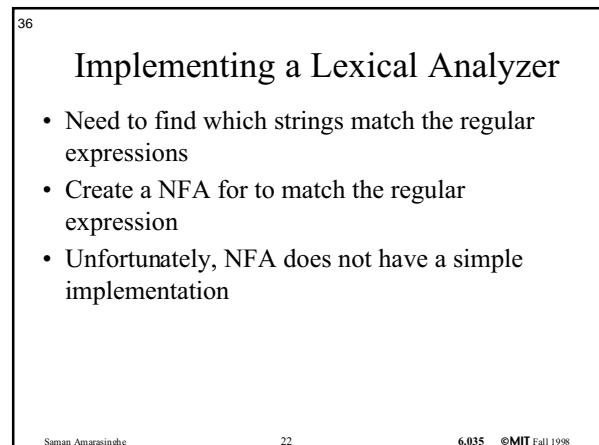
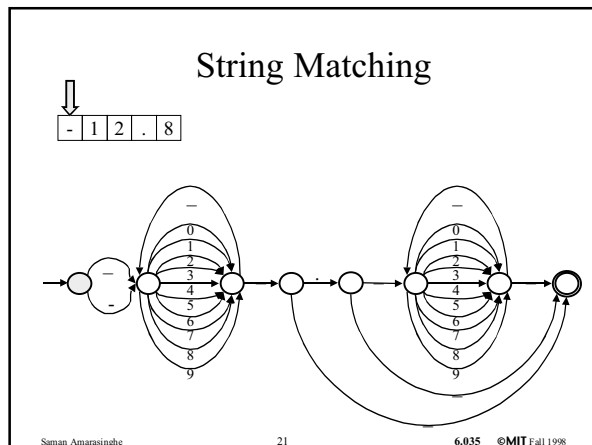
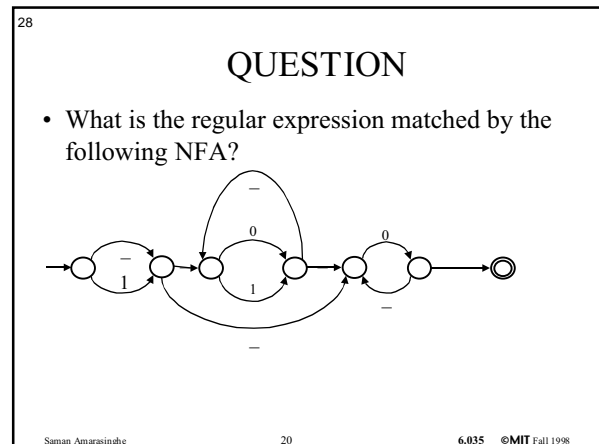
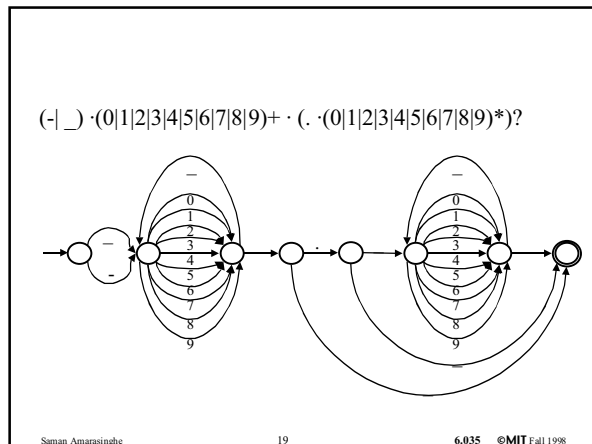


Constructing a NFA for a regular expression



Constructing a NFA for a regular expression





Constructing a DFA from a NFA

- Why do we need a DFA?
 - Easy to implement
 - Current state + input symbol uniquely identifies the next state
- How do you construct a DFA from a NFA?

1. Closure

- The closure of a state is the set of states that can be reached from that state without consuming any of the input
 - Closure(S) is the smallest set T such that

$$T = S \cup \bigcup_{s \in T} \text{edge}(s, \epsilon)$$

- Algorithm $T \mathcal{R} S$

repeat

$$T' \mathcal{R} T$$

$$T \mathcal{R} T \cup \bigcup_{s \in T'} \text{edge}(s, \epsilon)$$

until $T = T'$

$$S = \{1\}$$

$$T = \{1, 2\}$$

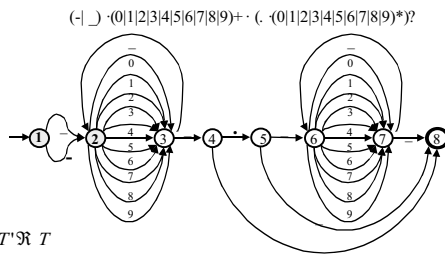
$$T' = \{1, 2\}$$

$T \mathcal{R} S$
repeat

$$T' \mathcal{R} T$$

$$T \mathcal{R} T \cup \bigcup_{s \in T'} \text{edge}(s, \epsilon)$$

until $T = T'$

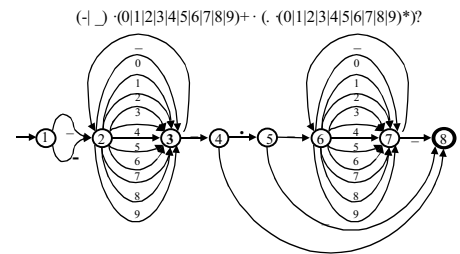


44

Question: What is closure(3)?

$$S = \{3\}$$

$$T = ??$$



46

2. DFAedge

- Given a symbol and a state, what states can you reach?

$$\text{DFAedge}(b, c) = \text{closure}(\text{edge}(\text{closure}(b), c))$$

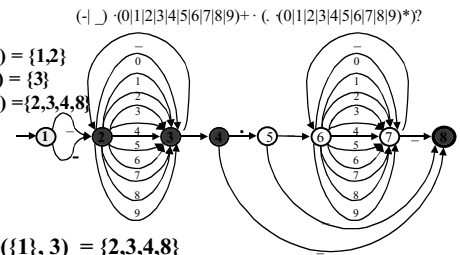
What is DFAedge({1}, 3)?

$$d = \{1, 2\}$$

$$\text{closure}(\{1\}) = \{1, 2\}$$

$$\text{edge}(\{2\}, 3) = \{3\}$$

$$\text{closure}(\{3\}) = \{2, 3, 4, 8\}$$



$$\text{DFAedge}(\{1\}, 3) = \{2, 3, 4, 8\}$$

49

Question: What is $\text{DFAedge}(\{3\}, .)$?

$d = \{3\}$

$(-|_) \cdot (0|1|2|3|4|5|6|7|8|9)^+ \cdot (-|_) \cdot (0|1|2|3|4|5|6|7|8|9)^*$

$\text{DFAedge}(\{3\}, .) = ??$

Saman Amarasinghe 31 6.035 ©MIT Fall 1998

51

3. DFA

```

states[0]  $\mathcal{R}$  1
states[1]  $\mathcal{R}$  closure(q)
p  $\mathcal{R}$  1
j  $\mathcal{R}$  0
while j  $\leq$  p
  foreach c  $\in$   $\Sigma$ 
    e  $\mathcal{R}$   $\text{DFAedge}(\text{states}[j], c)$ 
    if e = states[i] for some i  $\leq$  p
      then
        trans[j, c]  $\mathcal{R}$  i
    else
      p  $\mathcal{R}$  p + 1
      states[p]  $\mathcal{R}$  e
      trans[j, c]  $\mathcal{R}$  p
      j  $\mathcal{R}$  j + 1
  
```

Saman Amarasinghe 32 6.035 ©MIT Fall 1998

NFA

DFA

Saman Amarasinghe 33 6.035 ©MIT Fall 1998

54

In practice

- Uses automated tools to construct a lexical analyzer
 - Given a set of tokens defined using regular expressions, the tools will generate a character stream tokenizer by constructing a DFA
- Common scanner generator tools
 - lex in C
 - JLex in java
- Tonight, will talk about how to use JLex

Saman Amarasinghe 34 6.035 ©MIT Fall 1998

Summary

- Lexical analyzer create tokens out of a text stream
- Tokens are defined using regular expressions
- Regular expressions can be mapped to Nondeterministic Finite Automats (NFA)
 - by simple construction
- NFA is transformed to a DFA
 - Transformation algorithm
 - Executing a DFA is straightforward

Saman Amarasinghe 35 6.035 ©MIT Fall 1998

Next 3 Lectures

Saman Amarasinghe 36 6.035 ©MIT Fall 1998

Recitation on Programming and Collaboration Tools

- **TONIGHT!**
- From 7:30PM to 9:30PM
- Room: HERE!
- Will Discuss
 - CVS
 - Make
 - Athena Java tools
 - 6.035 locker
 - JLEX for the 1st assignment

Roadmap

Monday 9/4	Tuesday 9/5	Wednesday 9/6 Room 3-370 L1: Intro to compilers, course admin. info., lexical analysis	Thursday 9/7 Room 3-270 TODAY L2: LR(0) Paring Algorithms 7:30 – 9:30
Monday 9/11	Tuesday 9/12 Room 3-270 Next Lecture L3: Syntax analysis, bottom-up parsing	Wednesday 9/13 Room 3-370 L4: LR(0) Paring Algorithms and Parsing Tables	Thursday 9/14 Room 3-270 L5: LR(1) & LALR(1) Paring Algorithms

Scanner Segment assigned:

Project due on 9/11