# LEXICAL ANALYSIS
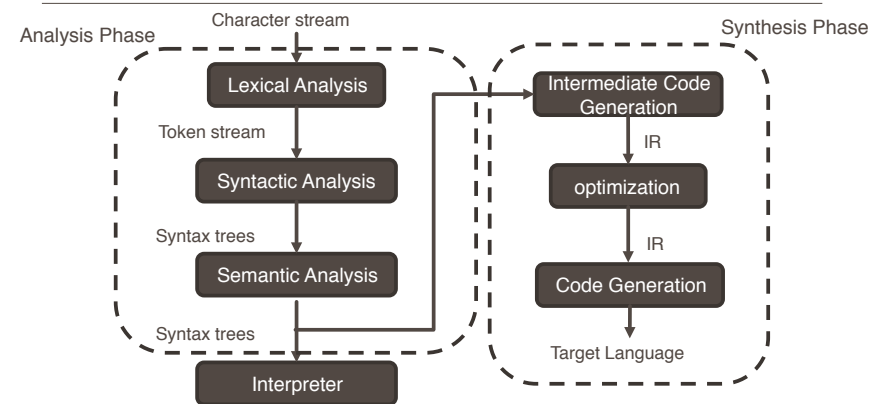
Baishakhi Ray
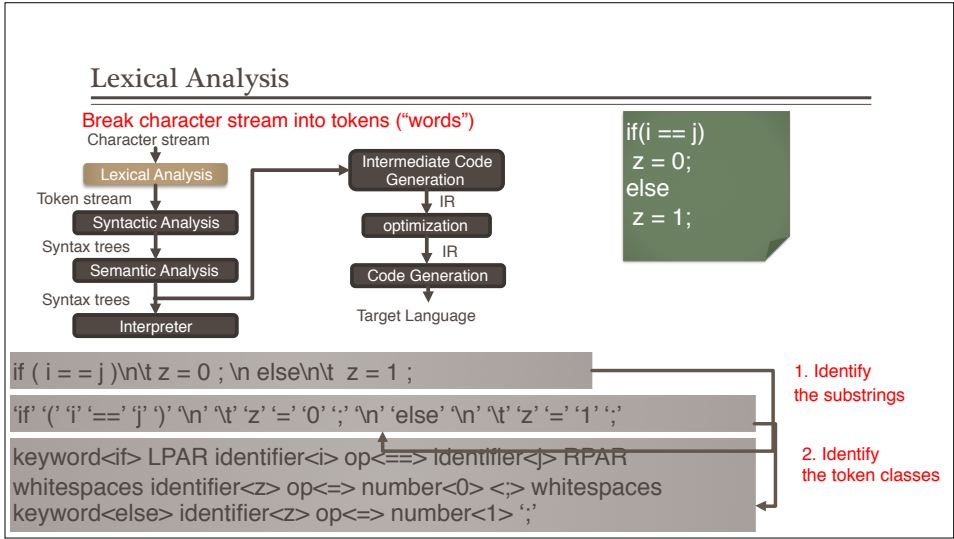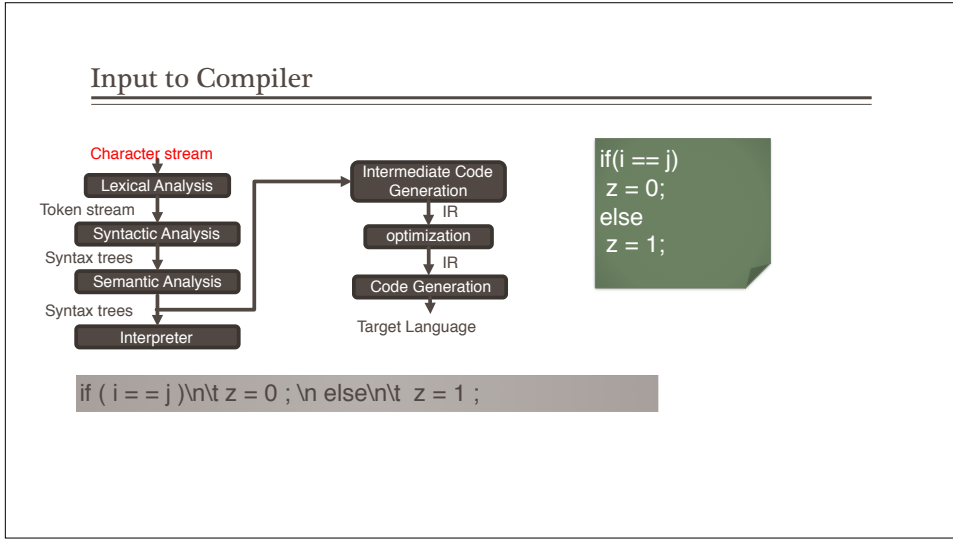
Fall 2019

These slides are motivated from Prof. Alex Aiken: Compilers (Stanford)

1

---

## Structure of a Typical Compiler

Analysis Phase                    Character stream                    Synthesis Phase

Lexical Analysis → Intermediate Code Generation

Token stream

Syntactic Analysis

optimization ← IR

Syntax trees

Semantic Analysis

IR

Code Generation

Syntax trees

Target Language

Interpreter

2

## Input to Compiler

Lexical Analysis

Token stream

Syntactic Analysis

Syntax trees

Semantic Analysis

Syntax trees

Interpreter

Intermediate Code Generation

IR

optimization

IR

Code Generation

Target Language

if(i == j)
 z = 0;
else
 z = 1;

if ( i = = j )\n\t z = 0 ; \n else\n\t  z = 1 ;

3

---

## Lexical Analysis

Break character stream into tokens ("words")

Character stream

Lexical Analysis

Token stream

Syntactic Analysis

Syntax trees

Semantic Analysis

Syntax trees

Interpreter

Intermediate Code Generation

IR

optimization

IR

Code Generation

Target Language

if(i == j)
 z = 0;
else
 z = 1;

if ( i = = j )\n\t z = 0 ; \n else\n\t  z = 1 ;

'if' '(' 'i' '==' 'j' ')' '\n' '\t' 'z' '=' '0' ';' '\n' 'else' '\n' '\t' 'z' '=' '1' ';'

keyword<if> LPAR identifier<i> op<==> identifier<j> RPAR whitespaces identifier<z> op<=> number<0> <;> whitespaces keyword<else> identifier<z> op<=> number<1> ';'

1. Identify the substrings

2. Identify the token classes
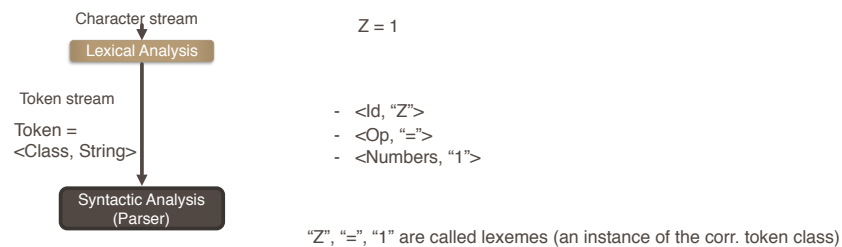
4

## Token Class

- In English?
  - Noun, verb, adjectives, …

- In Programming Language
  - keywords, identifiers, LPAR, RPAR, const, etc.

## Token Class

- Each class corresponds to a set of strings

- Identifier
  - Strings are letters or digits, starting with a letter
  - Eg:

- Numbers:
  - A non-empty strings of digits
  - Eg:

- Keywords
  - A fixed set of reserved words
  - Eg:

- Whitespace
  - A non-empty sequence of blanks, newlines, and tabs

## Lexical Analysis (Example)

- Classify program substrings according to roles (token class)
- Communicate tokens to parser

Character stream

Lexical Analysis

Token stream

Token =
<Class, String>

Syntactic Analysis
(Parser)

$Z = 1$

- <Id, "Z">
- <Op, "=">
- <Numbers, "1">

"Z", "=", "1" are called lexemes (an instance of the corr. token class)

## Lexical Analysis: HTML Examples

Here is a photo of <b> my house </b>
<p><img src="house.gif"/><br/>
see <a href="morePix.html">More Picture</a> if you liked that one.</p>

```
<text, "Here is a photo of">
<nodestart, b>
<text, "my house">
<nodeend, b>
<nodestart, p>
<selfendnode, img>
<selfendnode, br>
<text, "see">
<nodestart, a>
<text, "More Picture">
<nodeend, a>
<text, "if you liked that one.">
<nodeend, p>
```

## Exercise

```
x = p;
while ( x < 100 ) { x++ ; }
```

## Exercise

if(i == j)
 z = 0;
else
 z = 1;

==/=?

Keyword/Identifier?

## Lookahead

- Lexical analysis tries to partition the input string into the logical units of the language. This is implemented by reading left to right. "scanning", recognizing one token at a time.

- "Lookahead" is required to decide where one token ends and the next token begins.

```
if(i == j)
  z = 0;
else
  z = 1;
```

==/=?

Keyword/Identifier?

## Lookahead: Examples

- Usually, given the pattern describing the lexemes of a token, it is relatively simple to recognize matching lexemes when they occur on the input.

- However, in some languages, it is not immediately apparent when we have seen an instance of a lexeme corresponding to a token.

FORTRAN RULE: White Space is insignificant: VA R1 == VAR1

DO 5 I = 1,25

DO 5 I = 1.25

- Lexical analysis may require to "look ahead" to resolve ambiguity.
  - Look ahead complicates the design of lexical analysis
  - Minimize the amount of look ahead

## Lexical Analysis: Examples

- C++ template Syntax:
  - Foo<Bar>

- C++ stream Syntax:
  - cin >> var

- Ambiguity
  - Foo<Bar<Barq>>
  - cin >> var

## Summary So Far

- The goal of Lexical Analysis
  - Partition the input string to lexeme
  - Identify the token class of each lexeme

- Left-to-right scan => look ahead may require

  - In reality, lookahead is always needed

  - Our goal is to minimize thee amount of lookahead

# REGULAR LANGUAGES

- Lexical structure of a programming language is a set of token classes.
- Each token class consists of some set of strings.
- How to map which set of strings belongs to which token class?
  - Use regular languages
- Use Regular Expressions to define Regular Languages.

## Regular Expressions

- Single character
  - 'c' = {"c"}
- Epsilon
  - $\varepsilon$ = {""}
- Union
  - A + B = {a | a $\epsilon$ A} $\cup$ {b | b $\epsilon$ B}
- Concatenation
  - AB = {ab | a $\epsilon$ A ^ b $\epsilon$ B}
- Iteration (Kleene closure)
  - $A^{*} = \bigcup_{i>=0} A^{i} = $ A.....A (i times)
  - $A^{p} = \varepsilon$ (empty string)

## Regular Expressions

- Def: The regular expressions over $\Sigma$ are the smallest set of expressions including

  R = $\varepsilon$

  | 'c', 'c' $\epsilon$ $\Sigma$

  | R + R

  | RR

  | R*

## Regular Expression Example

- $\Sigma = \{p,q\}$
  - $q^*$
  - $(p+q)q$
  - $p^*+q^*$
  - $(p+q)^*$

- There can be many ways to write an expression

---

## Exercise

Choose the regular languages that are equivalent to the given regular language: $(p + q)^*q(p + q)^*$

A. $(pq + qq)^*(p + q)^*$

B. $(p + q)^*(qp + qq + q)(p + q)^*$

C. $(q + p)^*q(q + p)^*$

D. $(p + q)^*(p + q)(p + q)^*$

## Formal Languages

- Def: Let $\Sigma$ be a set of character (alphabet). A language over $\Sigma$ is a set of strings of characters drawn from $\Sigma$.
  - Regular languages is a formal language
- Alphabet = English character, Language = English Language
  - Is it formal language?
- Alphabet = ASCII, Language = C Language

## Formal Language

'c' = {"c"}

$\varepsilon$ = {""}

A + B = {a | a $\epsilon$ A} $\cup$ {b | b $\epsilon$ B}

AB = {ab | a $\epsilon$ A $\wedge$ b $\epsilon$ B}

$$A^* = \bigcup_{i>=0} A^i$$

expression

Set

## Formal Language

L('c') = {"c"}

$L(\varepsilon)$ = {""}

L(A + B) = {a | a $\epsilon$ L(A)} $\cup$ {b | b $\epsilon$ L(B)}

L(AB) = {ab | a $\epsilon$ L(A) ^ b $\epsilon$ L(B)}

$$L(A^*) = \bigcup_{i>=0} L(A^i)$$

expression

Set

L: Expressions -> Set of strings
- Meaning function L maps syntax to semantics
- Mapping is many to one
- Never one to many

## Lexical Specifications

- Keywords: "if" or "else" or "then" or "for" ….
  - Regular expression = 'i' 'f' + 'e' 'l' 's' 'e'
    = 'if' + 'else' + 'then'

- Numbers: a non-empty string of digits
  - digit = '1'+'0'+'2'+'3'+'4'+'5'+'6'+'7'+'8'+'9'
  - digit*
  - How to enforce non-empty string?
    - digit digit* = digit+

23

24

## Lexical Specifications

- Identifier: strings of letters or digits, starting with a letter
  - letter = 'a' + 'b' + 'c' + …. + 'z' + 'A' + 'B' + …. + 'Z'
    - = [a-zA-Z]
  - letter (letter + digit)*

- Whitespace: a non-empty sequence of blanks, newline, and tabs
  - (' ' + '\n' + '\t')+

## PASCAL Lexical Specification

- digit = '0'+'1'+'2'+'3'+'4'+'5'+'6'+'7'+'8'+'9'

- digits = digit+

- opt_fraction = ('.' digits) + $\varepsilon$ = ('.' digits)?

- opt_exponent = ('E' ('+' + '-' + $\varepsilon$) digits ) + $\varepsilon$

  = ('E' ('+' + '-')? digits )?

- num = digits opt_fraction opt_exponent

## Common Regular Expression

- At least one $A^+ \equiv AA^*$

- Union: $A \mid B \equiv A + B$

- Option: $A? \equiv A + \varepsilon$

- Range: 'a' + … + 'z' = [a-z]

- Excluded range: complement of [a-z] $\equiv$ [^a-z]

27

## Summary of Regular Languages

- Regular Expressions specify regular languages

- Five constructs
  - Two base expression
    - Empty and 1-character string

  - Three compound expressions
    - Union, Concatenation, Iteration

28

## Lexical Specification of a language

1. Write a regex for the lexemes of each token class
   - Number = digit$^+$
   - Keywords = 'if' + 'else' + ..
   - Identifiers = letter (letter + digit)*
   - LPAR = '('

## Lexical Specification of a language

2. Construct R, matching all lexemes for all tokens

   R = Number + Keywords + Identifiers + …

   = $R_1$ + $R_2$ + $R_3$ + …

3. Let input be $x_q \ldots x_n$.

   For $1 \leq i \leq n$, check $x_1 \ldots x_i \, \epsilon \, L(R)$

4. If successful, then we know that

   $x_1 \ldots x_i \, \epsilon \, L(R_j)$ for some j

5. Remove $x_1 \ldots x_i$ from input and go to step 3.

## Lexical Specification of a language

- How much input is used?
  - $x_1 \ldots x_i \in L(R)$
  - $x_1 \ldots x_j \in L(R)$, $i \neq j$
  - Which one do we want? (e.g., == or =)
  - Maximal munch: always choose the longer one

- Which token is used if more than one matches?
  - $x_1 \ldots x_i \in L(R)$ where $R = R_1 + R_2 + .. + R_n$
  - $x_1 \ldots x_i \in L(R_m)$
  - $x_1 \ldots x_i \in L(R_n)$, $m \neq n$
  - Eg: Keywords = 'if', Identifier = letter (letter + digit)*, if matches both
  - Keyword has higher priority
  - Rule of Thumb: Choose the one listed first

## Lexical Specification of a language

- What if no rule matches?
  - $x_1 \ldots x_i \notin L(R)$ … compiler typically tries to avoid this scenario
  - Error = [all strings not in the lexical spec]
  - Put it in last in priority

## Summary so far

- Regular Expressions are concise notations for the string patterns

- Use in lexical analysis with some extensions
  - To resolve ambiguities
  - To handle errors

- Implementation?
  - We will study next

## Finite Automata

- Regular Expression = specification

- Finite Automata = implementation

- A finite automaton consists of
  - An input Alphabet: $\Sigma$
  - A finite set of states: S

  - A start state: n

  - A set of accepting states: $F \subseteq S$

  - A set of transitions state: state1 $\xrightarrow{input}$ state2

a

## Transition

- s1 $\overset{a}{\to}$ s2 (state s1 on input a goes to state s2)

- If end of the input and in final state, the input is accepted

- Otherwise reject


- Language of FA = set of strings accepted by that FA

## Example Automata

- a finite automaton that accepts only "1"

## Example Automata

- A finite automaton that accepting any number of "1" followed by "0"

## Regular Expression to NFA

- For $\varepsilon$ (it's a choice)



- For input a

## Finite Automata

- Deterministic Finite Automata (DFA)
    - One transition per input per state
    - No $\varepsilon$-moves
    - Takes only one path through the state graph

- Nondeterministic Finite Automata (NFA)
    - Can have multiple transitions for one input in a given state
    - Can have $\varepsilon$-moves
    - Can choose which path to take
        - An NFA accepts if some of these paths lead to accepting state at the end of input.

## Finite Automata

- An NFA can get into multiple states



- Input:      1        0            0
- Output: {A}.     {A,B}        {A,B,C}

## NFA vs. DFA

- NFAs and DFAs recognize the same set of regular languages

- DFAs are faster to execute
  - No choices to consider

- NFAs are, in general, small

---

43



## Finite Automata

- For each kind of regex, define an equivalent NFA
  - Notation: NFA for regex M

44

## Regular Expression to NFA

- For $\varepsilon$



- For input a



45

## Regular Expression to NFA

- For AB



- For A + B



46

## Regular Expression to NFA

- For A*

## Example

- (q+p)*q

## NFA to DFA

```
Lexical Specification
        ↓
Regular Expressions
        ↓
       NFA  ──────→  DFA
          ↘         ↙
    Table driven implementation of
              automata
```

49

## ε-closure

- ε-closure of a state is all the state I can reach following ε move.



ε-closure(B) = {B,C,D}

50

## $\varepsilon$-closure

- $\varepsilon$-closure of a state is all the state I can reach following $\varepsilon$ move .



$\varepsilon$-closure(B) = {B,C,D}
$\varepsilon$-closure(G) = {A,B,C,D,G,H,I}

## NFA

- An NFA can be in many states at any time

- How many different states?
  - If NFA has N states, it reaches some subset of those states, say S
  - |S| $\leq$ $N$
  - There are $2^N$ – q possible subsets (finite number)

## NFA to DFA

### NFA

- States S
- Start s
- Final state F
- Transition state
  - $a(X) = \{y \mid x \in X \bigwedge x \xrightarrow{a} y)$
- $\varepsilon - closure$

### DFA

- States will be all possible subset of S except empty set
- Start state = $\varepsilon - closure(s)$
- Final state $\{X \mid X \cap F = \varnothing\}$
- $X \xrightarrow{a} Y$ if
  - $Y = \varepsilon - closure\ (a(X))$

## NFA to DFA
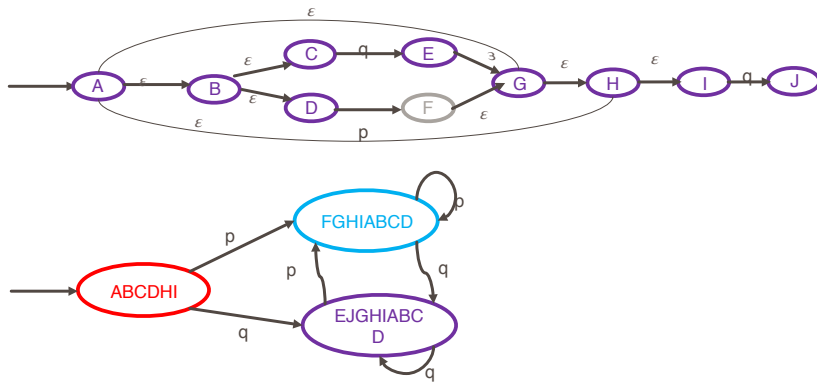
## NFA to DFA
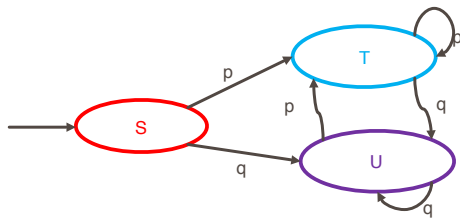
57

## NFA to DFA

58

## NFA to DFA



## Implementing DFA

- A DFA can be implemented by a 2D table T
  - One dimension is states
  - Another dimension is input symbol
  - For every transition $s_i ->^a s_k$: define $T[i,a] = k$

## Implementing DFA

Table A

| | p | q |
|---|---|---|
| S | T | U |
| T | T | U |
| U | T | U |



```
i = p;
state = p;
while(input[i])
    state = A[state,input[i]]
```

---

## Implementing DFA

Table A

| | p | q |
|---|---|---|
| S | T | U |
| T | T | U |
| U | T | U |

A lot of duplicate entries

Table B



| p | q |
|---|---|
| T | U |

Compact but need an extra indirection
- Inner loop will be slower

## Implementing DFA



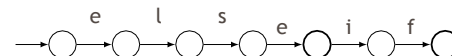| | p | q | |
|---|---|---|---|
| A | | | {B,H} |
| B | | | {C,D} |
| C | | {E} | |
| ... | | | |

Deal with set of states rather than single state-→    inner loop is complicated

## Deterministic Finite Automata: Example

```
{
    type token = ELSE | ELSEIF
}

rule token =
  parse "else"{ ELSE }
      | "elseif"{ ELSEIF }
```

# Deterministic Finite Automata

```
{ type token = IF | ID of string | NUM of string }

rule token =
  parse "if"{ IF }
      | ['a'-'z'] ['a'-'z' 'p'-'9'] as lit { ID(lit) }
      | ['p'-'9']+                   * as num { NUM(num) }
```