# Compilers

## Lexical Specification

Alex Aiken

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

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

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

- Range:   $'a' + 'b' + \ldots + 'z'$     $\equiv$   $[a\text{-}z]$

- Excluded range:

     complement of $[a\text{-}z]$   $\equiv$   $[\text{\^{}}a\text{-}z]$

Alex Aiken

- Last lecture: a specification for the predicate

$$s \in L(R)$$

- Not enough!

1. Write a rexp for the lexemes of each token class

- Number = digit$^+$

- Keyword = 'if' + 'else' + …

- Identifier = letter (letter + digit)*

- OpenPar = '( '

- …

Alex Aiken

2. Construct R, matching all lexemes for all tokens

R = Keyword + Identifier + Number + …

$\quad$ = $R_1$ + $R_2$ + …

3. Let input be $x_1...x_n$

   For $1 \leq i \leq n$ check

$$x_1...x_i \in L(R)$$

4. If success, then we know that

$$x_1...x_i \in L(R_j) \text{ for some } j$$

5. Remove $x_1...x_i$ from input and go to (3)

- How much input is used?

- Which token is used?

- What if no rule matches?

- Regular expressions are a concise notation for string patterns

- Use in lexical analysis requires small extensions
  - To resolve ambiguities
  - To handle errors

- Good algorithms known
  - Require only single pass over the input
  - Few operations per character (table lookup)

Alex Aiken