# Compilers

Lexical Specification

Alex Aiken

- At least one:  $A^+$         $\equiv$   $AA^*$
- Union:   $A \mid B$         $\equiv$   $A + B$
- Option: $A?$         $\equiv A + \varepsilon$
- Range:   $'a' +' b' +...+' z'$     $\equiv$   $[a\text{-}z]$
- Excluded range:

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

Alex Aiken

- Last lecture: a specification for the predicate

$$s \in L(R)$$

Sets of strings

- Not enough!

$$c_1 \, c_2 \, c_3 \mid c_4 \, c_5 \, c_6 \, c_7 \mid \; \cdots$$

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

- Number = $digit^+$

- Keyword = 'if' + 'else' + …

- Identifier = letter (letter + digit)*

- OpenPar = '( '

- …

2. Construct R, matching all lexemes for all tokens

R = Keyword + Identifier + Number + …

$\quad$ = $R_1$ + $R_2$ + …

3.  Let input be $x_1 \ldots x_n$

   For $1 \leq i \leq n$ check

$$x_1 \ldots x_i \in L(R) \ ?$$

4.  If success, then we know that

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

$$R = R_1 + R_2 + R_3 + \cdots$$

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

- How much input is used?

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

$$x_1 \dots x_j \in L(R)$$

$$i \neq j$$

"Maximal Munch"

Choose the one listed first

- Which token is used?

$$x_1 \ldots x_i \in L(R) \qquad R = R_1 + \ldots + R_N$$

$$x_1 \ldots x_i \in L(R_j)$$

$$x_1 \ldots x_i \in L(R_k)$$

$$if \in \begin{array}{l} L(Keywords) \\ L(Identifiers) \end{array} \begin{cases} Keywords = 'if' + 'else' + \ldots \\ Identifiers = letter(letter + digit)^* \end{cases}$$

- ## What if no rule matches?

$$X_1 \ldots X_i \notin L(R) \;] \; \text{No!}$$

$$\text{Error} = \left[ \begin{array}{l} \text{all strings not in the} \\ \text{lexical spec} \end{array} \right]$$

$$\rightarrow \text{put it last in priority}$$

Alex Aiken

- Regular expressions are a concise notation for string patterns

- Use in lexical analysis requires small extensions
  - To resolve ambiguities — *matches as long as possible*
  - To handle errors

    *hishest priority match*

- Good algorithms known

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

Alex Aiken