# Compilers

## Lexical Analysis

Alex Aiken

1. Lexical Analysis
2. Parsing
3. Semantic Analysis
4. Optimization
5. Code Generation

Alex Aiken

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

\tif (i == j)\n\t\tz = 0;\n\telse\n\t\tz = 1;

- Token Class (or Class)
  - In English:

    *Noun, verb, adjective, ...*
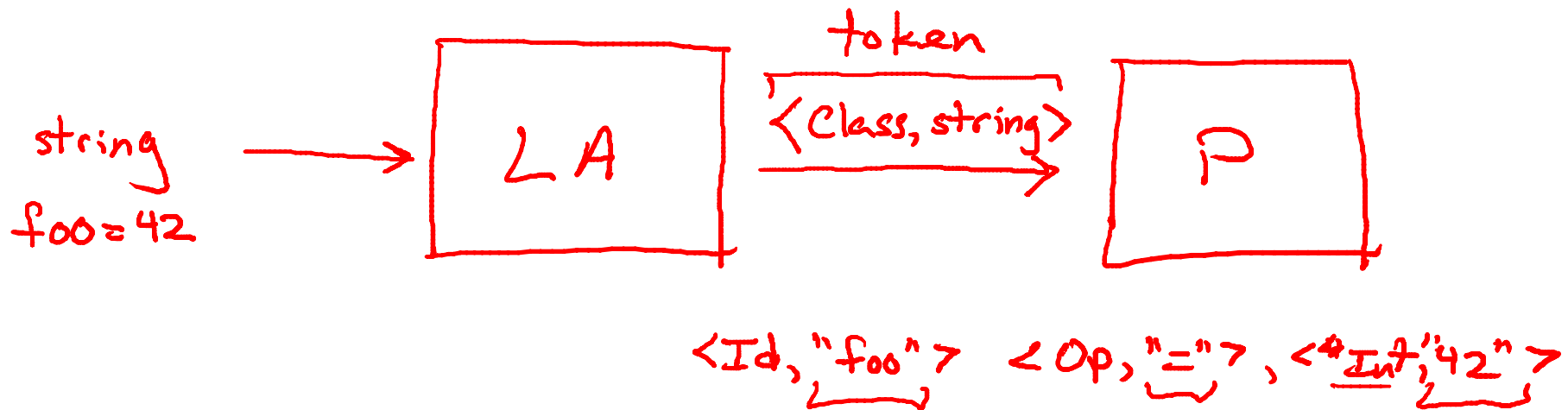
  - In a programming language:

    *Identifier, Keywords, '(', ')', Numbers, ...*

- Token classes correspond to sets of strings.

- Identifier:
  - *strings of letters or digits, starting with a letter*
- Integer:  _A1    _Foo    _B17
  - *a non-empty string of digits*    0    12    001    00
- Keyword:
  - *"else" or "if" or "begin" or …*
- Whitespace:
  - *a non-empty sequence of blanks, newlines, and tabs*

if _ _ _ ( → whitespace

Alex Aiken

- Classify program <u>substring</u>s according to <u>role</u>

  *token class*

- Communicate tokens to the parser

*token*

string

foo = 42

$\rightarrow$

LA

$\langle Class, string \rangle$

$\rightarrow$

P

$\langle Id, "foo" \rangle \quad \langle Op, "=" \rangle, \langle Int, "42" \rangle$

\t if (i==j)\n\t\tz = 0;\n\telse\n\t\tz = 1;

Operator
Whitespace
Keywords
Identifiers
Numbers

(
)
;
=

For the code fragment below, choose the correct number of tokens in each class that appear in the code fragment

$$x = 0;\backslash n\backslash twhile\ (x < 10)\ \{\backslash n\backslash tx++;\backslash n\}$$

○ W = 9; K = 1; I = 3; N = 2; O = 9

○ W = 11; K = 4; I = 0; N = 2; O = 9

○ W = 9; K = 4; I = 0; N = 3; O = 9

○ W = 11; K = 1; I = 3; N = 3; O = 9

W: Whitespace
K: Keyword
I: Identifier
N: Number
O: Other Tokens:
    { } ( ) < ++ ; =

- An implementation must do two things:

  1. Recognize substrings corresponding to <u>tokens</u>
     - The *lexemes*

  2. Identify the <u>token class</u> of each lexeme

$$\underbrace{< \text{token class}, \text{ lexeme}>}_{\text{token}}$$

Alex Aiken