

Token classes correspond to sets of strings.

Identifier:

- strings of letters or digits, starting with a letter.

Integer:

- a non-empty string of digits.

Keyword:

- "else" or "if" or "begin" or ...

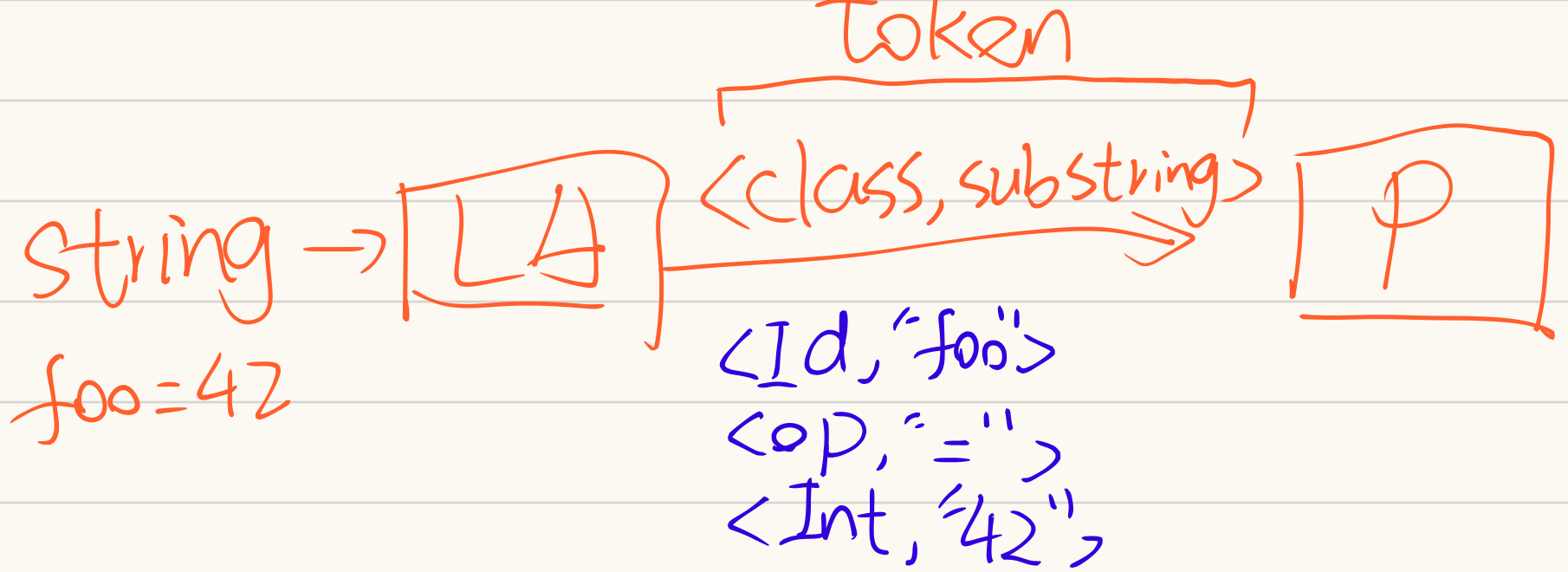
Whitespace:

- a non-empty sequence of blanks, newlines, and tabs.

Classify program substrings according to role.

token class

Communicate tokens to the parser.



An Implementation must do two things:

1. Recognize substrings corresponding to tokens
 - The lexemes 词素
2. Identify the token class of each lexeme.
 - $\langle \text{token class, lexeme} \rangle$
token

FORTRAN rule: Whitespace is insignifi-
-cant

VARI is the same as VA RI

1. The goal is to partition the string. This is implemented by reading left-to-right, recognizing one token at a time.

2. 'Lookahead' may be required to decide where one token ends and the next token begins.

The goal of LA is to

- Partition the input string into lexemes
- Identify the token of each lexeme

Left-to-right scan \Rightarrow lookahead sometimes required

Lexical structure = token classes

We must say what set of strings is in a token class.

- Use regular languages.

Regular Languages

- Single character

$$'c' = \{ "c" \}$$

two base cases

- Epsilon

$$\epsilon = \{ "" \}$$

three compound expressions

- Union

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

- Concatenation

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

- Iteration

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

$$A^i = \underbrace{A \dots A}_{i \text{ times}}$$

$$A^0 = \epsilon$$

Def. The regular expression over Σ are the smallest set of expressions including

Def. Let Σ be a set of characters (an alphabet) A language over Σ is a set of strings of characters drawn from Σ .

Alphabet = English characters

Language = English sentences

Alphabet = ASCII

Language = C program

Meaning function L maps syntax to semantics.

$L(e) = M$
 \downarrow \downarrow
reg exp set of strings

Why use a meaning function?

- Make clear what is syntax, what is semantics
- Allows us to consider notation as a separate issue.
- Because expressions and meanings are not 1-1

Meaning is many to one, and never
one to many!