

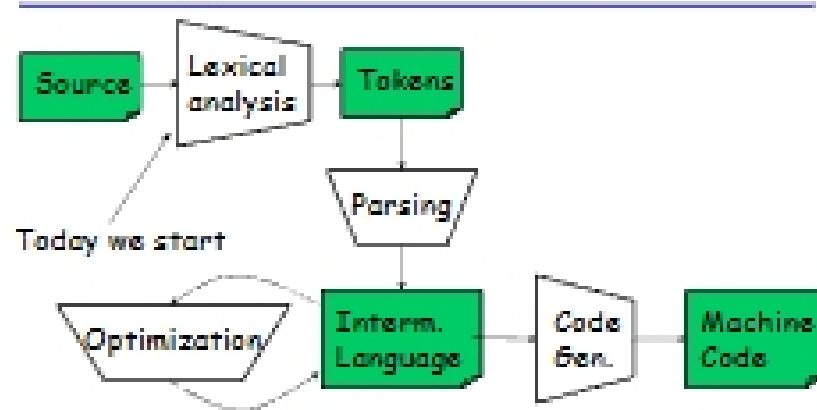
Lexical Analysis

1/1/2006

CS 685, November 11, 2006

1

The Structure of a Compiler



1/1/2006

CS 685, November 11, 2006

2

Lexical Analysis

- What do we want to do? Example:

```
if (i == j)
  z = 0;
else
  z = 1;
```
- The input is just a sequence of characters:

```
\tif (i == j)\n\tz = 0;\n\telse\n\tz = 1;
```
- Goal: Partition input string into substrings
- And classify them according to their role

1/1/2006

CS 685, November 11, 2006

3

What's a Token?

- Output of lexical analysis is a stream of tokens
- A token is a syntactic category
 - In English:
noun, verb, adjective, ...
 - In a programming language:
Identifier, (Keyword) Integer, Float, IF, THEN ...
- Parser relies on the token distinctions:
 - E.g., identifiers are treated differently than keywords

1/1/2006

CS-4375, Summer 11, 2006

4

Tokens

- A token corresponds to a string or a set of strings.
- Identifier: *strings of letters or digits, starting with a letter*
- Integer: *a non-empty string of digits*
- (Keyword) IF: *string if*
- Whitespace: *a non-empty sequence of blanks, newlines, and tabs*
- OpenPar: *a left-parenthesis*

1/1/2006

CS-4375, Summer 11, 2006

5

Lexical Analyzer: Implementation

- An implementation must do two things:
 1. Recognize substrings corresponding to tokens
 2. Return the value or lexeme of the token
 - The lexeme is the substring

1/1/2006

CS-4375, Summer 11, 2006

6

Example

- When reading the following:
`\tif (i == j)\n\t\tz = 0;\n\telse\n\t\tz = 1;`
- Token-lexeme pairs returned by the lexer:
 - (Whitespace, "\t")
 - ((Keyword) IF, "if")
 - (LeftPar, "(")
 - (Identifier, "i")
 - (Relation, "==")
 - (Identifier, "j")
 - ...

1/1/2006

CS 68A, Summer 11, 2006

7

Lexical Analyzer: Implementation

- The lexer usually discards "uninteresting" tokens that don't contribute to parsing.
- Examples: Whitespace, Comments
- Question: What happens if we remove all whitespace and all comments prior to lexing?

1/1/2006

CS 68A, Summer 11, 2006

8

Lookahead.

- Two important points:
 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
 - Even our simple example has lookahead issues
i vs. if
= vs. ==

1/1/2006

CS 68A, Summer 11, 2006

9
