

CGTM 106

David Gries

March 1968

Compiler Implementation System

Strings and input-output

This note describes the uses of strings in compiler-compiler language. Strings are introduced only to make input-output easier (error messages, etc.). This note supercedes any mention of strings in earlier notes.

1. Introduction
2. String variables
3. String assignments
4. <string>s
5. Related system names
6. Output

String manipulation

1. Introduction

Strings are intended to be used only for input-output functions — changing the input source program, inserting or deleting characters from it, printing the source program and error messages, etc. Internally, source program characters and atoms themselves are not used since they are represented by 16-bit (BYTE2) entities; thus string manipulation should not be necessary.

2. String variables

The lengths of string variables are fixed at meta-compile time and may not exceed 255.

A simple string may be defined by

```
<declaration> ::= STRING 

|                                                                                  |
|----------------------------------------------------------------------------------|
| $\langle \text{unsigned integer} \rangle$ $\langle \text{identifier} \rangle$    |
| $\langle \text{identifier} \rangle$ . = $\langle \text{string constant} \rangle$ |


```

With the first alternate, the $\langle \text{unsigned integer} \rangle$ specifies the length. In the second case, the $\langle \text{identifier} \rangle$ is initialized to $\langle \text{string constant} \rangle$ (see 3a); the length is the number of characters in the $\langle \text{string constant} \rangle$.

Tables (arrays) of strings are STATIC only, and are declared by

```
<declaration> ::= STRING ( $\langle \text{unsigned integer} \rangle_1$ ) TABLE ( $\langle \text{unsigned integer} \rangle_2$ )  
  
                   $\langle \text{identifier} \rangle$  . = 

|                                               |
|-----------------------------------------------|
| $\langle \text{string constant} \rangle$ [ ]* |
|-----------------------------------------------|

  
  
                  ::= STRING TABLE  $\langle \text{identifier} \rangle$  . = 

|                                               |
|-----------------------------------------------|
| $\langle \text{string constant} \rangle$ [ ]* |
|-----------------------------------------------|


```

In the first case, $\langle \text{unsigned integer} \rangle_1$ specifies the length of each string in the static table, while $\langle \text{unsigned integer} \rangle_2$ specifies the number of strings. If n $\langle \text{string constant} \rangle$ s appear, the first n elements are initialized to these $\langle \text{string constant} \rangle$ s.

In the second case, the number of strings in the table is the number of string constants; the length of each element is the length of the corresponding $\langle \text{string constant} \rangle$ to which it is initialized.

3. String assignments

A string assignment has the form

$\langle \text{assignment} \rangle ::= \langle \text{string variable} \rangle . = \langle \text{string} \rangle$

$\langle \text{assignment} \rangle ::= \text{SUBSTR} (\langle \text{string variable} \rangle, \langle \text{integer expression} \rangle_1, \langle \text{integer expr.} \rangle_2) . = \langle \text{string} \rangle$

In the second case, we must have $\text{length} (\langle \text{string variable} \rangle) \geq \langle \text{integer expr} \rangle_1 + \langle \text{integer expr} \rangle_2$.

In this case, $\langle \text{integer expression} \rangle_2$ characters of $\langle \text{string variable} \rangle$, beginning with the $\langle \text{integer expression} \rangle_1^{\text{th}}$, are replaced by the first

$\langle \text{integer expression} \rangle_2$ characters of $\langle \text{string} \rangle$. If $\langle \text{string} \rangle$ has too few

characters, blanks are added on the right; if too many, only the leftmost are used.

The default option for $\langle \text{integer expression} \rangle_2$ is $(\text{length} (\text{string variable}) - \langle \text{integer expression} \rangle_1 + 1)$.

The first form is equivalent to

$\text{SUBSTR} (\langle \text{string variable} \rangle, 1) . = \langle \text{string} \rangle$

4. $\langle \text{string} \rangle$ s

The following are $\langle \text{string} \rangle$ s.

- a) A $\langle \text{string constant} \rangle$ — of the form 'sequence of EBCDIC characters'

In the sequence the character ' must be represented by two apostrophes. Thus "B'C'" is the string " B'C' ".

- b) A string constant of the form X' (pairs of hexcharacters 0-9, A-F) '.

The value is the characters described by each pair of hex characters.

- c) TEXT (<expression>)

The <expression> must be the internal representation (BYTE2) of some atom or a pointer to some element of a DICT or the internal dictionary. The value of this <string> is the string of characters forming the atom on which the element is chained.

- d) A <syn identifier> - synonym for some source language atom.

- e) A variable declared as STRING.

- f) SUBSTR (<string>, <integer expression>₁, <integer expression>₂)

The value is the <integer expression>₂ characters of <string>, starting with the <integer expression>₁th character.

- g) BITTEXT (<expression>)

The value is the binary value of the expression as it appears in memory. (Thus BITTEXT (B'11010') is equal to '11010'.

- h) OCTTEXT (<expression>)

The value is the octal value of the expression: OCTTEXT (B'11010') is '32'.

- i) HEXTEXT (<expression>)

The value is the hex value of the expression: HEXTEXT (B'11010') is '1A'.

- j) INTTEXT (<expression>)

The expression is converted to a decimal integer and changed to character form.

- k) FLTTEXT (<expression>)

The value is the <expression> in some fixed, decimal floating point representation.

l) `<string>_1` CAT `<string>_2`

The result is the string `<string>_1` concatenated with the string `<string>_2`.

5. System names

`&INPUT` is a string variable which always contains the complete current source program line, as read in. It may not be assigned another value.

`&SCANNER` is a string variable containing the current source program line being used by the scanner (lexical analyzer). As characters are "used up" to form atoms they are deleted from `&SCANNER`.

`&POSIT`, type HWI, always contains the current position of the scanner in `&SCANNER` from the beginning of the original line. Thus if three characters on the current line have been "used up" to form atoms, `&POSIT` has the value 4.

`&LINE`, type HWI, contains the current output line number (the first card or line is number 1).

Strings can be added to the current beginning of the line in `&SCANNER` by the statement

```
&INSSCAN (<string>)
```

`&POSIT` is counted back the appropriate number of characters. This may be useful in error recovery procedures.

6. Output

One adds a string of characters to the end of the current line of

output by the statement

```
EOOUTPUT (<string>)
```

The statement

```
EOOUTPUT (<string>1,<string>2, ... ,<string>n)
```

is equivalent to

```
EOOUTPUT (<string>1); ... ; EOOUTPUT (<string>n)
```

When the current output line is filled up, it is written out and a new one started. Alternatively, one can cause the current line to be written out by executing

```
EOOUTPUT .
```

Execution of the statement

```
EOOUTDESCR(D)
```

where D is a DESCRIPTOR (see CGTM #104) or a pointer to a DESCRIPTOR causes the current line to be written out and the DESCRIPTOR D to be written out in some as yet undetermined format.