

User's Manual

FLOATING POINT INTERPRETER

093-000019-04

Basic Binary (5600-7600)	091-000012
Extended Binary (4100-7600)	091-000013
Extended Relocatable Binary	089-000046

Ordering No. 093-000019
© Data General Corporation, 1969, 1971, 1972, 1973
All Rights Reserved.
Printed in the United States of America
Rev. 04, March 1973

NOTICE

Data General Corporation (DGC) has prepared this manual for use by DGC personnel, licensees and customers. The information contained herein is the property of DGC and shall neither be reproduced in whole or in part without DGC prior written approval.

DGC reserves the right to make changes without notice in the specifications and materials contained herein and shall not be responsible for any damages (including consequential) caused by reliance on the materials presented, including but not limited to typographical or arithmetic errors.

Original Release - May 1969
First Revision - September 1971
Second Revision - December 1971
Third Revision - March 1972
Fourth Revision - March 1973

This revision of the Floating Point Interpreter User's Manual, 093-000019-04, supersedes 093-000019-03 and constitutes a minor revision. A list of changes is given at the back of the manual, following the index.

Most small, general purpose computers do not have hardware for the manipulation of floating point numbers. They have, therefore, implemented software packages to provide floating point functions. Two approaches can be taken to solve this software problem.

The first approach is to provide a number of subroutines, each of which performs a specific function, e.g. floating addition. This requires that the user pass inputs in a standard manner to the subroutine. Since many small computers have only two accumulators, only a single operand (represented in two words) can be passed in the accumulators. The address of the second operand must be provided with the subroutine call. A typical calling sequence requires two load accumulator instructions, a subroutine call, and an address word. One operand is always destroyed and replaced by the result (usually the accumulator operand). This means that an intermediate result cannot be tested without destroying an operand. A lengthy program requiring these subroutines cannot make efficient use of storage and the manipulation of operands becomes a tedious, cumbersome job.

The second approach is to provide an interpreter in which all the floating point functions are imbedded and which can simulate floating point registers in storage. The interpreter is given control by a subroutine call. Once in control, it accesses succeeding memory words and "interprets" a sixteen bit word in a manner similar to the hardware. These instructions are not executed in the hardware sense but are instead interpreted to provide extended machine features. The interpreter approach was adopted for the Nova-line floating point package.

The interpreter provides four floating point accumulators which can be addressed and manipulated in a manner similar to the actual hardware accumulators. Floating point instructions are syntactically similar to machine instructions* and are assembled in a similar way. For example, the instruction:

ADD# 0,2,SNR

which adds AC0 to AC2, skips the next instruction if the result is non-zero, and does not change AC0 or AC2, has a similar floating point version. The instruction

FADD# 0,2,FSNR

adds floating accumulator 0 (designated as FAC0) to FAC2, skips the next instruction if the result is non-zero, and does not destroy FAC0 or FAC2.

The interpreter can be implemented in Read Only Storage (ROS), and is completely reentrant. These properties, in addition to the floating point instruction set, are described in the body of this manual.

* The reader should be thoroughly familiar with machine instructions and the Nova assembler. Instructions are described in "How to Use the Nova Computers"; the absolute and relocatable assemblers are described in documents 093-000017 and 093-000040 respectively.

INTRODUCTION

Chapter 1 - Interpreter Structure and Floating Point Number Representation

Interpreter Structure	1-1
Floating Point Number Representation	1-1

Chapter 2 - Basic Interpreter Use and Command Structure

General	2-1
Basic and Extended Interpreters	2-1
The Writable Area	2-1
Initialization	2-2
Entering - Exiting the Interpreter Mode	2-2
Floating Point Instruction Set	2-2
Memory Reference Instructions	2-3
Arithmetic Instructions	2-4
Input/Output Instructions	2-6
Special Instructions	2-7
Illegal Instructions	2-7
Requirements for Reentrance	2-7

Chapter 3 - Extended Floating Point Interpreter

Features	3-1
Mathematical Functions	3-1
"F" Format Conversion	3-1

Chapter 4 - Extended Relocatable Floating Point Interpreter

Chapter 5 - Sample Programs

Square Root Newton Iteration	5-1
Polynomial Evaluation	5-1

Appendix A - Floating Point Instructions and Options

Appendix B - Floating Point Instruction Encoding

Appendix C - Extended Interpreter Mathematical Routines

Floating Point Arctangent	C-1
Floating Point Exponential	C-1
Floating Point Natural Logarithm	C-2
Floating Point Sine and Cosine	C-3
Floating Point Tangent	C-4
Floating Point Square Root	C-4

STRUCTURE OF THE INTERPRETER

The interpreter is non-destructive, i.e., no instruction is modified in any manner during execution. Further, it does not store temporary information within itself but uses a writable area which must be provided by the user.

Within the interpreter, subroutine linkage is via a push-down list, maintained in the writable area provided by the user. This property, coupled with the properties already mentioned, makes the interpreter reentrant. This means that if a second user routine requires the interpreter, it may interrupt the current routine, perform its function using the interpreter, and return to the first routine without affecting the state of the interpreter.

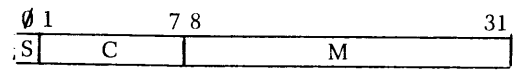
The absolute interpreters (basic binary and extended binary) require eight page 0 locations which the user must not destroy. These locations and their uses are:

- 004 contains the starting address of the interpreter.
- 005 contains the starting address of the initialization routine for the interpreter.
- 006 one temporary word used by the interpreter and saved by any routine reentering the package.
- 007 a word set up by the user to contain the address of a writable area for the interpreter.
- 040 a word containing the address of a user written "get character" routine.
- 041 a word containing the address of a user written "put character" routine.
- 042-043 words containing the linkage addresses to determine the location of extended code (if present).

In the relocatable version, the same functions are provided by six ZREL locations and two NREL locations. Those of interest to the user are ITMP, which corresponds to location 6 in the absolute versions and is declared as an entry by the interpreter, and WSA, GETC and PUTC, which correspond to locations 7, 40, and 41 respectively and which must be page zero relocatable. WSA, GETC and PUTC are defined as entries by the user program.

FLOATING POINT NUMBER REPRESENTATION*

Floating point numbers are internally stored in two consecutive 16-bit words. The form is:



S is the sign of the mantissa, M, in bits 8-31. The mantissa is considered to be a normalized six digit hexadecimal fraction, and the range of the magnitude of the mantissa is:

$$16^{*-1} \leq M \leq (1-16^{*-6}).$$

The characteristic, C, is the integer exponent of 16 in excess 64₁₀ code. The total range of magnitudes is:

$$16^{*-1} * 16^{*-63} \leq F \leq (1-16^{*-6}) * 16^{*63}$$

or approximately

$$2.4 * 10^{*-78} \leq F \leq 7.2 * 10^{*75}$$

Any operand having a zero mantissa is represented in true zero form, i.e., bits 0-31 are 0. Negative numbers are identical to their positive counterparts except S = 1 instead of 0.

The maximum error of a normalized mantissa is less than 16⁻⁶.

* See also "How to Use the NOVA Computers", Appendix C.

GENERAL

The interpreter provides four floating point accumulators. They are numbered 0, 1, 2, and 3 like the hardware accumulators. The designation FAC_n will be used for floating accumulator n . Arithmetic is performed accumulator-to-accumulator as with fixed point instructions. Operands can be accessed and stored using memory reference instructions. Instructions which reference floating point operands in memory should provide an address which points to the first word of the two word operand. If indexing is specified, the hardware index register is used. For example,

```
FLDA 1,4,2
```

loads FAC_1 from two consecutive words in memory whose first word address is $4 + C(AC_2)$.^{*} Certain instructions that manipulate hardware accumulators AC_2 and AC_3 will be described in the section 'Floating Point Instruction Set' on page 2-2. No facility is provided for manipulating AC_0 and AC_1 with the floating point instruction set.

BASIC AND EXTENDED INTERPRETERS

This chapter describes the features of the Basic Floating Point Interpreter. An extended version of the interpreter, in both binary^{**} and relocatable format, is available. The extended version has the same features described here; additional features of the Extended and Extended Relocatable Interpreters are described in Chapters 3 and 4 respectively.

THE WRITABLE AREA

The Basic Floating Point Interpreter requires 64 (decimal) words of contiguous writable memory. If the Extended Interpreter is used, 100 (decimal) words must be provided. The Extended Interpreter and its added features are described in detail in Chapter 3. Other than the difference in the length, the writable area of both versions is the same.

The address of the first word of the writable area provided by the user must be stored in location 007 of page 0 before any interpreter commands are executed. If a second routine (or third, etc.) may reenter the interpreter, the routine must provide an address pointing to a different memory area before reentrance is made.

^{*} $C(x)$ means 'contents of x '.

^{**} The Basic Interpreter (binary) requires locations 56008 - 76008, and the Extended Interpreter (binary) requires locations 41008 - 76008.

THE WRITABLE AREA (cont'd)

A number of flags are stored in this writable area and may be examined by the user. To access the flags, the user must exit from the interpretive mode (described on page 2-2). If an index register is loaded from location 007, the user may access the flag words by an instruction of the form:

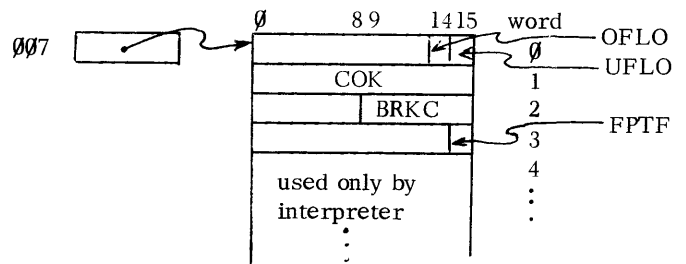
```
LDA 0,n,i
```

where: i is the index register (2 or 3).

n is a constant displacement as described below.

- $n = 0$: The first word contains the overflow/underflow flags. If the result of an operation is less than $16^{**}-64$, bit 15 of this word will be set to indicate underflow. If the result of an operation is greater than $(1-16^{**}-6)*16^{**}+63$, bit 14 of the word will be set to indicate overflow. Other conditions that may set overflow/underflow will be discussed where appropriate. This word is initially cleared. Once a bit is set it will remain set. It is the user's responsibility to reset these bits.
- $n = 1$: After an input conversion (see page 2-6), this word will be zero if no conversion was performed (because of an input error). Otherwise the word will be non-zero.
- $n = 2$: After an input conversion, this word will contain the 7-bit ASCII character, right justified, that served as the break character in the input stream.
- $n = 3$: After an input conversion, this word will be 0 if no decimal point (".") was encountered in the input stream. If a decimal point was seen, the word will be 1.

The diagram below summarizes these flags:



THE WRITABLE AREA (cont'd)

A convenient means of accessing the flags is to define four symbolic equivalences such as:

```

FLGS = 0 ; OVERFLOW/UNDERFLOW FLAGS
COK = 1 ; CONVERSION OK
BRKC = 2 ; BREAK CHARACTER
FPTF = 3 ; FLOATING POINT FLAG
    
```

Now if AC2 contains the address of the writable area, any flag may be accessed by statements like the following:

```

LDA 0, FLGS, 2
LDA 1, COK, 2
LDA 1, BRKC, 2
LDA 1, FPTF, 2
    
```

INITIALIZATION

The interpreter, or more correctly the writable area, must be initialized before floating point instructions are executed. Initialization should be given once for every writable area to be used. The command is

```

FINI
    
```

which generates the instruction

```

JSR @5
    
```

(Note that the initialization routine address is in location 005.) Location 007 must point to the writable area. This command destroys AC3, but preserves all other accumulators and Carry. The initialization routine clears the overflow/underflow flags and sets up linkage for the push-down list.

ENTERING-EXITING THE INTERPRETIVE MODE

To use the Floating Point Interpreter, it is necessary to distinguish between the processing of normal instructions (described in "How to Use the Nova Computers") and the processing of floating point instructions. Since floating point instructions are "interpreted" (not executed per se), the interpreter must be given control before floating point code will be executed properly. Whenever the interpreter is in control, this will be called the "interpretive mode." Otherwise, the machine will be referred to as in "normal mode." To enter the interpretive mode, the command is

```

FETR
    
```

ENTERING-EXITING THE INTERPRETIVE MODE (cont'd)

which generates the instruction:

```

JSR @4
    
```

As noted earlier, location 004 will contain the starting address of the interpreter. Once in the interpretive mode, only the floating point instruction set can be used. (Normal instructions will be decoded and executed as floating point instructions.) To return to normal processing, the command

```

FEXT
    
```

must be given.

FETR destroys AC3, but the contents of all other accumulators and Carry will be saved. AC2 and AC3 can be used for indexing in the interpretive mode. Certain instructions also enable modification of the contents of AC2 and AC3. FEXT will restore AC0, AC1, and Carry to their state before entering the interpreter. AC2 and AC3 will reflect any changes caused by floating point instructions that modified their contents.

The normal sequence of code using the interpreter would be of the form:

```

FINI ;INITIALIZE INTERPRETER
LDA 0, CNST
STA 0, TEMP ;NORMAL INSTRUCTIONS
FETR ;ENTER INTERPRETER
.
. ;FLOATING POINT INSTRUCTIONS
.
FEXT ;EXIT INTERPRETER
LDA 1, C5
    
```

FLOATING POINT INSTRUCTION SET

The set of instructions for the Basic Floating Point Interpreter obey the syntactic rules described in the Assembler Manual, 093-000017. For example, all floating ALC instructions require a source accumulator and a destination accumulator. All floating point instructions begin with an "F" to distinguish them from normal instructions. Appendix A summarizes the instructions and Appendix B gives their octal encoding.

Memory Reference Instructions

All memory reference instructions require an address, which is represented in the following discussion by "adr". The effective address is calculated exactly as in the standard Nova instruction set, except that the auto-incrementing, auto-decrementing properties of locations 20-37 have not been implemented. (For a discussion of effective address calculation, see "How to Use the Nova Computers", Section 2-1.)

The instruction

FLDA n, adr

causes FACn to be loaded with the two word operand at adr, adr+1.

The instruction

FSTA n, adr

causes FACn to be stored in memory at adr, adr+1.

The instruction

FJMP adr

causes control to be transferred to the floating point instruction at adr.

The instruction

FJSR adr

causes control to be transferred to the instruction at adr and AC3 to be set to the value of the current location counter + 1. For example, a floating point subroutine can be executed by the following:

```

                                ;MAIN PROGRAM
                                FLD3 0, LOC
                                FJMP SUBR
                                FSTA 0, RSLT
                                .
                                .
SUBR:                            ;SUBROUTINE
                                .
                                .
                                FJMP 0, 3    ;RETURN
    
```

The instruction

FFIX adr

causes the floating point number at adr, adr+1 to be

Memory Reference Instructions (cont'd)

converted to a fixed point, double precision integer (truncated) at adr, adr+1. If the conversion results in an integer whose absolute value is greater than $2^{24}-1$, the overflow flag will be set and $2^{24}-1$ will be returned as the magnitude. Note that while floating point numbers are represented in signed-magnitude format, fixed point values will always be represented in two's complement notation.

The instruction

FFLO adr

causes the fixed point, double precision integer at adr, adr+1 to be converted to a floating point number at adr, adr+1. Negative integers must be represented in two's complement format.

The instruction

FISZ adr

causes the contents of adr to be incremented by one and the next floating point instruction in sequence to be skipped if the result is zero.

The instruction

FDSZ adr

causes the contents of adr to be decremented by one and the next floating point instruction in sequence to be skipped if the result is zero. FISZ and FDSZ should be used with fixed point, single precision integers - not with floating point numbers.

The instruction

FST3 adr

causes AC3 to be stored at adr.

The instruction

FLD3 adr

causes AC3 to be loaded from the contents of adr. FST3 and FLD3 operate on real accumulator 3. FST3 is useful for saving the return address inside a floating point subroutine. Note that the return address should always be saved if a floating point subroutine exits and then enters the interpretive mode, since FETR destroys AC3. These two instructions also provide a means for initializing loop counts without leaving the interpretive mode.

Memory Reference Instructions (cont'd)

For example, the sequence:

```

LPI:      FEXT
          LDA 0,CNT
          STA 0,TEMP
          FETR
          :
          FEXT
          DSZ TEMP
          JMP LPI
          FETR
    
```

can be replaced by:

```

LPI:      FLD3 CNT
          FST3 TEMP
          :
          FDSZ TEMP
          FJMP LPI
    
```

Arithmetic Instructions (ALC)

Options

The floating point ALC instructions are similar to normal ALC instructions. Two floating accumulators must be specified. The first is the source accumulator, the second the destination accumulator.

Seven skip conditions are defined (in addition to the default "no skip"). These conditions are listed in Table 1. The conditions FSZR and FSNR should be used with caution. Since floating point arithmetic is inherently approximate, the probability of obtaining true zero is very low. The normal procedure is to test a result (or difference) against a small quantity, ϵ . For example, if we wish to test for the convergence of an iterative procedure we might use the following:

```

EPSLN:  1.0E-6          ;EPSILON IS 10**-5
          :
          :
          FLDA 0, ORSLT ;GET OLD RESULT
          FLDA 1, NRSLT ;GET NEW RESULT
          FSUB 1, 0     ;OLD-NEW
          FLDA 1,EPSLN ;EPSILON
          FPOS 0,0     ;ABS (OLD-NEW)
          FSUB 1,0,FSLE ;SKIP IF < EPSILON
    
```

Arithmetic Instructions (ALC) (cont'd)

Skip Mnemonic	Effect
FSLT	skip if result < 0
FSLE	skip if result ≤ 0
FSGT	skip if result > 0
FSGE	skip if result ≥ 0
FSNR	skip if result ≠ 0
FSZR	skip if result = 0
FSKP	unconditional skip

Table 1 - Skip Mnemonics

All ALC instructions permit the load/no load option. As with normal instructions, if a floating point instruction mnemonic is suffixed with "#", the results of the operation will not replace the contents of the destination register.

A further option is available with one class of ALC instructions. This option will prevent post-normalization of the result. The instructions described in the section following permit this option. It is called for by suffixing "U" (for unnormalized) to the instruction mnemonic. For example, FMOV 0,1 moves FAC0 to FAC1 and normalizes the result, while FMOVU 0,1 moves FAC0 to FAC1 without normalization.

ALC Instructions with Post-Normalize Option

The instruction

```
FMOV n,m
```

moves FAC_n to FAC_m.

The instruction

```
FPOS n,m
```

moves the absolute value of FAC_n to FAC_m.

The instruction

```
FMNS n,m
```

moves the negative of the absolute value of FAC_n to FAC_m.

ALC Instructions with Post-Normalize Option: (cont'd)

The instruction

FNEG n, m

moves the negative value of FAC_n to FAC_m.

The instruction

FRND n, m

rounds the value of FAC_n and moves it to FAC_m. By "round" we mean the following.

$$F = 16^{** -6} * \lfloor 16^{**6} * F + 1/2 \rfloor$$

The instruction

FADD n, m

adds FAC_n to FAC_m and moves the result to FAC_m. If underflow occurs, the underflow flag is set and true \emptyset is returned as the result. If overflow occurs, the overflow flag is set and a magnitude of $(1-16^{** -6}) * 16^{** +63}$ is returned as the result. The operands are assumed to be pre-normalized.

The instruction

FSUB n, m

subtracts FAC_n from FAC_m and moves the result to FAC_m. The overflow conditions are handled as with FADD. Prenormalized operands are assumed.

Table 2 summarizes the floating ALC instructions with post-normalize option.

Instruction	Effect
FMOV <u>n</u> , <u>m</u>	FAC _n → FAC _m
FPOS <u>n</u> , <u>m</u>	FAC _n → FAC _m
FMNS <u>n</u> , <u>n</u>	- FAC _n → FAC _m
FNEG <u>n</u> , <u>m</u>	- FAC _n → FAC _m
FRND <u>n</u> , <u>m</u>	rounded FAC _n → FAC _m
FADD <u>n</u> , <u>m</u>	FAC _n +FAC _m → FAC _m
FSUB <u>n</u> , <u>m</u>	FAC _m -FAC _n → FAC _m

Table 2 - ALC Instructions with Post-Normalize Option

ALC Instructions that Always Post-Normalize

This class of ALC instructions always post-normalizes the result. They assume pre-normalized operands. Overflow is checked and indicated by setting the overflow flag and returning a magnitude of $(1-16^{** -6}) * 16^{** +63}$ as the magnitude of the result. Underflow is checked and indicated by setting the underflow flag and returning true \emptyset as the result.

The instruction

FMPY n, m

multiplies FAC_n by FAC_m and moves the result to FAC_m.

The instruction

FDIV n, m

divides FAC_m by FAC_n and moves the result to FAC_m.

The instruction

FHLV n, m

halves FAC_n and moves the result to FAC_m.

Table 3 summarizes the ALC instructions that always post-normalize the result.

Instruction	Effect
FMPY <u>n</u> , <u>m</u>	FAC _m *FAC _n → FAC _m
FDIV <u>n</u> , <u>m</u>	FAC _m /FAC _n → FAC _m
FHLV <u>n</u> , <u>m</u>	FAC _n /2. → FAC _m

Table 3 - ALC Instructions that Always Post-Normalize

Floating ALC Instruction Examples

Some examples of legal floating point ALC instructions are:

FMPY 1, \emptyset
 FADD \emptyset , 1, FSGE
 FSUB# 1, \emptyset , FSLT
 FMOVU 3, \emptyset
 FMOV \emptyset , \emptyset , FSLT
 FNEG \emptyset , \emptyset , FSKP

\diamond $\lfloor x \rfloor$ gives the maximum integer K such that $K \leq x$.

Input/Output Instructions

The use of I/O instructions requires the user to provide two special routines. The first is an input routine which, when called, must return an ASCII input character, right justified in AC0 with bit 8 = 0. The address of this routine must be stored by the user in location 040 of page 0.

The second routine is an output routine which, when called, must accept an ASCII output character, right justified in AC0 with bit 8 = 0. The address of this routine must be stored in location 041 of page 0. All output messages to this routine will be terminated by a null (all zero) character.

These I/O routines must be reentrant for the interpreter to be reentrant. (If the routines which interrupt and use the interpreter do not use I/O instructions, the user I/O routines need not be reentrant.)

The instruction

FDFC n ;FLOATING POINT DECIMAL TO
;FLOATING CONVERT

will cause an ASCII character string in engineering notation to be converted to internal floating point form and loaded in FACn. The input characters must be provided by the user routine whose address is stored in location 408 of page 0. Numbers in the following form will be converted.

$$\left\{ \begin{array}{l} [+] \underline{n} \left[\dots \underline{n} [.] \underline{n} \dots \underline{n} \left[E \begin{array}{l} [+] \\ [-] \end{array} \underline{m} \left[\underline{m} \right] \right] \right] \begin{array}{l} b \\ r \\ e \\ a \\ k \end{array} \end{array} \right\}$$

where: n is the decimal mantissa (the first seven non-zero digits will be converted and the remaining digits ignored) and each m is a digit of the decimal characteristic.

The signs of the mantissa and characteristic are optional with the default assumed + .

The break character is any character other than

1. a decimal digit or
2. an "E", "+", "-", or "."

If the break character is a rubout (177), the entire string will be ignored and a new one must be given, i.e., the conversion starts over.

If the conversion results in a number less than $16^{*-1}16^{*-63}$, the underflow flag will be set and true zero will replace FACn. If the conversion results in a number whose magnitude is greater than

Input/Output Instructions (cont'd)

$(1-16^{*-6})16^{*-63}$, this latter magnitude will replace FACn and overflow will be set. As described on page 2-1, input conversion returns three additional words of information: conversion OK flag, the break character, and decimal point seen flag. These may be examined and used as necessary.

Examples of legal character strings are:

- 1 *
- 1. *
- 1 *
- +1 *
- 1E3 *
- 3.1415926 *
- 1.E+70 *
- 1.E-70 *

where * will be returned as the break character and conversion OK will be non-zero.

Some illegal character strings are:

- A (break character will be A)
- +* (break character will be *)
- +.! (break character will be !)

Conversion OK will be zero in all these illegal cases.

The instruction

FDFCI n ;FDFC WITH INDICATION

will provide the user with an indication before the conversion begins. The ASCII character "F" followed by a null character will be passed to the output routine whose address is given in location 041. For example, if the user has provided for I/O from the teletype, the use of FDFCI will print "F" on the page copy every time an input is required. In all other respects it is identical to FDFC.

The instruction

FFDC n ;FLOATING POINT FLOATING TO
;DECIMAL CONVERT

will convert the number in FACn to an ASCII character string in engineering notation. The output characters will be passed one at a time, right justified in AC0, to a user routine whose address is stored in location 418 of page 0. The output string will be of the form:

$$\left\{ \begin{array}{l} + \\ - \end{array} \right\} . \underline{n} \underline{n} \underline{n} \underline{n} \underline{n} \underline{n} \underline{n} E \left\{ \begin{array}{l} + \\ - \end{array} \right\} \underline{m} \underline{m}$$

Input/Output Instructions (cont'd)

where: each n represents a decimal digit of the mantissa.

each m represents a digit of the decimal characteristic.

The string will be terminated by a null character.

Special Instructions

Two special instructions are defined which modify the index registers.

The instruction

FIC2

causes AC2 to be incremented by two.

The instruction

FIC3

causes AC3 to be incremented by two. These instructions are useful for indexing through a table of floating point numbers. Use of FIC2 in indexing is shown in the second example in Chapter 5.

A third special instruction provides a HALT feature within the interpretive mode. The instruction

FHLT

will cause the interpreter to HALT. Hardware AC0 will contain the address of the FHLT instruction. The address lights will have no apparent relationship to the HALT, since the address is within the interpreter. The user may press CONTINUE to resume after this HALT.

Illegal Instructions

The proper encoding for all floating point instructions is given in Appendix B. The interpreter will HALT if an illegal instruction is encountered. Hardware AC0 will contain the address where the illegal instruction was found. This HALT will occur if extended instructions are used and only the Basic Interpreter is loaded, or on any bit configuration that cannot be decoded into a floating point instruction. The user cannot press CONTINUE to resume after this HALT.

REQUIREMENTS FOR REENTRANCE

A number of points regarding reentrance of the interpreter have been mentioned. This section explicitly defines the rules which must be followed by any routine which interrupts a base level routine and reenters the interpreter.

- 1) All hardware accumulators, Carry, and page 0 locations 006 and 007 must be saved.
- 2) A new writable area address must be provided in location 007.
- 3) An FINI must be issued after location 007 has been set up (only necessary the first time).
- 4) If I/O instructions are to be used, the user I/O routines must be reentrant. (Alternatively, locations 040 and 041 must be saved and addresses provided to different I/O routines).
- 5) Upon exit to the base level routine, the hardware accumulators, Carry and locations 006 and 007 must be restored.

FEATURES

The Extended Interpreter provides the instructions described in Chapter 2, plus a number of mathematical functions and "F" format output. If the Extended Interpreter is used, 1000 (decimal) words of writable storage must be provided by the user.

MATHEMATICAL FUNCTIONS

The mathematical functions are implemented using ALC instructions which always post-normalize. They permit the no load option as well as the floating skip options. Appendix C provides a detailed description of the methods used to implement these functions as well as a discussion of their accuracy. The following is a general description of each instruction.

The instruction

FALG n, m

computes the natural logarithm of the contents of FACn and moves the result to FACm. If the argument is less than 0, the overflow flag is set and $-(1-16^{*-6}) * 16^{*63}$ is returned as the result.

The instruction

FATN n, m

computes the arctangent of the contents of FACn and moves the result to FACm. The result is an angle expressed in radians in the range $-\pi/2 \leq \arctan(x) \leq \pi/2$.

The instruction

FCOS n, m

computes the cosine of the contents of FACn and moves the result to FACm. The argument is assumed to be an angle expressed in radians. If the argument is greater than 2^{*24} , the overflow flag is set, and the result will be incorrect.

The instruction

FSIN n, m

computes the sine of the contents of FACn and moves the result to FACm. The argument is assumed to be an angle expressed in radians. If the argument is greater than 2^{*24} , the overflow flag is set, and the result will be incorrect.

MATHEMATICAL FUNCTIONS (cont'd)

The instruction

FTAN n, m

computes the tangent of the contents of FACn and moves the result to FACm. The argument is assumed to be an angle expressed in radians. If the argument is greater than 2^{*24} , the overflow flag will be set, and the result will be incorrect.

The instruction

FEXP n, m

computes e raised to the power contained in FACn and moves the result to FACm. If the argument is less than -177.5, true 0 is returned, and the underflow flag is set. If the argument is greater than 174.673, $+(1-16^{*-6}) * 16^{*63}$ is returned, and the overflow flag is set.

The instruction

FSQR n, m

computes the square root of the argument in FACn and moves the result to FACm. If the argument is less than 0, the underflow flag is set, and true 0 is returned as the result.

Table 4 summarizes the math functions.

Instruction	Effect
FALG <u>n</u> , <u>m</u>	$\ln(\text{FACn}) \rightarrow \text{FACm}$
FATN <u>n</u> , <u>m</u>	$\arctan(\text{FACn}) \rightarrow \text{FACm}$
FCOS <u>n</u> , <u>m</u>	$\cosine(\text{FACn}) \rightarrow \text{FACm}$
FSIN <u>n</u> , <u>m</u>	$sine(\text{FACn}) \rightarrow \text{FACm}$
FTAN <u>n</u> , <u>m</u>	$tangent(\text{FACn}) \rightarrow \text{FACm}$
FEXP <u>n</u> , <u>m</u>	$e^{**}(\text{FACn}) \rightarrow \text{FACm}$
FSQR <u>n</u> , <u>m</u>	$(\text{FACn})^{*1/2} \rightarrow \text{FACm}$

Table 4 - Math Functions

"F" Format Conversion

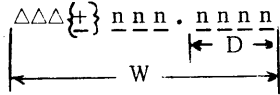
The Basic Interpreter provides floating point to decimal conversion with "E" format output. The extended version provides "F" format output as well.

The instruction

FFDCF n ;FFDC WITH "F" FORMAT OUTPUT

"F" Format Conversion (Cont'd)

will convert the floating point number in FAC_n to decimal and output a character string in "F" format via the user routine whose address is stored in location 041. The output will be of the form:



The width of the field (including sign and decimal point), W, and the number of places to be given after the decimal point, D, must be set up in the writable area before FFDCF is given. The displacements of these words are

W = 121₈
D = 122₈

They can be accessed in a manner similar to the flags described on page 2-1. Two conditions will cause the overflow flag to be set and no conversion to be performed,

"F" Format Conversion (cont'd)

1. W > 32 (entire width of field limited to 32 characters.)
2. W < D + 2 (W must be 2 greater than D to provide for sign and decimal point.)

If W is not large enough to accommodate the number, significant digits will be lost.

Assume W has been set to 12 (decimal) and D to 6.

Examples of "F" format versus "E" format outputs are:

"E" Format	"F" Format
+ .13746000E+02	+13.746000
- .7968433E-03	- .000796
+ 10000000E+04	+1000.000000
- 10000000E+05	10000.000000 (note sign lost)
+ 10000000E+06	00000.000000 (all lost)
- .35000000E-06	- .000000 (significance lost)
+ .4713279E-01	+ .047132

The Extended Relocatable Floating Point Interpreter is identical to the Extended Floating Point Interpreter, except for the following:

- 1) No absolute locations are used.
- 2) The interpreter requires 3 zero relocatable locations* and approximately 3500₈ normal relocatable locations.
- 3) Within any program that calls or initializes the interpreter, the appropriate normal external mnemonics must be declared. These are:

```
FENT      ;FLOATING INTERPRETER ENTER
FINT      ;FLOATING INTERPRETER
          ;INITIALIZE
```

FENT replaces the FETR command, used to call the absolute interpreter.

FINT replaces the FINI command, used to initialize the absolute interpreter.

For example:

```
.EXTN  FENT, FINT

.NREL

FINT      ;INITIALIZE
FENT      ;ENTER INTERPRETER
FLDA 0, FONE ;FLOATING INSTRUCTIONS
.
.
.
```

- 4) One ZREL location must be defined and declared as an entry (.ENT) with the label

WSA

Location WSA must contain a pointer to the writable area of the interpreter. WSA replaces location 007 of the absolute interpreter. (The writable area must be 10010 contiguous locations as in the absolute interpreter.)

- 5) If I/O instructions are to be used, two more ZREL locations must be defined and declared as entries with the labels:

```
GETC
PUTC
```

Location GETC must contain a pointer to the input character routine. GETC replaces location 040 of the absolute interpreter.

Location PUTC must contain a pointer to the output character routine. PUTC replaces location 041 of the absolute interpreter.

- 6) Requirements for Reentrance

- a) All hardware accumulators, Carry, and the contents of WSA and ITMP must be saved.
- b) A new writable area address must be provided in WSA.
- c) A FINI must be issued after location WSA has been set up. (This is necessary only once per writable area.)
- d) If I/O instructions are to be used, the user I/O routines must be reentrant (or locations GETC and PUTC must be saved and new addresses provided to other I/O routines.)
- e) Upon exit, all of the information saved above must be restored.

The Extended Relocatable Floating Point Interpreter can be loaded for stand-alone use or for use under the Disk Operating System. The relocatable version of the interpreter is sent in relocatable binary format, Tape #089-000046.

* Three ZREL locations are required in the interpreter, and three are provided by the user program (WSA, GETC, and PUTC).

SQUARE ROOT NEWTON ITERATION

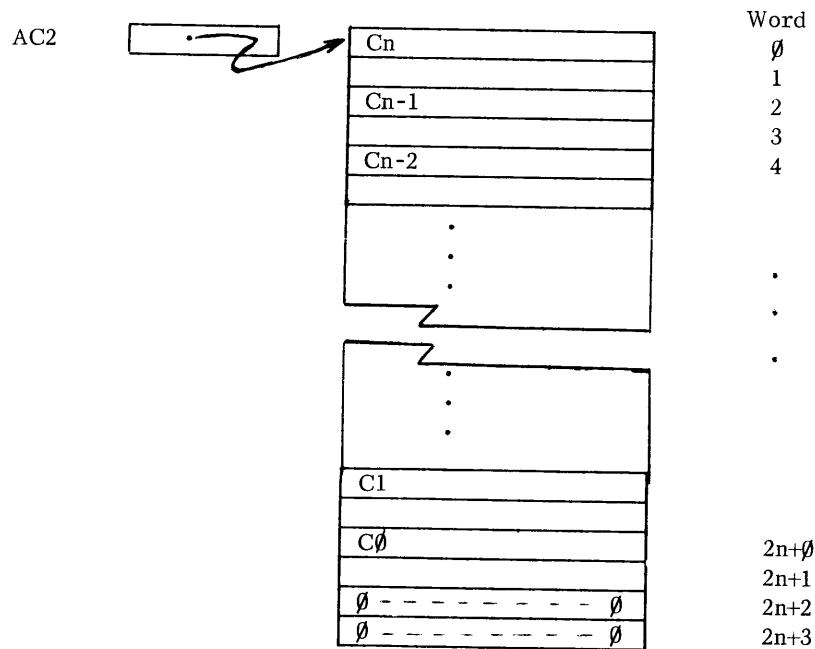
The following routine is an example of a floating point subroutine that performs a square root Newton iteration on the trial guess in FAC0 given the argument in FAC1.

```

;
;THIS ROUTINE PERFORMS A SQUARE ROOT NEWTON ITERATION
;FAC1 CONTAINS THE ARGUMENT AND IS NOT DESTROYED
;FAC0 CONTAINS THE PREVIOUS GUESSTIMATE
;COMPUTES (FAC0 + (FAC1/FAC0)) /2.
;CALLING SEQUENCE
;   FJSR NSR
;   RETURN : RESULT IN FAC0
;
NSR:   FMOV 1,2           ;SAVE ARGUMENT
       FDIV 0,2          ;FAC1/FAC0
       FADD 2,0          ;FAC0:FAC1/FAC0
       FHLV 0,0          ;(FAC0+FAC1/FAC0)/2.
       FJMP 0,3          ;RETURN
    
```

POLYNOMIAL EVALUATION

This polynomial evaluation routine requires AC2 to point to the first word of a table of floating point coefficients, ordered high order coefficient down and terminated by true 0. For example:



The routine uses Horner's method for evaluation, i. e.

$$f(x) = (\dots((x + C_n)x + C_{n-1})x + C_{n-2})x \dots + C_1)x + C_0$$

POLYNOMIAL EVALUATION (Continued)

```

;
;POLYNOMIAL EVALUATION
;FAC2 CONTAINS ARGUMENT (X)
;AC2 POINTS TO COEFFICIENT LIST TERMINATED BY 0,
;    AND ORDERED HIGH TO LOW.
;RESULT RETURNED IN FAC0
;FAC0, FAC1 DESTROYED
;CALLING SEQUENCE
;    FJSR          FPLY
;    RETURN
;

FZRO:  0
        0
FLPY:  FLDA          0, FZRO      ;CLEAR RESULT
FLPY1: FLDA          1, 0, 2      ;GET COEFFICIENT
        FMOV          1, 1, FSNR
        FJMP          0, 3        ;RETURN IF ZERO
        FMPY          2, 0        ;SUM * ARGUMENT
        FADD          1, 0        ;SUM * ARG. + COEF.
        FIC2          ;BUMP POINTER TO NEXT COEF.
        FJMP          FPLY1
    
```

Standard Instructions

FADD	Floating Add
FDFC	Floating Decimal to Floating Convert
FDFCI	FDFC with Indication
FDIV	Floating Divide
FDSZ	Floating Decrement and Skip if Zero
FETR	Floating Mode Enter
FEXT	Floating Mode Exit
FFDC	Floating Floating to Decimal Convert
FFIX	Floating to Fixed
FFLO	Fixed to Floating
FHLT	Floating Halt
FHLV	Floating Halve
FIC2	Floating Increment AC2
FIC3	Floating Increment AC3
FINI	Floating Initialize
FISZ	Floating Increment and Skip if Zero
FJMP	Floating Jump
FJSR	Floating Jump to Subroutine
FLD3	Floating Load AC3
FLDA	Floating Load Floating Accumulator
FMOV	Floating Move
FMNS	Floating Move Minus
FMPY	Floating Multiply
FNEG	Floating Negate
FPOS	Floating Move Positive
FRND	Floating Round
FST3	Floating Store AC3
FSTA	Floating Store Floating Accumulator
FSUB	Floating Subtract

Extended Instructions

FALG	Floating Natural Logarithm
FATN	Floating Arctangent
FCOS	Floating Cosine
FEXP	Floating Exponential
FFDCF	Floating Floating to Decimal Convert with "F" Format
FSIN	Floating Sine
FSQR	Floating Square Root
FTAN	Floating Tangent

Floating Point Options

#	No Load
FSGE	Floating Skip on Greater Than or Equal
FSGT	Floating Skip on Greater Than
FSKP	Floating Skip
FSLE	Floating Skip on Less Than or Equal
FSLT	Floating Skip on Less Than
FSNR	Floating Skip on Non Zero Result
FSZR	Floating Skip on Zero Result
U	Unnormalize (no post-normalization)

FLOATING POINT INTERPRETER

Appendix B - Floating Point Instruction Encoding

MEMORY REFERENCE WITHOUT ACCUMULATOR

INDIRECT															
0	FUNCTION				INDEX	DISPLACEMENT									
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0000	FJMP	1100	FFLO												
0001	FJSR	1101	FLD3												
0010	FISZ	1110	FST3												
0011	FDSZ	1111	FFIX												

MEMORY REFERENCE WITH ACCUMULATOR

INDIRECT															
0	FUNCTION	FAC ADDRESS	INDEX	DISPLACEMENT											
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
01	FLDA														
10	FSTA														

SPECIAL

1	0	0	FUNCTION	DISPLACEMENT											
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
00	FEXT														
01	FIC2														
01	FIC3														
11	FHLT														

CONVERSION

1	FUNCTION	FAC ADDRESS	DISPLACEMENT												TYPE
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
01	FDFC														FDFC
10	FFDC														1 I
															FFDC
															1 F

ARITHMETIC (OPTIONAL NORMALIZATION)

NO NORMALIZE															
1	FAC SOURCE ADDRESS	FAC DEST. ADDRESS	FUNCTION	0	0	0	NO LD.	SKIP							
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
			001	FNEG				001	FSGT						
			010	FMOV				010	FSLT						
			011	FPOS				011	FSNR						
			100	FMNS				100	FSZR						
			101	FSUB				101	FSGE						
			110	FADD				110	FSLE						
			111	FRND				111	FSKP						

ARITHMETIC (ALWAYS NORMALIZED)

1	FAC SOURCE ADDRESS	FAC DEST. ADDRESS	0	0	0	FUNCTION	NO LD.	SKIP							
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
			0001	FALG				001	FSGT						
			0010	FATN				010	FSLT						
			0011	FCOS				011	FSNR						
			0100	FMPY				100	FSZR						
			0101	FSIN				101	FSGE						
			0110	FTAN				110	FSLE						
			1000	FDIV				111	FSKP						
			1001	FEXP											
			1010	FSQR											
			1100	FHLV											

This appendix describes in some detail the function, method, and accuracy of the mathematical routines supplied with the Extended Interpreter. The accuracy of the routines is influenced by two factors:

1. The accuracy of the argument
2. The accuracy of the algorithm

These two factors will be mentioned for each routine. The accuracy of the algorithm itself assumes an argument that is exact, i. e., no argument error.

The relative and absolute errors of a function routine are defined as follows:

Let $f(x)$ = the true value of the function at x

Let $g(x)$ = the result returned by the function routine, given x

Now the absolute error of the result is,

$$\text{ABS}(f(x) - g(x))$$

and the relative error of the result is

$$\text{ABS}((f(x) - g(x))/f(x))$$

FLOATING POINT ARCTANGENT

Function:

To calculate the arctangent of x , where x is a floating point number, and return an angle in radians in the range $-\pi/2 \leq \arctan(x) \leq \pi/2$.

Method:

The range of x is immediately reduced to

$$0 \leq x \leq 1$$

by means of the identities

$$\arctan(-\text{abs}(x)) = -\arctan(\text{abs}(x))$$

and if $\text{abs}(x) > 1$,

$$\arctan(\text{abs}(x)) = \pi/2 - \arctan(1/\text{abs}(x))$$

For $x > \tan(\pi/12)$ the range is reduced to

$$\tan(-\pi/12) \leq y \leq \tan(\pi/12)$$

FLOATING POINT ARCTANGENT (Cont d)

by means of the identity

$$\arctan(x) = \arctan(x\emptyset) + \arctan[(x-x\emptyset)/(1-x*x\emptyset)],$$

For $x\emptyset = 1/3**\emptyset.5$, we obtain

$$\begin{aligned} \arctan(x) &= \pi/6 + \arctan[((x*3**\emptyset.5)-1)/(x+3**\emptyset.5)] \\ &= \pi/6 + \arctan(y) \end{aligned}$$

where y satisfies the range given above. The arctan is computed using the first four terms of a polynomial approximation of the form:

$$\arctan(x) \approx x * \sum_{i=0}^n C_i * x**(2*i)$$

Accuracy:

Argument Error

If x is the argument, the absolute error of the result is approximately

$$\epsilon/(1+x**2)$$

where ϵ is the absolute error in x . Thus for small values of x , the errors are almost equal, while as x becomes larger, the effect of the argument error decreases.

Maximum Relative Error

For the range

$$-\tan(\pi/12) \leq x \leq \tan(\pi/12),$$

the maximum relative error is approximately $10**-7.6$.

FLOATING POINT EXPONENTIAL

Function:

To calculate e to the power x , where x is a floating point number.

Method:

If $x < -177.5$, return \emptyset as the result and set underflow flag.

If $x > 174.673$, return the largest positive number as the result and set overflow flag.

FLOATING POINT EXPONENTIAL (cont d)

Otherwise, we use the equality

$$e^{**x} = 2^{**(\lfloor x \log_2(e) \rfloor)}$$

Let $x \log_2(e) = m \cdot 16^{**c}$

Further, let $q = \lfloor m \cdot 16^{**c} \rfloor$ and $f = (m \cdot 16^{**c}) - q$.
Therefore $0 \leq f < 1.0$.

Now $2^{**f} = 2^{**f} \cdot 2^{**q}$

Let us compute 2^{**f} .

The range of f can be further reduced if we let $g = f - 0.5$ where $-0.5 \leq g < 0.5$.

We compute 2^{**g} directly if $g \geq 0$, otherwise we compute 2^{**g} as

$$1/(2^{**\text{abs}(g)})$$

2^{**g} (where $0 \leq g < 0.5$) is computed using the first five terms of a polynomial approximation of the form:

$$2^{**g} \cong \sum_{i=0}^n C_i \cdot g^{**i}$$

Now $f = g + 0.5$

Therefore, $2^{**f} = 2^{**g} \cdot 2^{**0.5}$

But the answer is $2^{**f} \cdot 2^{**q}$

q is an integer and we let $q = 4i + j$

Now $2^{**f} \cdot 2^{**q} = 2^{**f} \cdot 16^{**i} \cdot 2^{**j}$

The characteristic of f is added to i to obtain $2^{**f} \cdot 16^{**i}$. This result is shifted j positions left if $i > 0$ or j positions right if $i < 0$. The result is, of course, e^{**x} .

Accuracy:

Argument Error

The relative error of the result is approximately equal to the absolute error of the argument. Thus for large values of x , substantial relative errors in the results can occur.

Maximum Relative Error

For the range $0 \leq x < 0.5$ the maximum relative error of 2^{**x} is approximately $10^{**-7.0}$.

FLOATING POINT NATURAL LOGARITHM

Function:

To calculate the natural logarithm of x , where x is a floating point number.

Method:

If $x < 0$, overflow flag is set and minus the largest floating point number is returned as the value. Otherwise, let $x = m \cdot 16^{**c}$. By means of a binary normalization, the range of m is reduced to

$$1/2 \leq m < 1,$$

and $x = m \cdot 16^{**p} \cdot 2^{**(-q)}$ where $q =$ number of left shifts required to normalize ($0 \leq q \leq 3$).

Now for $1/2 \leq m < 1/2 \cdot 16^{**0.5}$

let $a = 1/2$, $b = 1$

for $1/2 \cdot 16^{**0.5} \leq m < 1$

let $a = 1$, $b = 0$

Define $y = (m - a)/(m + a)$

then $m = a \cdot (y+1)/(-y+1)$

$$\begin{aligned} \text{Now } x &= 16^{**p} \cdot 2^{**(-q)} \cdot a \cdot (1+y)/(1-y) \\ &= 2^{**p} \cdot 2^{**(-q)} \cdot a \cdot (1+y)/(1-y) \end{aligned}$$

Using $\ln(x) = \ln(2) \cdot \log_2(x)$ we obtain

$$\ln(x) = \ln(2) \cdot [(4p-q-b) + \log_2((1+y)/(1-y))]$$

$$\ln(x) = \ln(2) \cdot (4p-q-b) + \ln((1+y)/(1-y))$$

From the above, we can determine that

$$2^{**-0.5} \leq (1+y)/(1-y) \leq 2^{**0.5}$$

The $\ln((1+y)/(1-y))$ is determined for the above range using the first three terms of a polynomial approximation of the form:

$$\ln(z) \cong z \cdot \sum_{i=0}^n C_i \cdot z^{**2 \cdot i}$$

where $z = (1+y)/(1-y)$.

FLOATING POINT NATURAL LOGARITHM (con t)

Accuracy:

Argument Error

The absolute error in the result is approximately equal to the relative error in the argument. Therefore, an argument close to 1 can give a large error since the function at this value is quite small.

Maximum Relative Error

For the range

$$1/2^{**-\theta}.5 \leq (1+x)/(1-x) \leq 2^{**\theta}.5$$

the maximum absolute error of ln(x) is approximately $10^{**-.7.6}$.

FLOATING POINT SINE AND COSINE

Function:

To calculate sin(x) or cos(x), where x is the floating point angle in radians.

Method:

Compute $p = \text{abs}(x) * 4 / \pi$

Let $q = \lfloor p \rfloor$, $f = p - q$ where $0 \leq f < 1$

Now q represents the half-quadrant in which the abs(x) falls. Using the following equalities.

$$\begin{aligned} \sin(x) &= -\sin(-x) \\ \cos(x) &= \cos(-x) \\ \cos(x) &= \sin(x + \pi/2) \end{aligned}$$

we define

$$\begin{aligned} q_1 &= q \text{ if sin is required and } x \geq 0 \\ q_1 &= q + 2 \text{ if cos is required} \\ q_1 &= q + 4 \text{ if sin is required and } x < 0 \end{aligned}$$

Then for all values of x, the computation has been reduced to

$$\sin(\pi/4 * (q_1 + f)) = \sin(t)$$

Since sine (and cosine) are periodic in $2 * \pi$, we take $q_1 = q_1 \text{ mod } 8$. Using the further equality

$$\sin(\pi/4 + x) = \cos(\pi/4 - x)$$

FLOATING POINT SINE AND COSINE . (cont d)

we finally can produce the table

q1	sin(t)
0	$\sin(\pi/4 * f)$
1	$\cos(\pi/4 * (1 - f))$
2	$\cos(\pi/4 * f)$
3	$\sin(\pi/4 * (1 - f))$
4	$-\sin(\pi/4 * f)$
5	$-\cos(\pi/4 * (1 - f))$
6	$-\cos(\pi/4 * f)$
7	$-\sin(\pi/4 * (1 - f))$

In all cases, the argument has been reduced to the range $0 \leq t < \pi / 4$. The sine is computed using the first four terms of a polynomial approximation of the form:

$$\sin(x) \approx x * \sum_{i=0}^n C_i * x^{** (2 * i)}$$

The cosine is computed using the first four terms of a polynomial approximation of the form:

$$\cos(x) \approx \sum_{i=0}^n C_i * x^{** (2 * i)}$$

Accuracy:

Argument Error

The absolute error of the result is approximately equal to the absolute error in the argument. Thus, the larger the argument, the larger the absolute error of the result.

Maximum Relative Error

For the range $0 \leq x < \pi/4$ the maximum relative error for sin(x) and cos(x) is approximately $10^{**-.7.4}$.

FLOATING POINT SQUARE ROOT

Function:

To calculate the square root of a floating point number, x.

Method:

Let $x = m * 16^{**c}$

FLOATING POINT SQUARE ROOT (cont'd)

If $m < \emptyset$, set underflow flag and return \emptyset . as the result.
 If $m = \emptyset$, return \emptyset . as the result.

Otherwise, let $c = 2 * p + q$ where p is an integer and $q = \emptyset$ or 1. Now if $q = \emptyset$, we have $x = m * 16 ** 2p$ and $x ** 1/2 = m ** 1/2 * 16 ** p$.

If $q = 1$, we have $x = m * 16 ** (2p+1)$
 $= m * 16 ** (2p+2) / 16$

and $x ** 1/2 = (m ** 1/2) / 4 * 16 ** (p+1)$

Therefore, the characteristic of the result is $p + q$, and the problem has been reduced to finding a suitable first guess for $m ** 1/2$ if $q = \emptyset$ or $(m ** 1/2) / 4$ if $q = 1$.

An initial guess is taken in the hyperbolic form $y\emptyset = a+b/(c+m)$

where for $q = \emptyset$,
 $a = 1.8\emptyset713$
 $b = -1.57727$
 $c = \emptyset.954182$

and for $q = 1$,
 $a = \emptyset.428795$
 $b = -\emptyset.343\emptyset368$
 $c = \emptyset.877552$

The initial guess is now

$$y1 = y\emptyset * 16 ** (p+q)$$

Two Newton iterations give us the result.

$$y2 = (y1+x/y1)/2$$

$$x ** 1/2 = y3 = (y2+x/y2)/2$$

Accuracy:

Argument Error

The relative error of the result is approximately half the relative error in the argument.

Maximum Relative Error

The maximum relative error for $x ** \emptyset.5$ is approximately $1\emptyset ** -6. \emptyset$.

FLOATING POINT TANGENT

Function:

To calculate the tangent of x , where x is the floating point angle in radians.

Method:

Compute $p = \text{abs}(x) * 4 / \pi$

Let $q = \lfloor p \rfloor$, $f = p - q$ where $\emptyset \leq f < 1$

Take $q_1 = q \bmod 4$ and in a manner similar to sine-cosine we can obtain the table below:

q_1	$\tan(x)$
\emptyset	$\tan(\pi/4 * f)$
1	$\cot(\pi/4 * (1-f))$
2	$-\cot(\pi/4 * f)$
3	$-\tan(\pi/4 * (1-f))$

In all cases, the argument has been reduced to the range $\emptyset \leq \text{arg} < \pi/4$.

The tangent is computed using the first six terms of a polynomial approximation of the form:

$$\tan(x) = x * \sum_{i=\emptyset}^n C_i * x ** (2 * i)$$

Accuracy:

Argument Error

The absolute error of the result is approximately equal to

$$\mathcal{E} * (1 + \tan(x) ** 2)$$

where \mathcal{E} is the absolute error of the argument. Thus if x is near an odd multiple of $\pi/2$, an argument error will produce a larger absolute error in the result.

Maximum Relative Error

For the range

$$\emptyset \leq x < \pi/4$$

the maximum relative error of $\tan(x)$ is approximately $1\emptyset ** -6.6$.

CHANGES FROM REVISION 3 TO REVISION 4 OF THE FLOATING POINT
INTERPRETER USER'S MANUAL

Page

- 1-1 Lower bound of the decimal range of a floating point number changed to 2.4×10^{-78} .
Spelling of during corrected.
- 2-1 Second note at bottom of page corrected.
- 2-4 The action taken by FMOV 0,1 is corrected to read "... FMOV 0,1 moves FAC0 to FAC1 ..."
- 2-7 Sentence on indexing example corrected.

- # option 2-4, 3-1
- ACCUMULATORS, floating point 2-1, 2-2, 2-3
- ARCTANGENT C-1, 3-1
- ARITHMETIC and logical instructions 2-4
- BASIC Interpreter Chapter 2
- BREAK character flag 2-1
- CHARACTER string conversion 2-6
- CONVERSION
 - character string 2-6
 - E format 3-2
 - F format 3-1
 - OK flag 2-1
- COSINE C-3, 3-1
- DECIMAL point seen flag 2-1
- E Format conversion 3-2
- EFFECTIVE address 2-3
- ENCODING of instructions App. B
- EXPONENTIAL C-1, 3-1
- EXTENDED Interpreter (binary) Chapter 3
- EXTENDED Interpreter (relocatable) Chapter 4
- F Format conversion 3-1
- FENT mnemonic 4-1
- FINT mnemonic 4-1
- FLAGS
 - break character 2-1
 - conversion OK 2-1
 - decimal point seen 2-1
 - overflow/underflow 2-1
- FSGE skip option 2-4
- FSGT skip option 2-4
- FSKP skip option 2-4
- FSLE skip option 2-4
- FSLT skip option 2-4
- FSNR skip option 2-4
- FSZR skip option 2-4
- GETC (get character) 4-1
- HALT 2-7
- INDEX register
 - modification 2-7
 - use in instructions 2-1
- INITIALIZATION 2-2
- INPUT/OUTPUT Instructions 2-6
- INSTRUCTIONS
 - arithmetic 2-4
 - illegal 2-7
 - formats in memory App. B
 - input/output 2-6
 - mathematical (extended) 3-1
 - memory reference 2-3
 - rules for 2-2
 - special 2-7
- INSTRUCTIONS, set of
 - FADD 2-5
 - FALG 3-1
 - FATN 3-1
 - FCOS 3-1
 - FDFC 2-6
 - FDFCI 2-5
 - FDIV 2-5
 - FDSZ 2-3
 - FETR 2-2, 2-3
 - FEXP 3-1
 - FEXT 2-2
 - FFDC 2-5
 - FFDCF 3-1
 - FFIX 2-3
 - FFLO 2-3
 - FHLT 2-7
 - FHLV 2-5
 - FIC2 2-7
 - FIC3 2-7
 - FINI 2-2, 2-7, 4-1
 - FISZ 2-3
 - FJMP 2-3
 - FJSR 2-3
 - FLD3 2-3
 - FLDA 2-1
 - FMOV 2-4
 - FMNS 2-4
 - FMPY 2-5
 - FNEG 2-5
 - FPOS 2-4
 - FRND 2-5
 - FSIN 3-1
 - FSQR 3-1
 - FST3 2-3
 - FSTA 2-3
 - FSUB 2-5
 - FTAN 3-1
- INTERPRETER
 - Basic 2-1, Chapter 2
 - Definition of i
 - entering the 2-2
 - exiting the 2-2
 - Extended 2-1, Chapter 3
 - Extended relocatable Chapter 4
 - general structure 2-1

FLOATING POINT INTERPRETER

Index

- INTERPRETIVE mode 2-2, 2-3
- ITMP mnemonic 4-1

- LOAD/no load option 2-4, 3-1

- MATHEMATICAL (extended) instructions 3-1, App. C
- MEMORY reference instructions 2-3

- NATURAL logarithm C-2, 3-1
- NORMAL mode 2-2
- NUMBER
 - conversion formats 3-1, 3-2
 - representation 1-1
 - rounding 2-5

- OVERFLOW/Underflow
 - flag 2-1
 - condition causing 2-1, 2-3, 2-5, 2-6, 3-1

- PAGE Zero locations used by interpreter 1-1
- POLYNOMIAL evaluation program 5-1
- POSTNORMALIZE option 2-4
- PUSH-DOWN list 1-1
- PUTC (put character) 4-1

- READ only storage (ROS) i, 1-1
- REENTRANCE
 - in I/O routines 2-6
 - linkage enabling 1-1
 - requirements for 2-7, 4-1
- ROUNDING a number 2-5

- SAMPLE programs Chapter 5
- SINE C-3, 3-1
- SKIP options 2-4, 3-1
- SPECIAL instructions 2-7
- SQUARE root C-4, 3-1
- SQUARE root Newton iteration program 5-1
- STARTING address
 - initialization 1-1
 - interpreter 1-1
 - GETC 1-1
 - PUTC 1-1
 - writable area 1-1, 2-1
- SUBROUTINE linkage 1-1, 2-3

- TANGENT C-4, 3-1
- TRIGONOMETRIC functions 3-1, App. C

- UNNORMALIZE (U) option 2-4

- WRITABLE area 1-1, 2-1, 4-1
- WSA (ZREL location) 4-1

DATA GENERAL CORPORATION
PROGRAMMING DOCUMENTATION
REMARKS FORM

DOCUMENT TITLE _____

DOCUMENT NUMBER (lower righthand corner of title page) _____

TAPE NUMBER (if applicable) _____

Specific Comments. List specific comments. Reference page numbers when applicable. Label each comment as an addition, deletion, change or error if applicable.

cut along dotted line

General Comments and Suggestions for Improvement of the Publication.

FROM: Name: _____ Date: _____
Title: _____
Company: _____
Address: _____

FOLD DOWN

FIRST

FOLD DOWN

FIRST
CLASS
PERMIT
No. 26
Southboro
Mass. 01772

BUSINESS REPLY MAIL

No Postage Necessary If Mailed In The United States

Postage will be paid by:

Data General Corporation

Southboro, Massachusetts 01772

ATTENTION: Programming Documentation

FOLD UP

SECOND

FOLD UP

STAPLE