

# The MINIFP package\*

Dan Luecking

2013/05/28

## Abstract

This package provides minimal fixed point exact decimal arithmetic operations. ‘Minimal’ means numbers are limited to eight digits on either side of the decimal point. ‘Exact’ means that when a number *can* be represented exactly within those limits, it will be.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>User macros</b>	<b>4</b>
2.1	Nonstack-based operations . . . . .	5
2.2	Commands to process numbers for printing . . . . .	8
2.3	Stack-based macros . . . . .	8
2.4	Errors . . . . .	12
<b>3</b>	<b>Implementation</b>	<b>12</b>
3.1	Utility macros . . . . .	12
3.2	Processing numbers and the stack . . . . .	14
3.3	The user-level operations . . . . .	19
3.4	The internal computations . . . . .	23
3.5	Commands to format for printing . . . . .	35
3.6	Miscellaneous . . . . .	38
<b>4</b>	<b>Extras</b>	<b>39</b>
4.1	Loading the extras . . . . .	40
4.2	Error messages . . . . .	41
4.3	Sine and Cosine . . . . .	42
4.4	Polar angle . . . . .	48
4.5	Logarithms . . . . .	51
4.6	Powers . . . . .	54
4.7	The square root . . . . .	60

---

\*This file has version number v0.95., last revised 2013/05/28. The code described here was developed by Dan Luecking.

# 1 Introduction

In working on an application that needed to be able to automatically generate numeric labels on the axes of a graph, I needed to be able to make simple calculations with real numbers. What  $\text{\TeX}$  provides is far too limited. In fact, its only native user-level support for real numbers is as factors for dimensions. For example one can “multiply”  $3.1 \times 0.2$  by `\dimen0=0.2pt \dimen0=3.1\dimen0` .

Unfortunately  $\text{\TeX}$  stores dimensions as integer multiples of the “scaled points” (`sp`) with `sp` =  $2^{-16}\text{pt}$ , and therefore `.2pt` is approximated by  $\frac{13107}{65536}$ , which is not exact. Then multiplying by 3.1 produces  $\frac{40631}{65536}$ . If we ask  $\text{\TeX}$  to display this, it produces `0.61998pt` and not the exact value 0.62. This is sufficiently accurate for positioning elements on a page, but not for displaying automatically computed axis labels if five digit accuracy is needed.

The MINIFP package was written to provide the necessary calculations with the necessary accuracy for this application. The implementation would have been an order of magnitude smaller and faster if only four digit accuracy were provided (and I may eventually do that for the application under consideration), but I have decided to clean up what I have produced and release it as is. The full MINIFP package provides nearly the same operations as a subset of the FP package, but the latter carries calculations to 18 decimal places, which is far more than necessary for my purposes. I want something small and fast to embed in the MFPIC drawing package.

I decided on eight digits on both sides of the decimal point essentially because I wanted at least five digits and the design I chose made multiples of four the easiest to work with.

MINIFP also provides a simple stack-based language for writing assembly language-like programs. Originally, this was to be the native calculation method, but it turned out to be too unwieldy for ordinary use. I left it in because it adds only about 10% overhead to the code.

But why *only* eight digits?  $\text{\TeX}$  only works with integers, and since the maximum integer allowed is about 2 000 000 000, the largest numbers that can be added are limited to about 999 999 999. It is very little trouble to add numbers by adding their fractional parts and integer parts separately as 9-digit integers. So it would seem multiples of 9 digits would be easy to implement.

However, something we have to do repeatedly in *division* is multiply the integer and fractional parts of a number by a one-digit number. For that purpose, nine digits would be too much, but eight digits is just right. For nine digits, we would have to inconveniently break the number into more than two parts. Limiting our numbers to eight-digit parts drastically simplifies division.

Another simplification: multiplication has to be done by breaking the number into parts.  $\text{\TeX}$  can multiply any two 4-digit integers without overflow, but it cannot multiply most pairs of 5-digit integers. Two 8-digit numbers conveniently break into four 4-digit parts. To get even nine digits of accuracy would require six parts (five, if we don’t insist on a separation occurring at the decimal point). The complexity of the multiplication process goes up as the square of the number of parts, so six parts would more than double the complexity.

A final simplification:  $\text{\TeX}$  places a limit of nine on the number of arguments a macro can have. Quite often the last argument is needed to clear out unused text to be discarded. Thus, a string of eight digits can quite often be processed with one execution of one nine-argument macro.

Addition and subtraction can be exact, multiplication and division can extend numbers past the 8-digit limit so they might be rounded. However, when the exact answer fits in the 8-digit limit, our code should produce it. Overflow (in the sense that the integer part can exceed the allowed eight digits) is always possible, but is much more likely with multiplication and division.

Multiplication is carried out internally to an exact answer, with 16 digits on each side of the decimal point. The underflow digits (places 9 through 16 after the decimal point) are used to round to an 8-digit result. Overflow digits (those to the left of the lowest 8 in the integer part) are discarded, usually without warning. Division is internally carried to nine digits after the decimal, which is then also rounded to an 8-digit result. Overflow digits are ignored for division also.

We supply two kinds of operations in this package. There are stack-based operations, in which the operands are *popped* from a stack and the results *pushed* onto it, and argument-based, in which the operands (and a macro to hold the result<sup>1</sup>) are arguments of a macro. Both types load the arguments into internal macros (think of them as “registers”), then call internal commands (think “microcode”) which return the results in internal macros. These results are then *pushed* onto the stack (stack-based operations) or stored in a supplied macro argument (think “variable”). The difference lies entirely in where the operands come from (arguments or stack) and where they go (macro or stack).

The stack is implemented as an internal macro which is redefined with each command. The binary operations act on the last two *pushed* objects in the order they were *pushed*. For example, the sequence “*push 5, push 3, subtract*” performs  $5 - 3$  by popping 3 and 5 into registers (thereby removing them from the stack), subtracting them and then pushing the result (2) onto the stack.

Our implementation of the *push* operation first prepares the number in a standard form. Thus, stack-based operations always obtain numbers in this form. The argument based operations will prepare the arguments in the same way. The internal commands will thus have a standard form to operate on. All results are returned in standard form.

The standard form referred to above is an integer part (one to eight digits with no unnecessary leading zeros nor unnecessary sign) followed by the decimal point (always a dot, which is ASCII 46), followed by exactly eight digits, all of this preceded by a minus sign if the number is negative. Thus,  $-0.25$  would be processed and stored as “0.25000000” and  $-0.333333$  as “-0.33333300”.

---

<sup>1</sup>Unlike most other packages for decimal arithmetic, MINIFP puts the macro to hold the result last. This allows the calculation to be performed before the macro is even read, and this makes it somewhat easier for the stack- and argument-based versions to share code.

## 2 User macros

MINIFP provides (so far) six binary operations (that act on a pair of numbers): addition, subtraction, multiplication, division, maximum and minimum, as well as fourteen unary operations (that act on one number): negation, absolute value, doubling, halving, integer part, fractional part, floor, ceiling, signum, squaring, increment, decrement and inversion. With the “extra” option, the unary operations sine, cosine, logarithm, powers and square root are available, and the binary operation angle. See section 4.

These extra operations are made available using the `extra` option in  $\LaTeX$ :

```
\usepackage[extra]{minifp}
```

In plain  $\TeX$ , they will be loaded if you give the macro `\MFPextra` a definition (any definition) before inputting `minifp.sty`:

```
\def\MFPextra{} \input minifp.sty
```

The extras can also be loaded by means of the command `\MFPlloadextra`, issued after `minifp.sty` is loaded. As of version 0.95 `mfpxextra` can be directly `\input`. It will detect whether `minifp.sty` has been loaded and input it if not. This will work only in plain  $\TeX$ .

If the extra operations are not needed, some memory and time might be saved by using `minifp.sty` alone. I have not seriously tried to keep `mfpxextra.tex` as small or fast as possible, but I do try to improve the accuracy when I can.

As previously mentioned, each of these operations come in two versions: a version that acts on operands and stores the result in a macro, and a version that acts on the stack. The former all have names that begin `\MFP` and the latter begin with `\R`. The former can be used anywhere, while the latter can only be used in a “program”. A program is started with `\startMFPprogram` and terminated with `\stopMFPprogram`. The `R` in the names is for ‘real’. This is because it is possible that stacks of other types will be implemented in the future.

For example, `\MFPadd{1.2}{3.4}\X` will add 1.20000000 to 3.40000000 and then define `\X` to be the resulting 4.60000000. These operand forms do not alter or even address the stack in any way. The stack-based version of the same operation would look like the following:

```
\Rpush{1.2}
\Rpush{3.4}
\Radd
\Rpop\X
```

which would *push* first 1.20000000 then 3.40000000 onto the stack, then replace them with 4.60000000, then remove that and store it in `\X`. Clearly the stack is intended for calculations that produce a lot of intermediate values and only the final result needs to be stored.

The command `\startMFPprogram` starts a group. That group should be ended by `\stopMFPprogram`. Changes to the stack and defined macros are local to that group. Thus the macro `\X` in the example above might seem to be useful only as a temporary storage for later calculations in the same program group. However, there are commands provided to force such a macro to survive the group, and

even to force the contents of the stack to survive the group (see the end of subsection 2.3). Do not try to turn a MINIFP program into a  $\text{\LaTeX}$  environment. The extra grouping added by environments would defeat the effects of these commands.

## 2.1 Nonstack-based operations

In the following tables, an argument designated  $\langle num \rangle$  can be any decimal real number with at most 8 digits on each side of the decimal point, or it can be a macro that contains such a number. If the decimal dot is absent, the fractional part will be taken to be zero, if the integer part or the fractional part is absent, it will be taken to be zero. (One consequence of these rules is that all the following arguments produce the same internal representation of zero:  $\{0.0\}$ ,  $\{0.\}$ ,  $\{.0\}$ ,  $\{0\}$ ,  $\{.\}$ , and  $\{ \}$ .) Spaces may appear anywhere in the  $\langle num \rangle$  arguments and are stripped out before the number is used. For example,  $\{3 . 1415 9265\}$  is a valid argument. Commas are not permitted. The decimal dot (period, fullstop) character must be inactivated if some babel language makes it a shorthand.

The  $\backslash macro$  argument is any legal macro. The result of using one of these commands is that the macro is defined (or redefined, there is no checking done) to contain the standard form of the result. If the  $\langle num \rangle$  is a macro, the braces surrounding it are optional.

### *Binary Operations*

Command	operation
$\backslash MFPadd\{\langle num_1 \rangle\}\{\langle num_2 \rangle\}\backslash macro$	Stores the result of $\langle num_1 \rangle + \langle num_2 \rangle$ in $\backslash macro$
$\backslash MFPsub\{\langle num_1 \rangle\}\{\langle num_2 \rangle\}\backslash macro$	Stores the result of $\langle num_1 \rangle - \langle num_2 \rangle$ in $\backslash macro$
$\backslash MFPmul\{\langle num_1 \rangle\}\{\langle num_2 \rangle\}\backslash macro$	Stores the result of $\langle num_1 \rangle \times \langle num_2 \rangle$ , rounded to 8 places after the decimal point, in $\backslash macro$
$\backslash MFPmpy\{\langle num_1 \rangle\}\{\langle num_2 \rangle\}\backslash macro$	Same as $\backslash MFPmul$
$\backslash MFPdiv\{\langle num_1 \rangle\}\{\langle num_2 \rangle\}\backslash macro$	Stores the result of $\langle num_1 \rangle / \langle num_2 \rangle$ , rounded to 8 places after the decimal point, in $\backslash macro$
$\backslash MFPmin\{\langle num_1 \rangle\}\{\langle num_2 \rangle\}\backslash macro$	Stores the smaller of $\langle num_1 \rangle$ and $\langle num_2 \rangle$ in $\backslash macro$
$\backslash MFPmax\{\langle num_1 \rangle\}\{\langle num_2 \rangle\}\backslash macro$	Stores the larger of $\langle num_1 \rangle$ and $\langle num_2 \rangle$ in $\backslash macro$

### Unary Operations

Command	operation
<code>\MFPchs{⟨num⟩}\macro</code>	Stores $-\langle num \rangle$ in <code>\macro</code> .
<code>\MFPabs{⟨num⟩}\macro</code>	Stores $ \langle num \rangle $ in <code>\macro</code> .
<code>\MFPdbl{⟨num⟩}\macro</code>	Stores $2 \times \langle num \rangle$ in <code>\macro</code> .
<code>\MFPhalve{⟨num⟩}\macro</code>	Stores $\langle num \rangle/2$ , rounded to 8 places after the decimal point, in <code>\macro</code> .
<code>\MFPint{⟨num⟩}\macro</code>	Replaces the part of $\langle num \rangle$ after the decimal point with zeros (keeps the sign unless the result is zero) and stores the result in <code>\macro</code> .
<code>\MFPfrac{⟨num⟩}\macro</code>	Replaces the part of $\langle num \rangle$ before the decimal point with zero (keeps the sign unless the result is zero) and stores the result in <code>\macro</code> .
<code>\MFPfloor{⟨num⟩}\macro</code>	Stores the largest integer not more than $\langle num \rangle$ in <code>\macro</code> .
<code>\MFPceil{⟨num⟩}\macro</code>	Stores the smallest integer not less than $\langle num \rangle$ in <code>\macro</code> .
<code>\MFPsgn{⟨num⟩}\macro</code>	Stores $-1$ , $0$ or $1$ (in standard form) in <code>\macro</code> according to whether $\langle num \rangle$ is negative, zero, or positive.
<code>\MFPsq{⟨num⟩}\macro</code>	Stores the square of $\langle num \rangle$ in <code>\macro</code> .
<code>\MFPinv{⟨num⟩}\macro</code>	Stores $1/\langle num \rangle$ , rounded to 8 places after the decimal point, in <code>\macro</code> .
<code>\MFPincr{⟨num⟩}\macro</code>	Stores $\langle num \rangle + 1$ in <code>\macro</code> .
<code>\MFPdecr{⟨num⟩}\macro</code>	Stores $\langle num \rangle - 1$ in <code>\macro</code> .
<code>\MFPzero{⟨num⟩}\macro</code>	Ignores $\langle num \rangle$ and stores $0.00000000$ in the <code>\macro</code> .
<code>\MFPstore{⟨num⟩}\macro</code>	Stores the $\langle num \rangle$ , converted to standard form, in <code>\macro</code>

The command `\MFPzero` is useful for “macro programs”. If you want to do something to a number depending on the outcome of a test, you may occasionally want to simply absorb the number and output a default result. This is more efficient than multiplying by zero (but less efficient than simply defining the `\macro` to be zero.)

Note that one could easily double, halve, square, increment, decrement or invert a  $\langle num \rangle$  using the binary versions of `\MFPadd`, `\MFPsub`, `\MFPmul` or `\MFPdiv`. The commands `\MFPdbl`, `\MFPhalve`, `\MFPsq`, `\MFPincr`, `\MFPdecr` and `\MFPinv` are designed to be more efficient versions, since they are used repeatedly in internal code.

Also, multiplication is far more efficient than division, so even if you use the two argument versions, `\MFPmul{⟨num⟩}{.5}` is faster than `\MFPdiv{⟨num⟩}{2}`.

There is one command that takes no argument and returns no value:

*Do Nothing*

Command	operation
<code>\MFPnoop</code>	Does nothing.

The following are not commands at all, but macros that contain convenient constants.

### Constants

Constant	value
<code>\MFPpi</code>	3.14159265, the eight-digit approximation to $\pi$ .
<code>\MFPe</code>	2.71828183, the eight-digit approximation to $e$ .
<code>\MFPphi</code>	1.61803399, the eight-digit approximation to the golden ratio $\phi$ .

There also exist commands to check the sign of a number and the relative size of two numbers:

```
\MFPchk{<num>}
\MFPcmp{<num1>}{<num2>}
```

These influence the behavior of six commands:

```
\IFneg{<true text>}{<false text>}
\IFzero{<true text>}{<false text>}
\IFpos{<true text>}{<false text>}
\IFlt{<true text>}{<false text>}
\IFeq{<true text>}{<false text>}
\IFgt{<true text>}{<false text>}
```

Issuing `\MFPchk{\X}` will check the sign of the number stored in the macro `\X`. Then `\IFneg{A}{B}` will produce ‘A’ if it is negative and ‘B’ if it is zero or positive. Similarly, `\MFPcmp{\X}{1}` will compare the number stored in `\X` to 1. Afterward, `\IFlt{A}{B}` will produce ‘A’ if `\X` is less than 1 and ‘B’ if `\X` is equal to or greater than 1.

If users finds it tiresome to type two separate commands, they can easily define a single command that both checks a value and runs `\IF...`. For example

```
\def\IFisneg#1{\MFPchk{#1}\IFneg}
```

Used like

```
\IFisneg{\X}{A}{B}
```

this will check the value of `\X` and run `\IFneg` on the pair of alternatives that follow.

The user might never need to use `\MFPchk` because every one of the operators provided by the MINIFP package runs an internal version of `\MFPchk` on the result of the operation before storing it in the `\macro`. For example, after `\MFPzero` the command `\IFzero` will always return the first argument. For this reason one should not insert any MINIFP operations between a check/compare and the `\IF...` command that depends on it.

Note: the behavior of all six `\IF...` commands is influenced by *both* `\MFPchk` and `\MFPcmp`. This is because internally `\MFPchk{\X}` (for example) and `\MFPcmp{\X}{0}` do essentially the same thing. In fact there are only three internal booleans that govern the behavior of the six `\IF...` commands. The different names are for clarity: `\IFgt` after a compare is less confusing than the entirely equivalent `\IFpos`.

It should probably be pointed out that the settings for the `\IF...` macros are local to any  $\text{\TeX}$  groups they are contained in.

## 2.2 Commands to process numbers for printing

After `\MFAdd{1}{2}\X` one can use `\X` anywhere and get 3.00000000. One might may well prefer 3.0, and so commands are provided to truncate a number or round it to some number of decimal places. Note: these are provided for printing and they will not invoke the above `\MFPchk`. They do not have any stack-based versions. The commands are

```
\MFPtruncate{<int>}{<num>}\macro
\MFPround{<int>}{<num>}\macro
\MFPstrip{<num>}\macro
```

where `<int>` is a whole number between  $-8$  and  $8$  (inclusive). The other two arguments are as before.

These commands merely process `<num>` and define `\macro` to produce a truncated or rounded version, or one stripped of trailing zeros, or one with added trailing zeros. Note that truncating or rounding a number to a number of digits greater than it already has will actually lengthen it with added zeros. For example, `\MFPround{4}{3.14159}\X` will cause `\X` to be defined to contain 3.1416, while `\MFPround{6}{3.14159}\X` will cause `\X` to contain 3.141590. If `\Y` contains 3.14159, then `\MFPtruncate{4}\Y\Y` will redefine `\Y` to contain 3.1415. Also, `\MFPstrip{1.20000000}\Z` will cause `\Z` to contain 1.2. All these commands first normalize the `<num>`, so any spaces are removed and redundant signs are discarded.

If `<int>` is negative, places are counted to the left of the decimal point and 0s are substituted for lower order digits. That is, `\MFPtruncate{-2}{1864.3}\X` will give `\X` the value 1800 and `\MFPround{-2}{1864}\X` will give `\X` the value 1900.

If the first argument of `\MFPround` or `\MFPtruncate` is zero or negative then the dot is also omitted from the result. If `\MFPstrip` is applied to a number with all zeros after the dot, then one 0 is retained. There is a star form where the dot and the zero are dropped.

For these three commands, the sign of the number is irrelevant. That is, the results for negative numbers are the negatives of the results for the absolute values. The processing will remove redundant signs along with redundant leading zeros: `\MFPtruncate{-3}{-+123.456}` will produce 0. The rounding rule is as follows: round up if the digit to the right of the rounding point is 5 or more, round down if the digit is 4 or less.

## 2.3 Stack-based macros

The stack-based macros can only be used in a MINIFP program group. This group is started by the command `\startMFPprogram` and ended by `\stopMFPprogram`. None of the stack-based macros takes an argument, but merely operate on values on the stack, replacing them with the results. There are also commands to manipulate the stack and save a value on the stack into a macro. Finally, since all changes to the stack (and to macros) are local and therefore lost after `\stopMFPprogram`, there are commands to selectively cause them to be retained.

To place numbers on the stack we have `\Rpush` and to get them off we have



`\Rpop`. The syntax is

`\Rpush{ $\langle num \rangle$ }`

`\Rpop\macro`

The first will preprocess the  $\langle num \rangle$  as previously discussed and put it on the stack, the second will remove the last number from the stack and define the given macro to have that number as its definition.

All the binary operations remove the last two numbers from the stack, operate on them in the order they were put on the stack, and *push* the result on the stack. Thus the program

`\Rpush{1.2}`

`\Rpush{3.4}`

`\Rsub`

will first put 1.20000000 and 3.40000000 on the stack and then replace them with -2.20000000. Note the order:  $1.2 - 3.4$ .

#### *Binary Operations*

<b>Command</b>	<b>operation</b>
<code>\Radd</code>	Adds the last two numbers on the stack.
<code>\Rsub</code>	Subtracts the last two numbers on the stack.
<code>\Rmul</code>	Multiplies the last two numbers on the stack, rounding to 8 decimal places.
<code>\Rmpy</code>	Same as <code>\Rmul</code> .
<code>\Rdiv</code>	Divides the last two numbers on the stack, rounding to 8 decimal places.
<code>\Rmin</code>	Replaces the last two elements on the stack with the smaller one.
<code>\Rmax</code>	Replaces the last two elements on the stack with the larger one.

The unary operations replace the last number on the stack with the result of the operation performed on it.

### *Unary Operations*

<b>Command</b>	<b>operation</b>
<code>\Rchs</code>	Changes the sign.
<code>\Rabs</code>	Obtains the absolute value.
<code>\Rdbl</code>	Doubles the value.
<code>\Rhalve</code>	Halves the value, rounding to 8 places.
<code>\Rint</code>	Replaces the fractional part with zeros. If the result equals 0.0, any negative sign will be dropped.
<code>\Rfrac</code>	Replaces the integer part with 0. If the result equals 0.0, any negative sign will be dropped.
<code>\Rfloor</code>	Obtains the largest integer not greater than the number.
<code>\Rceil</code>	Obtains the smallest integer not less than the number.
<code>\Rsgn</code>	Obtains $-1$ , $0$ or $1$ according to whether the number is negative, zero, or positive. These numbers are pushed onto the stack with the usual decimal point followed by 8 zeros.
<code>\Rsqr</code>	Obtains the square. Slightly more efficient than the equivalent <code>\Rdup\Rmul</code> . See below for <code>\Rdup</code> .
<code>\Rinv</code>	Obtains the reciprocal. Slightly more efficient than the equivalent division.
<code>\Rincr</code>	Increases by 1. Slightly more efficient than the equivalent addition.
<code>\Rdecr</code>	Decreases by 1. Slightly more efficient than the equivalent subtraction.
<code>\Rzero</code>	Replaces the number with zero. Slightly more convenient than the equivalent <code>\Rpop\X</code> followed by a <code>\Rpush{0}</code> .

There is one operation, which does not read the stack nor change it (nor do anything else).

### *Do Nothing*

<b>Command</b>	<b>operation</b>
<code>\Rnoop</code>	Does nothing.

There also exist commands to check the sign of the last number, and the relative size of the last two numbers on the stack:

`\Rchk`

`\Rcmp`

They do not remove anything from the stack. Just like the nonstack counterparts, they influence the behavior of six commands: `\IFneg`, `\IFzero`, `\IFpos`, `\IFlt`, `\IFeq` and `\IFgt`. Issuing `\Rchk` will check the sign of the last number on the stack, while `\Rcmp` will compare the last two in the order they were pushed. For example, in the following

```

\Rpush{1.3}
\Rpush{-2.3}
\Rcmp
\IFgt{\Radd}{\Rsub}
\Rpush\X

```

`\Rchk`  
`\IFneg{\Radd}{\Rsub}`

`\Rcmp` will compare 1.3 to  $-2.3$ . Since the first is greater than the second, `\IFgt` will be true and they will be added, producing  $-1.0$ . Following this the contents of the macro `\X` are pushed, it is examined by `\Rchk` and then either added to or subtracted from  $-1.0$ .

The user might never need to use `\Rchk` because every operator that puts something on the stack also runs `\Rchk`. In the above program, in fact, `\Rchk` is redundant since `\Rpush` will already have run it on the contents of `\X`.

There exist stack manipulation commands that allow the contents of the stack to be changed without performing any operations. These are really just conveniences, as their effects could be obtained with appropriate combinations of `\Rpop` and `\Rpush`. These commands, however, do not run `\Rchk` as `\Rpush` would.

#### *Stack Manipulations*

Command	operation
<code>\Rdup</code>	Puts another copy of the last element of the stack onto the stack.
<code>\Rexch</code>	Exchanges the last two elements on the stack.

After `\stopMFPprogram`, any changes to macros or to the stack are lost, unless arrangements have been made to save them. There are four commands provided. Two act on a macro which is the only argument, the other two have no arguments and act on the stack. The macro must simply contain a value, it cannot be more complicated and certainly cannot take an argument.

#### *Exporting changed values*

Command	operation
<code>\Export\macro</code>	Causes the definition of <code>\macro</code> to survive the program group.
<code>\Global\macro</code>	Causes the definition of <code>\macro</code> to be global.
<code>\ExportStack</code>	Causes the contents of the stack to survive the program group.
<code>\GlobalStack</code>	Causes the contents of the stack to be global.

The difference between `\Export` and `\Global` is solely in how *other* grouping is handled. If the program group is contained in another group (for example, inside an environment), then the result of `\Global\X` is that the definition of `\X` survives that group (and all containing groups) as well. On the other hand, after `\Export\X`, then the definition survives the program group, but not other containing groups.

If `TEX` grouping is used *inside* a program group, then using `\Export` inside that group has no effect at all, while `\Global` works as before.

The stack versions are implemented by running `\Export` or `\Global` on the internal macro that defines the stack, so they have the same behavior.

## 2.4 Errors

If one tries to *pop* from an empty stack, an error message will be issued. Ignoring the error causes the macro to have the value stored in the macro `\EndofStack`. Its default is 0.00000000.

If one tries to divide by zero, an error message will be issued. Ignoring the error causes the result to be one of the following:

- Dividing 0 by 0 gives a result whose integer part is stored in `\ZeroOverZeroInt` and whose fractional part is stored in `\ZeroOverZeroFrac`. The default is 0.00000000
- Dividing a nonzero  $x$  by 0 gives a result whose integer part is stored in `\xOverZeroInt` and whose fractional part is stored in `\xOverZeroFrac`. The defaults are both equal to 99999999. The sign of the result will be the sign of  $x$ .

You can change any of these macros, but make sure they produce a result which is a number in standard form (as described earlier). These macros are copied directly into the result without checking. Then further processing steps may require the result to be a number in standard form.

Error messages may result from trying to process numbers given in incorrect format. However, there are so many ways for numbers to be incorrect that this package does not even try to detect them. Thus, they will only be caught if some  $\text{\LaTeX}$  operation encounters something it cannot handle. (The  $\text{\LaTeX}$  manual calls these “weird errors” because the messages tend to be uninformative.) Incorrect numbers may even pass unnoticed, but leave random printed characters on the paper, or odd spacing.

## 3 Implementation

### 3.1 Utility macros

We announce ourself, and our purpose. We save the catcode of `@` and change it to letter. Several other catcodes are saved and set to other in this file. We also make provisions to load the extra definitions, either directly with `\MFPloadextra` or through a declared option in  $\text{\LaTeX}$ .

```
1 <*sty>
2 \expandafter
3 \ifx \csname MFP@finish\endcsname\relax
4 \else \expandafter\endinput \fi
5 \expandafter\edef\csname MFP@finish\endcsname{%
6   \catcode64=\the\catcode64 \space
7   \catcode46=\the\catcode46 \space
8   \catcode60=\the\catcode60 \space
9   \catcode62=\the\catcode62 \space}%
10 \ifx\ProvidesPackage\UndEfInEd
11   \newlinechar'\^^J%
12   \message{%
```

```

13      Package minifp: \MFPfiledate\space v\MFPfileversion. %
14      Macros for real number operations %
15      ^^Jand a stack-based programing language.^J}%
16 \else
17   \ProvidesPackage{minifp}[\MFPfiledate\space v\MFPfileversion. %
18     Macros for real number operations %
19     and a stack-based programing language.]%
20   \DeclareOption{extra}{\def\MFPextra{}}%
21   \ProcessOptions\relax
22 \fi
23 \catcode64=11
24 \ifx\MFPextra\UndEfInEd
25   \def\MFP@loadextra{}%
26 \else
27   \def\MFP@loadextra{\input mfpextra\relax}%
28 \fi
29 \def\MFPloadextra{\input mfpextra\relax}%
30 \catcode46=12
31 \catcode60=12
32 \catcode62=12

    We check for LATEX (ignoring LATEX209); \MFP@ifnoLaTeX...\MFP@end is
    skipped in LATEX and executed otherwise.
33 \long\def\gobbleto\MFP@end#1\MFP@end{}%
34 \def\MFP@end{\@empty}%
35 \ifx\documentclass\UndEfInEd
36   \def\MFP@ifnoLaTeX{}%
37 \else
38   \let\MFP@ifnoLaTeX\gobbleto\MFP@end
39 \fi

    We have LATEX's \zap@space. It pretty much must be used inside \edef or
    other purely expansion context. The rest of these are standard LATEX internals.
    Note that the token list that \zap@space is applied to should probably never
    contain braces or expandable tokens.
    Usage: \edef\X{\zap@space<tokens>\@empty}
    The space is necessary in case none exist; the \@empty terminates the loop.
40 \let\@xp\expandafter
41 \def\@XP{\@xp\@xp\@xp}%
42 \MFP@ifnoLaTeX
43   \def\@empty{}%
44   \long\def\@gobble#1{}%
45   \def\zap@space#1 #2{#1\ifx#2\@empty\else\@xp\zap@space\fi#2}%
46   \long\def\@ifnextchar#1#2#3{%
47     \let\reserved@d#1%
48     \def\reserved@a{#2}%
49     \def\reserved@b{#3}%
50     \futurelet\@let@token\@ifnch}%
51   \def\@ifnch{%
52     \ifx\@let@token\@sptoken

```

```

53     \let\reserved@c\@xifnch
54     \else
55         \ifx\@let@token\reserved@d
56             \let\reserved@c\reserved@a
57         \else
58             \let\reserved@c\reserved@b
59         \fi
60     \fi
61     \reserved@c}%
62 {%
63     \def\:{\global\let@sptoken= }\: %
64     \def\:{\@xifnch}\@xp\gdef\:{\futurelet\@let@token\@ifnch}%
65 }%
66 \def\@ifstar#1{\@ifnextchar*{\@firstoftwo{#1}}}%
67 \long\def\@firstofone #1{#1}%
68 \long\def\@firstoftwo #1#2{#1}%
69 \long\def\@secondoftwo#1#2{#2}%
70 \MFP@end

```

We need to divide by both  $10^4$  and  $10^8$  several times. I could have allocated two count registers, but have taken the approach of only using those for intermediate calculations.

```

71 \def\MFP@tttfour {10000}% ttt = Ten To The
72 \def\MFP@ttteight{100000000}%

```

These are for manipulating digits. The `\...ofmany` commands require a sequence of arguments (brace groups or tokens) followed by `\MFP@end`. The minimum number of required parameters is surely obvious. For example, `\MFP@ninthofmany` must be used like

```
\MFP@ninthofmany<9 or more arguments>\MFP@end
```

All these are fully expandable.

```

73 \def\MFP@oneofmany#1#2\MFP@end{#1}%
74 \def\MFP@fifthofmany#1#2#3#4#5#6\MFP@end{#5}%
75 \def\MFP@ninthofmany#1#2#3#4#5#6#7#8{\MFP@oneofmany}%
76 \def\MFP@eightofmany#1#2#3#4#5#6#7#8#9\MFP@end{#1#2#3#4#5#6#7#8}%

```

## 3.2 Processing numbers and the stack

Our stack stores elements in groups, like

```
{-1.234567890}{0.00001234}\MFP@eos
```

with an end marker. The purpose of the marker is to prevent certain parameter manipulations from stripping off braces. This means we can't use `\@empty` to test for an empty stack. At the moment, only `\Rpop` actually checks, but all other stack commands (so far) use `\Rpop` to get their arguments.

```

77 \let\MFP@eos\relax
78 \def\MFP@EOS{\MFP@eos}%
79 \def\MFP@initRstack{\def\MFP@Rstack{\MFP@eos}}%
80 \MFP@initRstack

```

Define some scratch registers for arithmetic operations. We don't care that these might be already in use, as we only use them inside a group. However, we need one counter that will not be messed with by any of our operations. I must be sure not to use commands that change `\MFP@loopctr` in code that depends on it.

```
81 \countdef \MFP@tempa 0
82 \countdef \MFP@tempb 2
83 \countdef \MFP@tempc 4
84 \countdef \MFP@tempd 6
85 \countdef \MFP@tempe 8
86 \countdef \MFP@tempf 10
87 \newcount \MFP@loopctr
```

The following can only be used where unrestricted expansion is robust. It will allow results obtained inside a group to survive the group, but not be unrestrictedly global. Example: the code

```
\MFP@endgroup@after{\def\noexpand\MFP@z@Val{\MFP@z@Val}}
```

becomes

```
\edef\x{\endgroup\def\noexpand\MFP@z@Val{\MFP@z@Val}}\x
```

which gives, upon expansion of `\x`,

```
\endgroup\def\MFP@z@Val{\expansion-of-\MFP@z@Val}
```

which defines `\MFP@z@Val` outside the current group to equal its expansion within the current group (provided the group was started with `\begingroup`).

We define a `\MFP@returned@values` to make all the conceivable produced values survive the group. The `\MFP@curr@Sgn` part is to permit testing the sign of the result and allow conditional code based on it.

I have been lax at making sure `\MFP@z@Ovr` is properly initialized and properly checked whenever it could be relevant, and properly passed on. I think every internal command `\MFP@Rxxx` should ensure it starts being zero and ends with a numerical value. At one time division could leave it undefined.

`\MFP@subroutine` executes its argument (typically a single command) with a wrapper that initializes all the macros that might need initializing, and returns the necessary results.

```
88 \def\MFP@endgroup@after#1{\edef\x{\endgroup#1}\x}%
89 \def\MFP@afterdef{\def\noexpand}%
90 \def\MFP@returned@values{%
91   \MFP@afterdef\MFP@z@Val{\MFP@z@Sign\MFP@z@Int.\MFP@z@Frc}%
92   \MFP@afterdef\MFP@z@Ovr{\MFP@z@Ovr}%
93   \MFP@afterdef\MFP@z@Und{\MFP@z@Und}%
94   \MFP@afterdef\MFP@curr@Sgn{\MFP@z@Sgn}}%
95 \def\MFP@subroutine#1{%
96   \begingroup
97   \MFP@Rzero
98   \def\MFP@z@Ovr{0}%
99   \def\MFP@z@Und{0}%
100   #1%
101   \MFP@endgroup@after\MFP@returned@values}%
102 \def\MFP@Rzero{%
```

```

103 \def\MFP@z@Sgn{0}%
104 \def\MFP@z@Int{0}%
105 \def\MFP@z@Frc{00000000}%

\EndofStack      We define here the error messages: popping from an empty stack and dividing
                  by zero. In addition to the error messages, we provide some default values that
                  hopefully allow some operations to continue.
                  We also have a warning or two.

106 \def\MFP@errmsg#1#2{%
107 \begingroup
108 \newlinechar'\^^J\let~\space
109 \def\MFP@msgbreak{^^J~~~~~}%
110 \edef\reserved@a{\errhelp{#2}}\reserved@a
111 \errmessage{MiniFP error: #1}%
112 \endgroup}%
113 \def\MFP@popempty@err{%
114 \MFP@errmsg{cannot POP from an empty stack}%
115 {There were no items on the stack for the POP operation. %
116 If you continue, ^^Jthe macro will contain the %
117 value \EndofStack.}%
118 \def\EndofStack{0.00000000}%
119 \def\MFP@dividebyzero@err{%
120 \MFP@errmsg{division by zero}%
121 {You tried to divide by zero. What were you thinking? %
122 If you continue, ^^Jthe value assigned will be either %
123 \ZeroOverZeroInt.\ZeroOverZeroFrac~(numerator=0) or %
124 ^^J+/-\xOverZeroInt.\xOverZeroFrac~(numerator<>0).}%
125 \def\MFP@warn#1{%
126 \begingroup
127 \newlinechar'\^^J\let~\space
128 \def\MFP@msgbreak{^^J~~~~~}%
129 \immediate\write16{^^JMiniFP warning: #1.^^J}%
130 \endgroup}%

\MaxRealInt      These are the largest possible integer and fractional parts of a real number.
\MaxRealFrac      They are returned for division by zero, for logarithm of zero, and when overflow
                  is detected in the exponential function.

131 \def\MaxRealInt {99999999}%
132 \def\MaxRealFrac {99999999}%

\xOverZeroInt     These are the results returned when trying to divide by zero. Two are used
\xOverZeroFrac     when dividing a nonzero number by zero and two when trying to divide zero
\ZeroOverZeroInt   by zero.
\ZeroOverZeroFrac

133 \def\xOverZeroInt {\MaxRealInt}%
134 \def\xOverZeroFrac {\MaxRealFrac}%
135 \def\ZeroOverZeroInt {0}%
136 \def\ZeroOverZeroFrac{00000000}%

```

These macros strip the spaces, process a number into sign, integer and fractional parts, and pad the fractional part out to eight decimals. They are used in



*push* so that the stack will only contains reals in a normalized form. Some of them are also used to preprocess the reals in the operand versions of commands

The `\MFP@*@Int` and `\MFP@*@Frc` parts are always positive, the sign being stored in `\MFP@*@Sgn` as  $-1$ ,  $0$  or  $1$ .

We strip the spaces and pad the fractional parts separately because they are unnecessary when processing *poped* reals (though they wouldn't hurt).

The number to be parsed is `#4` and the macros to contain the parts are the first three arguments. Since we normally call `\MFPparse@real` with one of two sets of macros, we have two shortcuts for those cases.

```

137 \def\MFPparse@real#1#2#3#4{%
138   \MFPnospace@def\MFPtemp@Val{#4}%
139   \MFPprocess@into@parts\MFPtemp@Val#1#2#3%
140   \MFP@padtoeight#3}%
141 \def\MFPparse@x{\MFPparse@real\MFP@x@Sgn\MFP@x@Int\MFP@x@Frc}%
142 \def\MFPparse@y{\MFPparse@real\MFP@y@Sgn\MFP@y@Int\MFP@y@Frc}%

```

This macro strips all spaces out of the number (not just before and after). It takes a macro that will hold the result, followed by the number (as a macro or a group of actual digits).

```

143 \def\MFPnospace@def#1#2{%
144   \edef#1{#2\space}\edef#1{\@xp\zap@space#1\@empty}}%

```

This is the process that splits a number into parts. The biggest difficulty is obtaining the sign. All four arguments are macros, with the first one holding the number. Following that are the macros to hold the sign, integer and fractional parts.

```

145 \def\MFPprocess@into@parts#1#2#3#4{%
146   \@xp\MFPsplit@dot#1..\MFP@end #3#4%

```

At this point `#3` holds the part before the dot (or the whole thing if there was no dot) and `#4` holds the part after the dot, (or nothing). Now is the first place where having at most eight digits simplifies things. Note that `#3` could contain any number of consecutive signs followed by up to eight digits. It could be zero or empty, so to avoid losing the sign we append a `1` (for up to nine digits). We temporarily define the sign based on the result, but may need to drop it if both the integer and fractional parts are zero.

Prepending a zero to the fractional part permits it to be empty. In the final `\edef`, `#3` is made positive.

```

147   \ifnum#31<0 \def#2{-1}%
148   \else \def#2{1}%
149   \fi
150   \ifnum #30=0
151     \def#3{0}%
152     \ifnum 0#4=0 \def#2{0}\fi
153   \fi
154   \edef#3{\number \ifnum #2<0 -\fi#3}}%

```

This only copies the parts before and after the dot, `#1` and `#2`, into macros `#4` and `#5`.

```
155 \def\MFPsplit@dot#1.#2.#3\MFP@end#4#5{\edef#4{#1}\edef#5{#2}}%
```

This is used to pad the fractional part to eight places with zeros. If a number with more than eight digits survives to this point, it gets truncated.

```
156 \def\MFP@padtoeight#1{%
157 \edef#1{\xp\MFP@eightofmany#100000000\MFP@end}}%
```

These take operands off the stack. We know already that there are no spaces and that the fractional part has eight digits.

```
158 \def\MFPgetoperand@x{\Rpop\MFP@x@Val
159 \MFPprocess@into@parts\MFP@x@Val\MFP@x@Sgn\MFP@x@Int\MFP@x@Frc}%
160 \def\MFPgetoperand@y{\Rpop\MFP@y@Val
161 \MFPprocess@into@parts\MFP@y@Val\MFP@y@Sgn\MFP@y@Int\MFP@y@Frc}%

```

Concatenate an argument (or two) to the front of stack. The material must already be in correct format.

```
162 \def\MFP@Rcat#1{\edef\MFP@Rstack{#1}\MFP@Rstack}%
163 \def\MFP@Rcattwo#1#2{\edef\MFP@Rstack{#1}{#2}\MFP@Rstack}}%
```

Convert from a signum (a number) to a sign (— or nothing):

```
164 \def\MFP@Sign#1{\ifnum#1<0 -\fi}%
165 \def\MFP@x@Sign{\MFP@Sign\MFP@x@Sgn}%
166 \def\MFP@y@Sign{\MFP@Sign\MFP@y@Sgn}%
167 \def\MFP@z@Sign{\MFP@Sign\MFP@z@Sgn}%

```

Sometimes only parts of the number needs changing (used in CHS, ABS). This copies the integer and fractional parts of  $x$  into  $z$ .

```
168 \def\copyMFP@x{\edef\MFP@z@Int{\MFP@x@Int}\edef\MFP@z@Frc{\MFP@x@Frc}}%
```

We use `\MFPpush@result` to put the result of internal operations onto the stack. For convenience, we also have it set the sign flags.

```
169 \def\MFPpush@result{\MFP@Rchk\MFPcurr@Sgn\MFP@Rcat\MFP@z@Val}%

```

When *pop* encounters an empty stack it gobbles the code that would perform the *pop* (#1) and defines the macro (#2) to contain `\EndofStack`. The default meaning for this macro is 0.

```
170 \def\if@EndofStack{%
171 \ifx\MFP@EOS\MFP@Rstack
172 \xp@firstoftwo
173 \else
174 \xp@secondoftwo
175 \fi}%

```

The macro `\Rpop` calls `\MFP@popit` followed by the contents of the stack, the token `\MFP@end` and the macro to *pop* into. If the stack is not empty, `\doMFP@popit` will read the first group #1 into that macro #3, and then redefine the stack to be the rest of the argument #2. If the stack is empty, `\doMFP@EOS` will equate the macro to `\EndofStack` (initialized to 0.00000000) after issuing an error message.

```
176 \def\MFP@popit{\if@EndofStack\doMFP@EOS\doMFP@popit}%
177 \def\doMFP@EOS#1\MFP@end#2{\MFP@popempty@err\let#2\EndofStack}%
178 \def\doMFP@popit#1#2\MFP@end#3{\edef\MFP@Rstack{#2}\edef#3{#1}}%
```

### 3.3 The user-level operations

All operations that can be done on arguments as well as the stack will have a common format: The stack version pops the requisite numbers and splits them into internal macros (`\MFPgetoperand*`), runs an internal command that operates on these internal macros, then “pushes” the result returned. The internal commands take care to return the result in proper form so we don’t actually run `\Rpush`, but only `\MFPpush@result`.

The operand version processes the operands into normalized form (as if pushed, using `\MFPparse*`), then proceeds as in the stack version, but copies the result into the named macro instead of to the stack (`\MFPstore@result`).

For unary operations we process one argument or stack element. We call it *x* and use the *x* version of all macros. All internal commands (*#1*) return the results in *z* versions.

`\MFPchk` The `\MFPchk` command examines its argument and sets a flag according to its sign.

```
179 \def\MFPchk#1{%
180   \MFPparse@x{#1}%
181   \MFP@Rchk\MFP@x@Sgn}%
```

We make `\MFP@Rchk` a little more general than is strictly needed here, by giving it an argument (instead of only examining `\MFP@x@Sgn`). This is so we can apply it to the results of operations (which would be in `\MFPcurr@Sgn`).

```
182 \def\MFP@Rchk#1{%
183   \MFPclear@flags
184   \ifnum#1>0 \MFP@postrue
185   \else\ifnum#1<0 \MFP@negtrue
186   \else \MFP@zerotrue
187   \fi\fi}%
188 \def\MFPclear@flags{\MFP@zerofalse \MFP@negfalse \MFP@posfalse}%
```

`\IFzero` These are the user interface to the internal `\ifMFP@zero`, `\ifMFP@neg`,

`\IFneg` `\ifMFP@pos`

```
\IFpos 189 \def\IFzero{\ifMFP@zero\@xp\@firstoftwo\else\@xp\@secondoftwo\fi}%
190 \def\IFneg {\ifMFP@neg \@xp\@firstoftwo\else\@xp\@secondoftwo\fi}%
191 \def\IFpos {\ifMFP@pos \@xp\@firstoftwo\else\@xp\@secondoftwo\fi}%
192 \newif\ifMFP@zero \newif\ifMFP@neg \newif\ifMFP@pos
```

`\MFPcmp` Our comparison commands parallel the check-sign commands. They even reuse the same internal booleans. The differences: the internal `\MFP@Rcmp` doesn’t take arguments and the comparison test is a little more involved. We could simply subtract, which automatically sets the internal booleans, but it is way more efficient to compare sizes directly.

```
193 \newif\ifMFPdebug
194 \def\MFPcmp#1#2{\MFPparse@x{#1}\MFPparse@y{#2}\MFP@Rcmp}%
195 \def\MFP@Rcmp{\MFPclear@flags
196   \ifnum \MFP@x@Sign\MFP@x@Int>\MFP@y@Sign\MFP@y@Int\relax
197     \MFP@postrue
198   \else\ifnum \MFP@x@Sign\MFP@x@Int<\MFP@y@Sign\MFP@y@Int\relax
```

```

199 \MFP@negtrue
200 \else\ifnum \MFP@x@Sign\MFP@x@Frc>\MFP@y@Sign\MFP@y@Frc\relax
201 \MFP@postrue
202 \else\ifnum \MFP@x@Sign\MFP@x@Frc<\MFP@y@Sign\MFP@y@Frc\relax
203 \MFP@negtrue
204 \else
205 \MFP@zerotrue
206 \fi\fi\fi\fi}%
207 \let\IFeq\IFzero\let\IFlt\IFneg \let\IFgt\IFpos

```

Given an operation (*pop*, *chs*, or whatever), the stack version will have the same name with “R” (for “real”) prepended. The operand versions will have the same name with “MFP” prepended. The internal version has the same name as the stack version, with an “MFP@” prepended.

The unary operations are:

**chs** change sign of  $x$ .  
**abs** absolute value of  $x$ .  
**dbl** double  $x$ .  
**halve** halve  $x$ .  
**sgn** +1, -1 or 0 depending on the sign of  $x$ .  
**sq** square  $x$ .  
**int** zero out the fractional part of  $x$ .  
**frac** zero out the integer part of  $x$ .  
**floor** largest integer not exceeding  $x$ .  
**ceil** smallest integer not less than  $x$ .

The binary operations are ( $x$  represents the first and  $y$  second):

**add** add  $x$  and  $y$ .  
**sub** subtract  $y$  from  $x$ .  
**mul** multiply  $x$  and  $y$ .  
**div** divide  $x$  by  $y$ .

There are also some operations that do not actually change any values, but may change the stack or the state of some boolean:

**cmp** compare  $x$  and  $y$  (stack version does not change stack).  
**chk** examine the sign of  $x$  (stack version does not change stack).  
**dup** stack only, duplicate the top element of the stack.  
**push** stack only, put a value onto the top of the stack.  
**pop** stack only, remove the top element of the stack, store it in a variable.  
**exch** stack only, exchange top two elements of the stack.

**\startMFPprogram** The purpose of **\startMFPprogram** is to start the group, inside of which all the  
**\stopMFPprogram** stack operations can be used. The ensuing **\stopMFPprogram** closes the group.

```

208 \def\startMFPprogram{%
209 \begingroup

```

`\Rchs` Then we give definitions to all the stack-based macros. These definitions are  
`\Rabs` all lost after the group ends.  
`\Rdbl` We start with the unary operations. Note that all they do is call a wrapper  
`\Rhalve` macro `\MFP@stack@Unary` with an argument which is the internal version of the  
`\Rsgn` command.  
`\Rsqr` 210 `\def\Rchs {\MFP@stack@Unary\MFP@Rchs}%`  
`\Rinv` 211 `\def\Rabs {\MFP@stack@Unary\MFP@Rabs}%`  
`\Rint` 212 `\def\Rdbl {\MFP@stack@Unary\MFP@Rdbl}%`  
`\Rfrac` 213 `\def\Rhalve{\MFP@stack@Unary\MFP@Rhalve}%`  
`\Rfloor` 214 `\def\Rsgn {\MFP@stack@Unary\MFP@Rsgn}%`  
`\Rceil` 215 `\def\Rsqr {\MFP@stack@Unary\MFP@Rsqr}%`  
`\Rinv` 216 `\def\Rinv {\MFP@stack@Unary\MFP@Rinv}%`  
`\Rincr` 217 `\def\Rint {\MFP@stack@Unary\MFP@Rint}%`  
`\Rdecr` 218 `\def\Rfrac {\MFP@stack@Unary\MFP@Rfrac}%`  
`\Rzero` 219 `\def\Rfloor{\MFP@stack@Unary\MFP@Rfloor}%`  
220 `\def\Rceil {\MFP@stack@Unary\MFP@Rceil}%`  
221 `\def\Rincr {\MFP@stack@Unary\MFP@Rincr}%`  
222 `\def\Rdecr {\MFP@stack@Unary\MFP@Rdecr}%`  
223 `\def\Rzero {\MFP@stack@Unary\MFP@Rzero}%`  
  
`\Radd` Then the binary operations, which again call a wrapper around the internal  
`\Rsub` version.  
`\Rmul` 224 `\def\Radd {\MFP@stack@Binary\MFP@Radd}%`  
`\Rmpy` 225 `\def\Rmul {\MFP@stack@Binary\MFP@Rmul}%`  
`\Rdiv` 226 `\let\Rmpy\Rmul`  
`\Rmin` 227 `\def\Rsub {\MFP@stack@Binary\MFP@Rsub}%`  
`\Rmax` 228 `\def\Rdiv {\MFP@stack@Binary\MFP@Rdiv}%`  
229 `\def\Rmin {\MFP@stack@Binary\MFP@Rmin}%`  
230 `\def\Rmax {\MFP@stack@Binary\MFP@Rmax}%`  
  
`\Rnoop` And finally some special commands. There is a no-op and commands for com-  
`\Rcmp` paring, checking, and manipulation of the stack.  
`\Rchk` 231 `\let\Rnoop\relax`  
`\Rpush` 232 `\def\Rcmp{%`  
`\Rpop` 233 `\MFPgetoperand@y\MFPgetoperand@x`  
`\Rexch` 234 `\MFP@Rcat\MFP@x@Val\MFP@Rcat\MFP@y@Val`  
`\Rdup` 235 `\MFP@Rcmp}%`  
236 `\def\Rchk{%`  
237 `\MFPgetoperand@x`  
238 `\MFP@Rcat\MFP@x@Val`  
239 `\MFP@Rchk\MFP@x@Sgn}%`  
240 `\def\Rpush##1{%`  
241 `\MFPparse@x{##1}%`  
242 `\edef\MFP@z@Val{\MFP@x@Sign\MFP@x@Int.\MFP@x@Frc}%`  
243 `\edef\MFP@curr@Sgn{\MFP@x@Sgn}%`  
244 `\MFPpush@result}%`  
245 `\def\Rpop{\@xp\MFP@popit\MFP@Rstack\MFP@end}%`  
246 `\def\Rexch{%`  
247 `\Rpop\MFP@x@Val\Rpop\MFP@y@Val`  
248 `\MFP@Rcattwo\MFP@y@Val\MFP@x@Val}%`

```

249 \def\Rdup{%
250   \Rpop\MFP@x@Val
251   \MFP@Rcattwo\MFP@x@Val\MFP@x@Val}%

    If mfpxextra.tex is input, then \MFP@Rextra makes the additional commands
    in that file available to an MINIFP program.

\Global      The last four commands allow computed values to be made available outside
\GlobalStack the program group
\Export      252 \MFP@Rextra
\ExportStack 253 \let\Global\MFP@Global
              254 \let\GlobalStack\MFP@GlobalStack
              255 \let\Export\MFP@Export
              256 \let\ExportStack\MFP@ExportStack}%
              257 \def\stopMFPprogram{\@xp\endgroup\MFPprogram@returns}%
              258 \let\MFP@Rextra\@empty
              259 \let\MFPprogram@returns\@empty

\MFPchs      Now we define the operand versions. These also are defined via a wrapper
\MFPabs      command that executes the very same internal commands as the stack versions.
\MFPdbl      First the unary operations.
\MFPhalve    260 \def\MFPchs {\MFP@op@Unary\MFP@Rchs}%
\MFPsgn      261 \def\MFPabs {\MFP@op@Unary\MFP@Rabs}%
\MFPsq       262 \def\MFPdbl {\MFP@op@Unary\MFP@Rdbl}%
\MFPinv      263 \def\MFPhalve{\MFP@op@Unary\MFP@Rhalve}%
\MFPint      264 \def\MFPsgn {\MFP@op@Unary\MFP@Rsgn}%
\MFPfrac     265 \def\MFPsq {\MFP@op@Unary\MFP@Rsq}%
\MFPfloor    266 \def\MFPinv {\MFP@op@Unary\MFP@Rinv}%
\MFPceil     267 \def\MFPint {\MFP@op@Unary\MFP@Rint}%
\MFPincr     268 \def\MFPfrac {\MFP@op@Unary\MFP@Rfrac}%
\MFPdecr     269 \def\MFPfloor{\MFP@op@Unary\MFP@Rfloor}%
\MFPzero     270 \def\MFPceil {\MFP@op@Unary\MFP@Rceil}%
\MFPstore    271 \def\MFPincr {\MFP@op@Unary\MFP@Rincr}%
              272 \def\MFPdecr {\MFP@op@Unary\MFP@Rdecr}%
              273 \def\MFPzero {\MFP@op@Unary\MFP@Rzero}%
              274 \def\MFPstore{\MFP@op@Unary\MFP@Rstore}%

\MFPadd      Then the binary operations.
\MFPsub      275 \def\MFPadd{\MFP@op@Binary\MFP@Radd}%
\MFPmul      276 \def\MFPmul{\MFP@op@Binary\MFP@Rmul}%
\MFPmpy      277 \let\MFPmpy\MFPmul
\MFPdiv      278 \def\MFPsub{\MFP@op@Binary\MFP@Rsub}%
\MFPmin      279 \def\MFPdiv{\MFP@op@Binary\MFP@Rdiv}%
\MFPmax      280 \def\MFPmin{\MFP@op@Binary\MFP@Rmin}%
              281 \def\MFPmax{\MFP@op@Binary\MFP@Rmax}%

    A nullary operation is one that produces a result with no operand. Thus, it
    could return a fixed constant, or it could perform calculations that obtain input
    from the system (e.g., current time). At the moment we don't define any.

282 \def\MFP@stack@Nullary#1{%
283   \MFP@subroutine{#1}\MFPpush@result}%
284 \def\MFP@op@Nullary#1{%

```

```
285 \MFP@subroutine{#1}\MFPstore@result}%
```

These are the wrappers for unary operations. The operand versions have a second argument, the macro that stores the result. But this will be the argument of `\MFPstore@result`.

```
286 \def\MFP@stack@Unary#1{%
287   \MFPgetoperand@x
288   \MFP@subroutine{#1}\MFPpush@result}%
289 \def\MFP@op@Unary#1#2{%
290   \MFPparse@x{#2}%
291   \MFP@subroutine{#1}\MFPstore@result}%
292 \def\MFPstore@result#1{\MFP@Rchk\MFPcurr@Sgn\edef#1{\MFP@z@Val}}%
```

These are the wrappers for binary operations. The top level definitions are almost identical to those of the unary operations. The only difference is they *pop* or parse two operands.

```
293 \def\MFP@stack@Binary#1{%
294   \MFPgetoperand@y \MFPgetoperand@x
295   \MFP@subroutine{#1}\MFPpush@result}%
296 \def\MFP@op@Binary#1#2#3{%
297   \MFPparse@x{#2}\MFPparse@y{#3}%
298   \MFP@subroutine{#1}\MFPstore@result}%
\MFPnoop
```

We end with a traditional, but generally useless command, the no-op, which does nothing. It doesn't even have a wrapper.

```
299 \let\MFPnoop\relax
```

### 3.4 The internal computations

To change the sign or get the absolute value, we just need to set the value of `\MFP@x@Sgn`.

```
300 \def\MFP@Rabs{%
301   \copy\MFP@x \edef\MFP@z@Sgn{\ifnum\MFP@x@Sgn=0 0\else1\fi}}%
302 \def\MFP@Rchs{\copy\MFP@x \edef\MFP@z@Sgn{\number-\MFP@x@Sgn}}%
```

The doubling and halving operations are more efficient ways to multiply or divide a number by 2. For doubling, copy  $x$  to  $y$  and add. For halving, we use basic  $\text{\TeX}$  integer division, more efficient than multiplying by 0.5 and far more than using `\MFP@Rdiv`.

In `\MFP@Rhalve`, we add 1 to the fractional part for rounding purposes, and we move any odd 1 from the end of the integer part to the start of the fractional part.

```
303 \def\MFP@Rdbl{\MFP@Rcopy xy\MFP@Radd}%
304 \def\MFP@Rhalve{%
305   \MFP@tempa\MFP@x@Int
306   \MFP@tempb\MFP@x@Frc\relax
307   \ifodd\MFP@tempb
308     \def\MFP@z@Und{5}%
309     \advance\MFP@tempb 1
310     \ifnum\MFP@ttteight=\MFP@tempb
```

```

311      \MFP@tempb0 \advance\MFP@tempa1
312      \fi
313      \fi
314      \ifodd \MFP@tempa
315        \advance\MFP@tempb \MFP@ttteight\relax
316      \fi
317      \divide\MFP@tempa 2
318      \divide\MFP@tempb 2
319      \MFP@Rloadz\MFP@x@Sgn\MFP@tempa\MFP@tempb}%

```

The signum is 0.0, 1.0 or  $-1.0$  to match the sign of  $x$ .

```

320 \def\MFP@Rsgn{\MFP@Rloadz \MFP@x@Sgn{\ifnum\MFP@x@Sgn=0 0\else1\fi}0}%

```

The squaring operation just calls `\MFP@Rmul` after copying  $x$  to  $y$ . Its gain in efficiency over a multiplication is that it can skip preprocessing of the second (identical) operand.

```

321 \def\MFP@Rsqr{\MFP@Rcopy xy\MFP@Rmul}%

```

The inversion operation just calls `\MFP@Rdiv` after copying  $x$  to  $y$  and 1 to  $x$ . Its advantage over a divide is it skips the preprocessing of 1 as an operand.

```

322 \def\MFP@Rinv{\MFP@Rcopy xy\MFP@Rload x110\MFP@Rdiv}%

```

Integer part: replace fractional part with zeros.

```

323 \def\MFP@Rint{%
324   \MFP@Rloadz {\ifnum\MFP@x@Int=0 0\else\MFP@x@Sgn\fi}\MFP@x@Int 0}%

```

Fractional part: replace integer part with a zero.

```

325 \def\MFP@Rfrac{%
326   \MFP@Rloadz {\ifnum\MFP@x@Frc=0 0\else\MFP@x@Sgn\fi}0\MFP@x@Frc}%

```

To increment and decrement by 1, except in border cases, we need only address the integer part of a number. This doesn't seem so simple written out but, even so, it is more efficient than full-blown addition. It would be very slightly more efficient if `\MFP@Rdecr` did not call `\MFP@Rincr`, but instead was similarly coded.

```

327 \def\MFP@Rincr{%
328   \ifnum\MFP@x@Sgn<0
329     \ifcase\MFP@x@Int
330       \MFP@tempa\MFP@ttteight
331       \advance\MFP@tempa -\MFP@x@Frc\relax
332       \MFP@Rloadz 10\MFP@tempa
333     \or
334       \MFP@Rloadz{\ifnum\MFP@x@Frc=0 0\else -1\fi}0\MFP@x@Frc
335     \else
336       \MFP@tempa\MFP@x@Int
337       \advance\MFP@tempa -1
338       \MFP@Rloadz{-1}\MFP@tempa\MFP@x@Frc
339     \fi
340   \else
341     \MFP@tempa\MFP@x@Int
342     \advance\MFP@tempa 1
343     \MFP@Rloadz 1\MFP@tempa\MFP@x@Frc
344   \fi}%

```



```

345 \def\MFP@Rdecr{%
346   \edef\MFP@x@Sgn{\number -\MFP@x@Sgn}\MFP@Rincr
347   \edef\MFP@z@Sgn{\number -\MFP@z@Sgn}}%
348 \def\MFP@Rstore{\MFP@Rcopy xz}%

```

The floor of a real number  $x$  is the largest integer not larger than  $x$ . The ceiling is the smallest integer not less than  $x$ . For positive  $x$ , floor is the same as integer part. Not true for negative  $x$ . Example:  $\text{int}(-1.5) = -1$  but  $\text{floor} = -2$

We use the same code to get floor or ceiling, the appropriate inequality character being its argument.

```

349 \def\MFP@Rfloorceil#1{%
350   \MFP@tempa\MFP@x@Int\relax
351   \ifnum \MFP@x@Sgn #10
352     \ifnum\MFP@x@Frc=0
353       \else
354         \advance\MFP@tempa1
355       \fi
356     \fi
357   \MFP@Rloadz{\ifnum\MFP@x@Int=0 0\else\MFP@x@Sgn\fi}\MFP@tempa0}%
358 \def\MFP@Rfloor{\MFP@Rfloorceil<}%
359 \def\MFP@Rceil {\MFP@Rfloorceil>}%

```

For multiplication, after the usual break into integer and fractional parts, we further split these parts into 4-digit pieces with `\MFP@split`. The first argument (#1) holds the eight digit number, then #2 is a macro that will hold the top four digits and #3 will hold the bottom four.

```

360 \def\MFP@split#1#2#3{%
361   \begingroup
362     \MFP@tempa#1\relax
363     \MFP@tempb\MFP@tempa
364     \divide\MFP@tempb by\MFP@tttfour
365     \edef#2{\number\MFP@tempb}%
366     \multiply\MFP@tempb by\MFP@tttfour
367     \advance\MFP@tempa-\MFP@tempb
368   \MFP@endgroup@after{%
369     \MFP@afterdef#2{#2}%
370     \MFP@afterdef#3{\number\MFP@tempa}%
371   }%
372 %
373 \def\MFP@@split{%
374   \MFP@split\MFP@x@Int\MFP@x@Int@ii\MFP@x@Int@i
375   \MFP@split\MFP@x@Frc\MFP@x@Frc@i\MFP@x@Frc@ii
376   \MFP@split\MFP@y@Int\MFP@y@Int@ii\MFP@y@Int@i
377   \MFP@split\MFP@y@Frc\MFP@y@Frc@i\MFP@y@Frc@ii}%

```

We will store the intermediate and final products in `\MFP@z@*`. Each one is ultimately reduced to four digits, like the parts of  $x$  and  $y$ . As each base-10000 digit of  $y$  is multiplied by a digit of  $x$ , we add the result to the appropriate digit of the partial result  $z$ .

The underflow ends up in `\MFP@z@Frc@iv` and `\MFP@z@Frc@iii`. Overflow will be in `\MFP@z@Int@iii`. Unlike the rest, it can be up to eight digits because we do not need to carry results out of it.

This command prepends zeros so a number fills four slots. Here `#1` is a macro holding the value and it is redefined to contain the result. A macro that calls this should ensure that `#1` is not empty and is less than 10,000.

```
378 \def\makeMFP@fourdigits#1{%
379   \edef#1{\@xp\MFP@fifthofmany\number#1}{0}{00}{000}\MFP@end\number#1}}%
```

This is the same, but produces eight digits. Similarly `#1` should be nonempty and less than 100,000,000.

```
380 \def\makeMFP@eightdigits#1{%
381   \edef#1{\@xp\MFP@ninthofmany\number#1%
382     {}{0}{00}{000}{0000}{00000}{000000}{0000000}\MFP@end\number#1}}%
```

The following macros implement carrying. The macros `\MFP@carrya` and `\MFP@carrym` should be followed by two macros that hold numbers. The first number can have too many digits. These macros remove extra digits from the front and add their value to the number in the second macro (the “carry”). Both act by calling `\MFP@carry`, which is told the number of digits to keep via `#1` (10,000 for four digits, 100,000,000 for eight). The “a” in `\MFP@carrya` is for addition and “m” is for multiplication, which indicates where these will mainly be used.

```
383 \def\MFP@carrya{\MFP@carry\MFP@ttteight}%
384 \def\MFP@carrym{\MFP@carry\MFP@tttfour}%
385 \def\MFP@carry#1#2#3{%
386   \begingroup
387     \MFP@carryi{#1}#2#3%
388   \MFP@endgroup@after{%
389     \MFP@afterdef#3{\number\MFP@tempa}%
390     \MFP@afterdef#2{\number\MFP@tempb}%
391   }}%
```

This is the “internal” carry. `#1`, `#2`, and `#3` are as in `\MFP@carry`. Its advantage is that it can be used where `#2` and `#3` are not macros, leaving the result in `\MFP@tempa` and `\MFP@tempb` with `\MFP@tempb` in the correct range,  $[0, \#1)$ . Its disadvantage is it does not protect temporary registers. Warning: do not use it with `#2=\MFP@tempa` and do not use it without grouping if you want to preserve the values in these temporary count registers.

```
392 \def\MFP@carryi#1#2#3{%
393   \MFP@tempa=#3\relax
394   \MFP@tempb=#2\relax
395   \MFP@tempc=\MFP@tempb
396   \divide \MFP@tempc #1\relax
397   \advance \MFP@tempa \MFP@tempc
398   \multiply\MFP@tempc #1\relax
399   \advance \MFP@tempb -\MFP@tempc}%%
```

This adds `#1` to `#2`, the result goes into macro `#3`. This does no checking. It is basically used to add with macros instead of registers.

```

400 \def\MFP@addone#1#2#3{%
401   \begingroup
402     \MFP@tempa#1%
403     \advance\MFP@tempa#2\relax
404   \MFP@endgroup@after{%
405     \MFP@afterdef#3{\number\MFP@tempa}%
406   }}%

```

Multiply **#1** by `\MFP@tempb` and add to **#2**. `\MFP@tempb` is one digit (base=10000) of  $y$  in multiplying  $x \times y$ , **#1** (usually a macro) holds one digit of  $x$ . **#2** is a macro that will hold one digit of the final product  $z$ . The product is added to it (overflow is taken care of later by the carry routines).

```

407 \def\MFP@multiplyone#1#2{%
408   \MFP@tempa#1%
409   \multiply\MFP@tempa\MFP@tempb
410   \advance\MFP@tempa#2%
411   \edef#2{\number\MFP@tempa}}%

```

This does the above multiplication-addition for all four “digits” of  $x$ . This is where `\MFP@tempb` is initialized for `\MFP@multiplyone`. The first argument represents a digit of  $y$ , the remaining four arguments are macros representing the digits of  $z$  that are involved in multiplying the digits of  $x$  by **#1**.

```

412 \def\MFP@multiplyfour#1#2#3#4#5{%
413   \MFP@tempb #1\relax
414   \MFP@multiplyone\MFP@x@Int@ii #2%
415   \MFP@multiplyone\MFP@x@Int@i #3%
416   \MFP@multiplyone\MFP@x@Frc@i #4%
417   \MFP@multiplyone\MFP@x@Frc@ii #5}%

```

Now we begin the internal implementations of the binary operations. All four expect macros `\MFP@x@Sgn`, `\MFP@x@Int`, `\MFP@x@Frc`, `\MFP@y@Sgn`, `\MFP@y@Int` and `\MFP@y@Frc` to be the normalized parts of two real numbers  $x$  and  $y$ .

`\MFP@Rsub` just changes the sign of  $y$  and then calls `\MFP@Radd`.

`\MFP@Radd` checks whether  $x$  and  $y$  have same or different signs. In the first case we need only add absolute values and the sign of the result will match that of the operands. In the second case, finding the sign of the result is a little more involve (and “borrowing” may be needed).

```

418 \def\MFP@Rsub{\edef\MFP@y@Sgn{\number-\MFP@y@Sgn}\MFP@Radd}%
419 \def\MFP@Radd{%
420   \MFP@tempa\MFP@x@Sgn
421   \multiply\MFP@tempa\MFP@y@Sgn\relax
422   \ifcase\MFP@tempa
423     \ifnum \MFP@x@Sgn=0
424       \MFP@Rcopy yz%
425     \else
426       \MFP@Rcopy xz%
427     \fi
428   \or
429     \@xp\MFP@Radd@same
430   \else

```

```

431 \xp\MFP@Radd@diff
432 \fi}%

```

\MFP@Radd@same adds two numbers which have the same sign. The sign of the result is the common sign. The fractional and integer parts are added separately, then a carry is invoked. The overflow (\MFP@z@Ovr) could be only a single digit 0 or 1.

```

433 \def\MFP@Radd@same{%
434 \MFP@addone\MFP@x@Frc\MFP@y@Frc\MFP@z@Frc
435 \MFP@addone\MFP@x@Int\MFP@y@Int\MFP@z@Int
436 \MFP@carrya\MFP@z@Frc\MFP@z@Int
437 \MFP@carrya\MFP@z@Int\MFP@z@Ovr
438 \makeMFP@eightdigits\MFP@z@Frc
439 \edef\MFP@z@Sgn{\MFP@x@Sgn}}%

```

We are now adding two numbers with opposite sign. Since  $x \neq 0$  this is the same as  $\text{sgn}(x)(|x| - |y|)$ . So we subtract absolute values, save the result in \MFP@z@Sgn, \MFP@z@Int and \MFP@z@Frc (with the last two nonnegative, as usual), then change the sign of \MFP@z@Sgn if \MFP@x@Sgn is negative. Since the difference between numbers in  $[0, 10^8)$  has absolute value in that range, there is no carrying. However, there may be borrowing.

```

440 \def\MFP@Radd@diff{%
441 \MFP@addone\MFP@x@Frc{-\MFP@y@Frc}\MFP@z@Frc
442 \MFP@addone\MFP@x@Int{-\MFP@y@Int}\MFP@z@Int

```

Now we need to establish the sign and arrange the borrow. The sign of the result is the sign of \MFP@z@Int unless it is 0; in that case it, is the sign of \MFP@z@Frc. There must be a simpler coding, though.

```

443 \MFP@tempa=\MFP@z@Int
444 \MFP@tempb=\MFP@z@Frc\relax
445 \ifnum\MFP@tempa=0 \else \MFP@tempa=\MFP@Sign\MFP@tempa 1 \fi
446 \ifnum\MFP@tempb=0 \else \MFP@tempb=\MFP@Sign\MFP@tempb 1 \fi
447 \ifnum\MFP@tempa=0 \MFP@tempa=\MFP@tempb \fi

```

Now we have the sign of  $|x| - |y|$  in \MFP@tempa, and we multiply that sign by the sign of  $x$  to get \MFP@z@Sgn. Then we multiply the current value of  $z$  by that sign to get the absolute value, stored in \MFP@tempa and \MFP@tempb.

```

448 \edef\MFP@z@Sgn{\number\MFP@x@Sign\MFP@tempa}%
449 \MFP@tempb\MFP@tempa
450 \multiply\MFP@tempa \MFP@z@Int
451 \multiply\MFP@tempb \MFP@z@Frc\relax

```

What we should have now is a positive number which might still be represented with a negative fractional part. A human being performing the subtraction would have borrowed first. Being a computer, we do it last, and we're done.

```

452 \ifnum\MFP@tempb<0
453 \advance\MFP@tempb\MFP@ttteight
454 \advance\MFP@tempa-1
455 \fi
456 \edef\MFP@z@Int{\number\MFP@tempa}%

```

```

457 \edef\MFP@z@Frc{\number\MFP@tempb}%
458 \makeMFP@eightdigits\MFP@z@Frc}%

\MFP@Rmul first computes the (theoretical) sign of the product: if it is zero, re-
turn zero, otherwise provisionally set the sign of the product and call \MFP@@Rmul.
459 \def\MFP@Rmul{%
460 \ifnum\MFP@x@Sgn=0 \MFP@Rzero
461 \else\ifnum\MFP@y@Sgn=0 \MFP@Rzero
462 \else \edef\MFP@z@Sgn{\number\MFP@x@Sign\MFP@y@Sgn}%
463 \XP\MFP@@Rmul
464 \fi\fi}%

```

\MFP@@Rmul first initializes the macros that will hold the base-10000 digits of  $z$ . Then it splits the four expected macros into eight macros that hold the base-10000 digits for each of  $x$  and  $y$ . Then each digit of  $y$  is used to multiply the four digits of  $x$  and the results are added to corresponding digits of  $z$ .

```

465 \def\MFP@@Rmul{%
466 \def\MFP@z@Frc@iv {0}\def\MFP@z@Frc@iii{0}%
467 \def\MFP@z@Frc@ii {0}\def\MFP@z@Frc@i {0}%
468 \def\MFP@z@Int@i {0}\def\MFP@z@Int@ii {0}%
469 \def\MFP@z@Int@iii{0}%
470 \MFP@@split
471 \MFP@multiplyfour \MFP@y@Frc@ii \MFP@z@Frc@i
472 \MFP@z@Frc@ii \MFP@z@Frc@iii \MFP@z@Frc@iv
473 \MFP@multiplyfour \MFP@y@Frc@i \MFP@z@Int@i
474 \MFP@z@Frc@i \MFP@z@Frc@ii \MFP@z@Frc@iii
475 \MFP@multiplyfour \MFP@y@Int@i \MFP@z@Int@ii
476 \MFP@z@Int@i \MFP@z@Frc@i \MFP@z@Frc@ii
477 \MFP@multiplyfour \MFP@y@Int@ii \MFP@z@Int@iii
478 \MFP@z@Int@ii \MFP@z@Int@i \MFP@z@Frc@i

```

Now apply the carry routines on the underflow digits...

```

479 \MFP@carrym\MFP@z@Frc@iv\MFP@z@Frc@iii
480 \MFP@carrym\MFP@z@Frc@iii\MFP@z@Frc@ii

```

...and pause to round the lowest digit that will be kept...

```

481 \ifnum\MFP@z@Frc@iii<5000 \else
482 \MFP@tempb\MFP@z@Frc@ii
483 \advance\MFP@tempb1
484 \edef\MFP@z@Frc@ii{\number\MFP@tempb}%
485 \fi

```

...and continue carrying.

```

486 \MFP@carrym\MFP@z@Frc@ii\MFP@z@Frc@i
487 \MFP@carrym\MFP@z@Frc@i \MFP@z@Int@i
488 \MFP@carrym\MFP@z@Int@i \MFP@z@Int@ii
489 \MFP@carrym\MFP@z@Int@ii\MFP@z@Int@iii

```

To end, we arrange for all macros to hold four digits (except \MFP@z@Int@ii and \MFP@z@Int@iii which don't need leading 0s) and load them into the appropriate 8-digit macros. The underflow digits are stored in \MFP@z@Und in case we ever need to examine them (we now do: in our unit conversion routine \MFP@Dpmul),

and the overflow in \MFP@z@Ovr in case we ever want to implement an overflow error. Theoretically  $z \neq 0$ , but it is possible that  $z = 0$  after reducing to eight places. If so, we must reset \MFP@z@Sgn.

```

490 \makeMFP@fourdigits\MFP@z@Frc@iv
491 \makeMFP@fourdigits\MFP@z@Frc@iii
492 \makeMFP@fourdigits\MFP@z@Frc@ii
493 \makeMFP@fourdigits\MFP@z@Frc@i
494 \makeMFP@fourdigits\MFP@z@Int@i
495 \edef\MFP@z@Int{\number\MFP@z@Int@ii\MFP@z@Int@i}%
496 \edef\MFP@z@Frc{\MFP@z@Frc@i\MFP@z@Frc@ii}%
497 \edef\MFP@z@Ovr{\number\MFP@z@Int@iii}%
498 \edef\MFP@z@Und{\MFP@z@Frc@iii\MFP@z@Frc@iv}%
499 \ifnum\MFP@z@Int>0
500 \else\ifnum\MFP@z@Frc>0
501 \else \def\MFP@z@Sgn{0}%
502 \fi\fi}%

```

For division, we will obtain the result one digit at a time until the 9th digit after the decimal is found. That 9th will be used to round to eight digits (and stored as underflow). We normalize the denominator by shifting left until the integer part is eight digits. We do the same for the numerator. The integer quotient of the integer parts will be one digit (possibly a 0). If the denominator is shifted  $d$  digits left and the numerator  $n$  digits left, the quotient will have to be shifted  $n - d$  places right or  $d - n$  places left. Since the result is supposed to have 9 digits after the dot, our quotient needs  $9 + d - n + 1$  total digits. Since  $d$  can be as high as 15 and  $n$  as low as 0, we could need 25 repetitions. However, that extreme would put 15 or 16 digits in the integer part, a 7 or 8 digit overflow. (It can be argued that only 16 significant digits should be retained in any case.) If  $d$  is 0 and  $n$  is 15 we would need  $-5$  digits. That means the first nonzero digit is in the 15th or 16th place after the dot and the quotient is effectively zero.

Here I explain why we normalize the parts in this way. If a numerator has the form  $n_1.n_2$  and the denominator has the form  $d_1.d_2$  then  $\text{\TeX}$  can easily obtain the integer part of  $n_1/d_1$ , because these are within its range for integers. The resulting quotient (let's call it  $q_1$ ) is the largest integer satisfying  $q_1 d_1 \leq n_1$ . What we seek, however is the largest integer  $q$  such that  $q(d_1.d_2) \leq n_1.n_2$ . It can easily be shown that  $q \leq q_1$ . It is true, but not so easily shown, that  $q \geq q_1 - 1$ . This is only true if  $d_1$  is large enough, in our case it has to be at least five digits. Thus we only have to do one simple division and decide if we need to reduce the quotient by one. If we arrange for  $d_1$  to have eight digits, then  $q_1$  will be one digit and the test for whether we need to reduce it becomes easier.

This test is done as follows. The first trial quotient,  $q_1$ , will work if

$$q_1 d_1 (10)^8 + q_1 d_2 \leq n_1 (10)^8 + n_2$$

This means

$$0 \leq (n_1 - q_1 d_1)(10)^8 + n_2 - q_1 d_2. \quad (1)$$

Since  $d_2$  is no more than eight digits,  $q_1 d_2$  is less than  $9(10)^8$ . Inequality (1) is therefore satisfied if  $n_1 - q_1 d_1 \geq 9$ . If that is not the case then the right side of (1)

is computable within  $\text{\TeX}$ 's integer ranges and we can easily test the inequality. If the inequality holds, then  $q = q_1$ , otherwise  $q = q_1 - 1$ .

Note also that when  $q = q_1$ , then both terms in (1) (ignoring the  $10^8$  factor) will be needed to calculate the remainder. If  $q = q_1 - 1$ , we simply add  $d_1$  and  $d_2$  to the respective parts. Thus we will save these values for that use.

Now I need to get it organized.  $\text{\MFP@Rdiv}$  will have  $\text{\MFP@x@*}$  and  $\text{\MFP@y@*}$  available. One step (could be first or last). Is to calculate the sign. Let's do it first (because we need to check for zero anyway).

We invoke an error message upon division by zero, but nevertheless return a value. By default it is 0 for  $0/0$  and the maximum possible real for  $x/0$  when  $x$  is not zero. If the numerator is zero and the denominator not, we do nothing as  $z$  was initialized to be zero.

If neither is zero, we calculate the sign of the result and call  $\text{\MFP@@Rdiv}$  to divide the absolute values.

```

503 \def\MFP@Rdiv{%
504   \ifnum\MFP@y@Sgn=0 \MFP@dividebyzero@err
505   \ifnum\MFP@x@Sgn=0
506     \edef\MFP@z@Int{\ZeroOverZeroInt}%
507     \edef\MFP@z@Frc{\ZeroOverZeroFrc}%
508   \else
509     \edef\MFP@z@Int{\xOverZeroInt}%
510     \edef\MFP@z@Frc{\xOverZeroFrc}%
511   \fi
512   \edef\MFP@z@Sgn{\MFP@x@Sgn}%
513   \else\ifnum\MFP@x@Sgn=0 \MFP@Rzero
514   \else \edef\MFP@z@Sgn{\number\MFP@x@Sign\MFP@y@Sgn}\MFP@@Rdiv
515   \fi\fi}%

```

Now we have two positive values to divide. Our first step is to shift the denominator ( $y$ ) left and keep track of how many places. We store the shift in  $\text{\MFP@tempa}$ . This actually changes the value of  $y$ , but knowing the shift will give us the correct quotient in the end.

We first arrange that  $\text{\MFP@y@Int}$  is nonzero by making it  $\text{\MFP@y@Frc}$  if it is zero (a shift of eight digits). Then the macro  $\text{\MFP@numdigits@toshift}$  computes 8 minus the number of digits in  $\text{\MFP@y@Int}$ , which is how many positions left  $y$  will be shifted. We then call  $\text{\MFP@doshift@y}$  on the concatenation of the digits in the integer and fractional parts (padded with zeros to ensure there are at least 16). All this macro does is read the first eight digits into  $\text{\MFP@y@Int}$  and the next eight into  $\text{\MFP@y@Frc}$ .

```

516 \def\MFP@@Rdiv{%
517   \ifnum\MFP@y@Int=0
518     \edef\MFP@y@Int{\number\MFP@y@Frc}%
519     \def\MFP@y@Frc{00000000}%
520     \MFP@tempa=8
521   \else
522     \MFP@tempa=0
523   \fi
524   \advance\MFP@tempa\MFP@numdigits@toshift\MFP@y@Int\relax

```

```
525 \@XP\MFP@doshift@y\@xp\MFP@y@Int\MFP@y@Frc0000000\MFP@end
```

We repeat all that on the numerator  $x$ , except shifting its digits left means the final outcome will need a corresponding *right* shift. We record that fact by reducing `\MFP@tempa`, which ends up holding the net shift necessary.

This has the advantage that we know the result will be in the range  $[0.1, 10)$ . It also means we can reduce the number of places we will need to shift left as well as reduce the number of iterations of the loop that calculates the digits.

```
526 \ifnum\MFP@x@Int=0
527   \edef\MFP@x@Int{\number\MFP@x@Frc}%
528   \def\MFP@x@Frc{00000000}%
529   \advance\MFP@tempa -8
530 \fi
531 \advance\MFP@tempa-\MFP@numdigits@toshift\MFP@x@Int\relax
532 \@XP\MFP@doshift@x\@xp\MFP@x@Int\MFP@x@Frc0000000\MFP@end
```

Since our result will have at most one digit in the integer part, a rightward shift of 10 places will make every digit 0, including the rounding digit, so we do nothing (returning 0).

```
533 \ifnum\MFP@tempa<-9
534 \else
```

Now we perform the division, which is a loop repeated  $10 + \text{\MFP@tempa}$  times. Therefore, we add 10 to `\MFP@tempa` in `\MFP@tempf`, our loop counter. We also initialize the macro that will store the digits and then, after the division, shift and split it into parts.

```
535 \MFP@tempf\MFP@tempa
536 \advance\MFP@tempf 10
537 \def\MFP@z@digits{}%
538 \MFP@Rdivloop
539 \MFP@shiftandsplit@z@digits
```

The last remaining step is to round and carry and get the fractional part in the appropriate 8-digit form..

```
540 \ifnum\MFP@z@Und>4
541   \MFP@addone\MFP@z@Frc1\MFP@z@Frc
542   \MFP@carrya\MFP@z@Frc\MFP@z@Int
543   \MFP@carrya\MFP@z@Int\MFP@z@Ovr
544   \makeMFP@eightdigits\MFP@z@Frc
545 \fi
546 \fi}%
```

If #1 of `\MFP@numdigits@toshift`, has  $n$  digits then `\MFP@numdigits@toshift` picks out the value  $8 - n$ . `\MFP@doshift@x` reads the first eight digits into `\MFP@x@Int` and then pulls out eight more from the rest (#9) inside `\MFP@x@Frc`. The same with `\MFP@doshift@y`.

```
547 \def\MFP@numdigits@toshift#1{\@xp\MFP@ninthofmany#101234567\MFP@end}%
548 \def\MFP@doshift@x#1#2#3#4#5#6#7#8#9\MFP@end{%
549   \def\MFP@x@Int{#1#2#3#4#5#6#7#8}%
550   \edef\MFP@x@Frc{\MFP@eightofmany#9\MFP@end}}%
```



```

551 \def\MFP@doshift@y#1#2#3#4#5#6#7#8#9\MFP@end{%
552 \def\MFP@y@Int{#1#2#3#4#5#6#7#8}%
553 \edef\MFP@y@Frc{\MFP@eightofmany#9\MFP@end}}%

```

The loop counter is \MFP@tempf, \MFP@tempa is reserved for the shift required later, the quotient digit will be \MFP@tempb. The remainder will be calculated in \MFP@tempc and \MFP@tempd. \MFP@tempe will hold the value whose size determines whether the quotient needs to be reduced.

```

554 \def\MFP@Rdivloop{%
555 \MFP@tempb\MFP@x@Int          % \MFP@tempb = n_1
556 \MFP@tempc\MFP@y@Int          % \MFP@tempc = d_1
557 \divide\MFP@tempb \MFP@tempc  % \MFP@tempb = n_1/d_1 = q_1
558 \multiply \MFP@tempc \MFP@tempb % \MFP@tempc = q_1 d_1
559 \MFP@tempd \MFP@y@Frc          % \MFP@tempd = d_2
560 \multiply \MFP@tempd \MFP@tempb % \MFP@tempd = q_1 d_2
561 \MFP@tempe \MFP@tempc
562 \advance \MFP@tempe -\MFP@x@Int\relax % \MFP@tempe = -n_1 + q_1 d_1
563 \ifnum \MFP@tempe > -9          % n_1 - q_1 d_1 < 9
564 \multiply \MFP@tempe\MFP@tteight % -(n_1 - q_1 d_1)(10)^8
565 \advance \MFP@tempe \MFP@tempd   % add q_1 d_2
566 \advance \MFP@tempe -\MFP@x@Frc\relax % add -n_2
567 \ifnum \MFP@tempe>0             % Crucial inequality fails
568 \advance\MFP@tempb -1           % new q = q_1 - 1
569 \advance\MFP@tempc -\MFP@y@Int  % q_1 d_1 - d_1 = q d_1
570 \advance\MFP@tempd -\MFP@y@Frc\relax % q_1 d_2 - d_2 = q d_2
571 \fi
572 \fi
573 \edef\MFP@z@digits{\MFP@z@digits\number\MFP@tempb}%

```

It remains to:

- Do the carry from \MFP@tempd to \MFP@tempc. Then \MFP@tempc.\MFP@tempd will represent  $q \cdot y$ .
- Subtract them from \MFP@x@Int and \MFP@x@Frc (i.e. remainder =  $x - qy$ ).
- Borrow, if needed, and we will have the remainder in \MFP@x@Int.\MFP@x@Frc.

Then we decrement the loop counter, and decide whether to repeat this loop. If so, we need to shift the remainder right one digit (multiply by 10). We don't use \MFP@carrya since it requires macros; its internal code, \MFP@carryi just leaves the results in \MFP@tempa.\MFP@tempb.

```

574 \begingroup
575 \MFP@carryi\MFP@tteight\MFP@tempd\MFP@tempc
576 \MFP@endgroup@after{%
577 \MFP@tempc=\number\MFP@tempa
578 \MFP@tempd=\number\MFP@tempb\relax
579 }%
580 % subtract
581 \MFP@addone\MFP@x@Int{-\MFP@tempc}\MFP@x@Int
582 \MFP@addone\MFP@x@Frc{-\MFP@tempd}\MFP@x@Frc
583 % borrow
584 \ifnum\MFP@x@Frc<0

```

```

585 \MFP@addone\MFP@x@Frc\MFP@ttteight\MFP@x@Frc
586 \MFP@addone\MFP@x@Int{-1}\MFP@x@Int
587 \fi
588 \advance\MFP@tempf -1
589 \ifnum\MFP@tempf>0
590 \edef\MFP@x@Int{\MFP@x@Int0}%
591 \edef\MFP@x@Frc{\MFP@x@Frc0}%
592 \MFP@carrya\MFP@x@Frc\MFP@x@Int
593 \@xp\MFP@Rdivloop
594 \fi}%

```

Now `\MFPshiftandsplit@z@digits`. At this point, the digits of the quotient are stored in `\MFP@z@digits`. We need to shift the decimal `\MFP@tempa` places left, and perform the rounding. There are `\MFP@tempa + 10` digits. This could be as little as 1 or as great as 25. In the first case `\MFP@tempa` is  $-9$ , and this (rightward) shift produces 0 plus a rounding digit. In the latter case `\MFP@tempa` is 15, and the shift produces 8 digits overflow, an 8-digit integer part, an 8-digit fractional part and a rounding digit. In the example 0123456, `\MFP@tempa + 10` is 7, so `\MFP@tempa` is  $-3$ . The shift produces 0.0001 2345 6. The rounding digit (6) makes the answer 0.0001 2346.

We take two cases:

- `\MFP@tempa  $\leq$  7`, prepend  $7 - \text{\MFP@tempa}$  zeros. The first 8 digits will become the integer part, and there should be exactly 9 more digits.
- `\MFP@tempa  $>$  7`, pluck `\MFP@tempa - 7` digits for overflow, the next 8 for integer part, leaving 9 more digits

In either case, the 9 last digits will be processed into a fractional part (with possible carry if the rounding increases it to  $10^8$ ).

After this, we will return to `\MFP@Rdiv` so overwriting `\MFP@temp*` won't cause any problems.

```

595 \def\MFPshiftandsplit@z@digits{%
596 \advance \MFP@tempa -7
597 \ifnum\MFP@tempa>0
598 \def\MFP@z@Ovr{}%
599 \@xp\MFPget@Ovrdigits\MFP@z@digits\MFP@end
600 \else
601 \ifnum\MFP@tempa<-7
602 \edef\MFP@z@digits{00000000\MFP@z@digits}%
603 \advance\MFP@tempa8
604 \fi
605 \ifnum\MFP@tempa<-3
606 \edef\MFP@z@digits{0000\MFP@z@digits}%
607 \advance\MFP@tempa4
608 \fi
609 \edef\MFP@z@digits{%
610 \ifcase-\MFP@tempa\or
611 0\or
612 00\or
613 000\or

```

```

614         0000\else
615         00000%
616         \fi \MFP@z@digits}%
617     \xp\MFPget@Intdigits\MFP@z@digits\MFP@end
618     \fi}%

```

The macro `\MFPget@Ovrdigits` is a loop that loads the first `\MFP@tempa` digits of what follows into `\MFP@z@Ovr`. It does this one digit (`#1`) at a time. Once the counter reaches 0, we call the macro that processes the integer part digits.

```

619 \def\MFPget@Ovrdigits#1{%
620     \edef\MFP@z@Ovr{\MFP@z@Ovr#1}%
621     \advance\MFP@tempa -1
622     \ifnum\MFP@tempa>0
623         \xp\MFPget@Ovrdigits
624     \else
625         \xp\MFPget@Intdigits
626     \fi}%

```

The macro `\MFPget@Intdigits` should have exactly 17 digits following it. It puts eight of them in `\MFP@z@Int`, then calls `\MFPget@Frcdigits` to read the fractional part. That requires exactly nine digits follow it, putting eight in `\MFP@z@Frc` and the last in `\MFP@z@Und`. Still, to allow a graceful exit should there be more, we gobble the rest of the digits.

```

627 \def\MFPget@Intdigits#1#2#3#4#5#6#7#8{%
628     \def\MFP@z@Int{\number#1#2#3#4#5#6#7#8}%
629     \MFPget@Frcdigits}%
630 \def\MFPget@Frcdigits#1#2#3#4#5#6#7#8#9{%
631     \def\MFP@z@Frc{\number#1#2#3#4#5#6#7#8}%
632     \def\MFP@z@Und{\number#9}\gobbleto\MFP@end}%

```

The max and min operations simply run the compare operation and use and use the resultant booleans to copy  $x$  or  $y$  to  $z$ .

```

633 \def\MFP@Rmax{%
634     \MFP@Rcmp \ifMFP@neg \MFP@Rcopy yz\else\MFP@Rcopy xz\fi}%
635 \def\MFP@Rmin{%
636     \MFP@Rcmp \ifMFP@pos \MFP@Rcopy yz\else\MFP@Rcopy xz\fi}%

```

### 3.5 Commands to format for printing

`\MFPtruncate` This first runs the parsing command so the fractional part has exactly eight digits. These become the arguments of `\MFP@@Rtrunc`, which just keeps the right number. For negative truncations we prepend zeros to the integer part so it too is exactly eight digits. These become the arguments of `\MFP@@iRtrunc`, which substitutes 0 for the last `-\MFP@tempa` of them.

The macro to store the result in follows `#2`. It is read and defined by either `\MFP@Rtrunc` or `\MFP@iRtrunc`.

```

637 \def\MFPtruncate#1#2{%
638     \begingroup
639         \MFP@tempa#1\relax

```

```

640 \MFPparse@x{#2}%
641 \ifnum\MFP@tempa<1
642 \xp\MFP@iRtrunc
643 \else
644 \xp\MFP@Rtrunc
645 \fi}%
646 \def\MFP@Rtrunc#1{%
647 \edef\MFP@x@Frc{\xp\MFP@@Rtrunc\MFP@x@Frc\MFP@end}%
648 \ifnum\MFP@x@Int=0
649 \ifnum\MFP@x@Frc=0
650 \def\MFP@x@Sgn{0}%
651 \fi
652 \fi
653 \MFP@endgroup@after{%
654 \MFP@afterdef#1{\MFP@x@Sgn\MFP@x@Int.\MFP@x@Frc}}}%
655 \def\MFP@@Rtrunc#1#2#3#4#5#6#7#8#9\MFP@end{%
656 \ifcase\MFP@tempa\or
657 #1\or
658 #1#2\or
659 #1#2#3\or
660 #1#2#3#4\or
661 #1#2#3#4#5\or
662 #1#2#3#4#5#6\or
663 #1#2#3#4#5#6#7\else
664 #1#2#3#4#5#6#7#8\fi}%
665 \def\MFP@iRtrunc#1{%
666 \makeMFP@eightdigits\MFP@x@Int
667 \edef\MFP@x@Val{\number\MFP@x@Sgn\xp\MFP@@iRtrunc\MFP@x@Int\MFP@end}%
668 \MFP@endgroup@after{\MFP@afterdef#1{\MFP@x@Val}}}%
669 \def\MFP@@iRtrunc#1#2#3#4#5#6#7#8#9\MFP@end{%
670 \ifcase-\MFP@tempa
671 #1#2#3#4#5#6#7#8\or
672 #1#2#3#4#5#6#70\or
673 #1#2#3#4#5#600\or
674 #1#2#3#4#5000\or
675 #1#2#3#40000\or
676 #1#2#300000\or
677 #1#2000000\or
678 #10000000\else
679 00000000\fi}%

```

**\MFPround** For rounding we simply add the appropriate fraction and truncate. The macro in which to store the result will follow #2, and be picked up by the **\MFPtruncate** command.

```

680 \def\MFPround#1#2{%
681 \begingroup
682 \MFP@tempa#1\relax
683 \ifnum 0>\MFP@tempa
684 \edef\MFP@y@Tmp{%
685 \ifcase-\MFP@tempa\or

```

```

686         5\or
687         50\or
688         500\or
689         5000\or
690         50000\or
691         500000\or
692         5000000\else
693         50000000\fi
694     }%
695 \else
696     \edef\MFP@y@Tmp{%
697         \ifcase\MFP@tempa
698         .5\or
699         .05\or
700         .005\or
701         .0005\or
702         .00005\or
703         .000005\or
704         .0000005\or
705         .00000005\else
706         0\fi
707     }%
708 \fi
709 \MFPchk{#2}\ifMFP@neg\edef\MFP@y@Tmp{-\MFP@y@Tmp}\fi
710 \MFPadd{#2}\MFP@y@Tmp\MFP@z@Tmp
711 \MFP@endgroup@after{\MFP@afterdef\MFP@z@Tmp{\MFP@z@Tmp}}%
712 \MFPtruncate{#1}\MFP@z@Tmp}%

```

**\MFPstrip** Stripping zeros from the right end of the fractional part. The star form differs only in the handling of a zero fractional part. So we check whether it is zero and when it is, we either append ‘.0’ or nothing. The rest of the code grabs a digit at a time and stops when the rest are zero.

```

713 \def\MFPstrip{%
714     \@ifstar{\MFP@strip{}}{\MFP@strip{.0}}}%
715 \def\MFP@strip#1#2#3{%
716     \MFPparse@x{#2}%
717     \ifnum \MFP@x@Frc=0
718         \edef#3{\MFP@x@Sign\MFP@x@Int#1}%
719     \else
720         \edef#3{\MFP@x@Sign\MFP@x@Int.\@xp\MFP@@strip\MFP@x@Frc\MFP@end}%
721     \fi}%
722 \def\MFP@@strip#1#2\MFP@end{%
723     #1%
724     \ifnum 0#2>0
725         \@xp\MFP@@strip
726     \else
727         \@xp\gobbleto\MFP@end
728     \fi#2\MFP@end}%

```

### 3.6 Miscellaneous

Here is the code that allows definitions to survive after `\stopMFPprogram`. The `\Global` variants are easiest.

```
729 \def\MFP@Global#1{\toks@{\xp{#1}\xdef#1{\the\toks@}}}%
730 \def\MFP@GlobalStack{\MFP@Global\MFP@Rstack}%
```

The `\Export` command adds the command and its definition to a macro that is executed after the closing group of the program.

```
731 \def\MFP@Export#1{%
732   \begingroup
733   \toks@{\xp{\MFPprogram@returns}}%
734   \MFP@endgroup@after{%
735     \MFP@afterdef\MFPprogram@returns{\the\toks@ \MFP@afterdef#1{#1}}%
736   }%
737 \def\MFP@ExportStack{\MFP@Export\MFP@Rstack}%
```

The various operations `\MFP@R...` together make up a “microcode” in terms of which the stack language and the operand language are both defined. As a language in its own right, it lacks only convenient ways to move numbers around, as well as a few extra registers for saving intermediate results. In this language, numbers are represented by a three part data structure, consisting of a signum, an integer part and a fractional part.

Here we define extra commands to remedy this lack, starting with a way to load a number (or rather, a three part data structure representing a number) directly into a register. Here `#1` is a register name (we always use a single letter) and the remaining arguments are the signum, the integer part and the fractional part (automatically normalized to 8 digits). The “register” is just a set of three macros created from the name given.

We make loading a number into a register a little more general than strictly needed, allowing the parts to be specified as anything `TEX` recognizes as a number and allowing any register name. This generality might reduce efficiency but it simplifies code. Because register `z` is by far the most common one to load, we make more efficient version of it.

```
738 \def\MFP@Rload #1#2#3#4{%
739   \xp\edef\csname MFP@#1@Sgn\endcsname{\number#2}%
740   \xp\edef\csname MFP@#1@Int\endcsname{\number#3}%
741   \xp\edef\csname MFP@#1@Frc\endcsname{\number#4}%
742   \xp\makeMFP@eightdigits\csname MFP@#1@Frc\endcsname}%
743 \def\MFP@Rcopy#1#2{%
744   \MFP@Rload #2{\csname MFP@#1@Sgn\endcsname}%
745               {\csname MFP@#1@Int\endcsname}%
746               {\csname MFP@#1@Frc\endcsname}}%
747 \def\MFP@Rloadz#1#2#3{%
748   \edef\MFP@z@Sgn{\number#1}%
749   \edef\MFP@z@Int{\number#2}%
750   \edef\MFP@z@Frc{\number#3}%
751   \makeMFP@eightdigits\MFP@z@Frc}%
```

`\MFPpi`      These are some miscellaneous constants. The 8-digit approximation to  $\pi$ , is

`\MFPe` `\MFPhi` and the constant mathematicians call  $e$  is `\MFPe`. Finally, the golden ratio `\MFPhi` (often called  $\phi$ ) is obtained by `\MFPhi`.

```

752 \def\MFPhi{3.14159265}%
753 \def\MFPe{2.71828183}%
754 \def\MFPhi{1.61803399}%
Load (conditionally) mfpextra.tex.
755 \MFPe@loadextra
756 \MFPe@finish
757 </sty>

```

## 4 Extras

The extras consist so far of sine, cosine, angle, logarithm, powers, and square root. For completeness, here is the table of user-level commands available.

*Operand versions*

Command	operation
<code>\MFPe@sin{&lt;num&gt;}\macro</code>	Stores $\sin(\langle num \rangle)$ in <code>\macro</code> , where $\langle num \rangle$ is an angle in degrees.
<code>\MFPe@cos{&lt;num&gt;}\macro</code>	Stores $\cos(\langle num \rangle)$ in <code>\macro</code> , where $\langle num \rangle$ is an angle in degrees.
<code>\MFPe@angle{&lt;x&gt;}{&lt;y&gt;}\macro</code>	Stores in <code>\macro</code> the polar angle coordinate $\theta$ of the point $(x, y)$ , where $-180 < \theta \leq 180$ .
<code>\MFPe@rad{&lt;num&gt;}\macro</code>	The angle $\langle num \rangle$ in degrees is converted to radians, and result is stored in <code>\macro</code> .
<code>\MFPe@deg{&lt;num&gt;}\macro</code>	The angle $\langle num \rangle$ in radians is converted to degrees, and result is stored in <code>\macro</code> .
<code>\MFPe@log{&lt;num&gt;}\macro</code>	Stores $\log(\langle num \rangle)$ in <code>\macro</code> (base 10 logarithm).
<code>\MFPe@ln{&lt;num&gt;}\macro</code>	Stores $\ln(\langle num \rangle)$ in <code>\macro</code> (natural logarithm).
<code>\MFPe@exp{&lt;num&gt;}\macro</code>	Stores $\exp(\langle num \rangle)$ (i.e., $e^x$ ) in <code>\macro</code> .
<code>\MFPe@sqrt{&lt;num&gt;}\macro</code>	Stores the square root of $\langle num \rangle$ in <code>\macro</code> .
<code>\MFPe@pow{&lt;num&gt;}{&lt;int&gt;}\macro</code>	Stores the $\langle int \rangle$ power of $\langle num \rangle$ in <code>\macro</code> . The second operand must be an integer (positive or negative).

Command	operation
<code>\Rsin</code>	The number is interpreted as degrees, and its sine is computed.
<code>\Rcos</code>	The number is interpreted as degrees, and its cosine is computed.
<code>\Rangle</code>	The top two numbers are interpreted as coordinates of a point $P$ in the order they were pushed. The polar angle coordinate $\theta$ of $P$ , with $-180 < \theta \leq 180$ is computed.
<code>\Rrad</code>	The number of degrees is converted to radians.
<code>\Rdeg</code>	The number of radians is converted to degrees.
<code>\Rlog</code>	Computes the base-10 logarithm.
<code>\Rln</code>	Computes the natural logarithm.
<code>\Rexp</code>	Computes the exponential of the number (i.e., $e^x$ ).
<code>\Rsqrt</code>	Computes the square root of the number.
<code>\Rpow</code>	Computes $x^y$ . The last number pushed ( $y$ ) must be an integer.

The user could easily convert between radians and degrees using multiplication and/or division. One could similarly convert between natural logarithms and base ten logarithms. The commands `\Rdeg`, `\Rrad`, `\Rlog` and `\Rln` (and their `\MFP...` counterparts) aim for more accurate results.

#### 4.1 Loading the extras

```

\Rsin We start mfpextra with the hook \MFP@Rextra that \startMFPprogram will call to
\Rcos make available the extra operations defined here. If minifp.sty has been loaded,
\Rangle this macro is \@empty, otherwise it should be undefined. If it is undefined we load
\Rrad minifp.sty. If it is then not \@empty we assume mfpextra.tex was previously
\Rdeg loaded and end input here.
\Rlog 758 \<extra>
\Rln 759 % check if mfpextra already loaded:
\Rexp 760 \expandafter\ifx\csname MFP@xfinish\endcsname\relax
\Rsqrt 761 \else \expandafter\endinput\fi
\Rpow 762 \expandafter\edef\csname MFP@xfinish\endcsname{%
763   \catcode64=\the\catcode64 \space
764   \catcode46=\the\catcode46 \space
765   \catcode60=\the\catcode60 \space
766   \catcode62=\the\catcode62 \space}%
767 \catcode64=11 % @
768 \catcode46=12 % . (period)
769 \catcode60=12 % <
770 \catcode62=12 % >
771 \ifx\MFP@Rextra\UndEfInEd \input minifp.sty \fi
772 \ifx\MFP@Rextra\@empty
773 \else
774   \immediate\write16{mfpextra.tex: already loaded.^^J}%
775   \MFP@xfinish
776   \expandafter\endinput
777 \fi

```



```

778 \immediate\write16{%
779     mfpextra.tex: extra operations for MiniFP package.^^J}%
780 \def\MFP@Rextra{%
781     \def\Rcos    {\MFP@stack@Unary\MFP@Rcos }%
782     \def\Rsin    {\MFP@stack@Unary\MFP@Rsin }%
783     \def\Rangle  {\MFP@stack@Binary\MFP@Rangle}%
784     \def\Rad     {\MFP@stack@Unary\MFP@Rad }%
785     \def\Rdeg    {\MFP@stack@Unary\MFP@Rdeg }%
786     \def\Rlog    {\MFP@stack@Unary\MFP@Rlog }%
787     \def\Rln     {\MFP@stack@Unary\MFP@Rln }%
788     \def\Rexp    {\MFP@stack@Unary\MFP@Rexp }%
789     \def\Rsqrtn  {\MFP@stack@Unary\MFP@Rsqrtn}%
790     \def\Rpow    {\MFP@stack@Binary\MFP@Rpow}}%
\MFPsin      Then the wrappers for the operand versions.
\MFPcos      791 \def\MFPcos    {\MFP@op@Unary\MFP@Rcos }%
\MFPprad     792 \def\MFPsin    {\MFP@op@Unary\MFP@Rsin }%
\MFPdeg      793 \def\MFPangle  {\MFP@op@Binary\MFP@Rangle}%
\MFPlog      794 \def\MFPprad   {\MFP@op@Unary\MFP@Rad }%
\MFPln       795 \def\MFPdeg    {\MFP@op@Unary\MFP@Rdeg }%
\MFPexp      796 \def\MFPlog    {\MFP@op@Unary\MFP@Rlog }%
\MFPsqrtn    797 \def\MFPln     {\MFP@op@Unary\MFP@Rln }%
\MFPexp      798 \def\MFPexp    {\MFP@op@Unary\MFP@Rexp }%
\MFPpow      799 \def\MFPsqrtn  {\MFP@op@Unary\MFP@Rsqrtn}%
800 \def\MFPpow     {\MFP@op@Binary\MFP@Rpow}%

```

## 4.2 Error messages

These extra commands come with a few possible new warnings and errors.

`\LogOfZeroInt`      Trying to take the logarithm of zero will result in an error message. If one  
`\LogOfZeroFrac`    allows  $\TeX$  to continue, the returned value will be negative, with an integer part  
whose absolute value is equal to the contents of `\LogOfZeroInt` and a fractional  
part equal to the contents of `\LogOfZeroFrac`. The defaults are both 99999999.

Trying to take the logarithm of a negative number will produce the warning

MFP warning: Log of a negative number is complex.  
Only the real part will be computed.

The log of the absolute value is returned.

Trying to take the square root of a negative number has similar behavior. It produces a warning and returns 0.

Trying to take the exponential of a number larger than about 18.42 will cause an error and the number returned has integer part 99999999 and fractional part 99999999.

Trying to take a negative power of 0 produces an error and returns the same value as trying to divide 1 by 0.

Messages for errors related to impossible powers and logarithms.

```

801 \def\MFP@logofzero@err{%
802     \MFP@errmsg{logarithm of zero}%
803     {You tried to take the logarithm of zero. What were you %

```

```

804     thinking? If you ^^Jcontinue, the value %
805     assigned will be -\LogOfZeroInt.\LogOfZeroFrac.}}%
806 \def\LogOfZeroInt {\MaxRealInt}%
807 \def\LogOfZeroFrac{\MaxRealFrac}%
808 \def\MFP@expoverflow@err{%
809   \MFP@errmsg{Power too large}%
810   {The power you tried to calculate is too large for %
811     8 digits. If you continue, ^^Jthe value assigned will be %
812     \MaxRealInt.\MaxRealFrac.}}%
813 \def\MFP@badpower@err{%
814   \MFP@errmsg{negative power of zero}%
815   {You tried to take a negative power of zero. What were you
816     thinking? If you ^^Jcontinue, the value assigned will be %
817     \xOverZeroInt.\xOverZeroFrac.}}%

```

A debugging utility, \MFPshowreg displays the contents of a register.

```

818 \def\MFPshowreg #1{%
819 \ifMFPdebug
820 \begingroup
821   \edef\theregister{%
822     #1 = \expandafter \MFP@Sign
823         \csname MFP@#1@Sgn\endcsname %
824         \csname MFP@#1@Int\endcsname.%
825         \csname MFP@#1@Frc\endcsname}%
826   \show\theregister
827 \endgroup
828 \fi}%

```

### 4.3 Sine and Cosine

For iterated code, the most common register to copy is  $z$  and the most common place to copy it is to  $x$  or  $y$  so we make single commands to do those.

```

829 \def\MFP@Rcopyz#1{\MFP@Rload {#1}\MFP@z@Sgn\MFP@z@Int\MFP@z@Frc}%
830 \def\MFP@Rcopyzx{\MFP@Rcopyz x}%
831 \def\MFP@Rcopyzy{\MFP@Rcopyz y}%

```

Our code assumes the number  $x$  is an angle in degrees. To get sine and cosine of numbers as radians, simply convert your radians to degrees using \MFPdeg or \Rdeg. Then find the sine or cosine of the result. For example, if \X holds the angle in radians and you want the result to be stored in \S:

```
\MFPdeg\X\Y \MFPsin\Y\S
```

For unit conversions such as radian to degree we try to be more accurate than a multiplication by an eight-digit conversion factor allows. If  $x$  is large and the factor is off by  $0.5 \times 10^{-8}$ , then the result can be significantly off. But if we are able to give the conversion factor 16 digits precision, then only the imprecision of  $x$  will significantly affect the result.

We express the conversion factor as an integer part and two eight-digit fractional parts. We multiply  $x$  by the integer and first fractional part (#1 and #2) with a normal \MFP@Rmul, but we save the underflow digits and undo the rounding that

occured at the 8th digit. Together these give us an essentially exact result. Then we multiply  $x$  by the second fractional part (#3) and add the saved underflow to the result. Finally, we round and add the result to the first product. Argument #3, as well as the underflow digits, represent numbers less than  $10^{-8}$ , so we effectively scale them up by  $10^8$ , round the result to an integer and scale that back down.

The registers  $w$  and  $v$  are used to save intermediate results. The “DP” in  $\backslash\text{MFP@DPmul}$  refers to the fact that we are multiplying by a “double precision” real.

```

832 \def\MFP@DPmul#1#2#3{%
833   \ifnum\MFP@x@Sgn=0
834     \MFP@Rzero
835   \else
836     \MFP@Rcopy xv%
837     \MFP@Rload y1{#1}{#2}\MFP@Rmul
838     \edef\MFP@w@Und{\MFP@z@Und}%
839     \ifnum\MFP@z@Frc@iii>4999
840       \MFP@tempa\MFP@z@Frc \advance\MFP@tempa-1
841       \edef\MFP@z@Frc{\number\MFP@tempa}%
842       \makeMFP@eightdigits\MFP@z@Frc
843     \fi
844     \MFP@Rcopyz w%
845     \MFP@Rcopy vx\MFP@Rload y10{#3}\MFP@Rmul
846     \MFP@Rcopyzx\MFP@Rload y\MFP@v@Sgn 0{\MFP@w@Und}\MFP@Radd
847     \MFP@tempa\MFP@z@Int\relax
848     \ifnum\MFP@z@Frc<50000000 \else \advance\MFP@tempa 1 \fi
849     \ifnum\MFP@tempa<\MFP@ttteight\relax
850       \MFP@Rload x{\ifnum\MFP@tempa>0 \MFP@z@Sgn\else0\fi}0\MFP@tempa
851     \else
852       \MFP@Rload x\MFP@z@Sgn10%
853     \fi
854     \MFP@Rcopy wy\MFP@Radd
855   \fi}%

```

Conversion factors:

- radians to degrees: 57.2957795130823209
- degrees to radians: 0.0174532925199433
- natural log to common log: 0.4342944819032518
- common log to natural log: 2.3025850929940457

Note that the comparatively large size of the first number means that the  $\pm 0.5 \cdot 10^{-8}$  imprecision that  $x$  implicitly carries will be multiplied to approximately  $\pm 29.6 \cdot 10^{-8}$  in the result. The only way around this would be to operate with higher precision internally. We do that in the code for computing angles.

```

856 \def\MFP@Rdeg{\MFP@DPmul{57}{29577951}{30823209}}%
857 \def\MFP@Rrad{\MFP@DPmul{0}{01745329}{25199433}}%
858 \def\MFP@RbaseX{\MFP@DPmul{0}{43429448}{19032518}}%
859 \def\MFP@RbaseE{\MFP@DPmul{2}{30258509}{29940457}}%

```

There are very few angles that are expressible in eight digits whose sine or cosine can be expressed exactly in eight digits. For these, we do obtain an exact

result. Other values produce inexact results. It would be nice if we could at least obtain these correctly rounded to eight decimals, but unfortunately our methods will often produce a result off by 1 in the eighth decimal from the correctly rounded value. Anything that involves the addition of two or more rounded results can have this problem. The only way to get correctly rounded results is to carry out all operations internally to additional places. Even then, there will be the occasional .4999... that should round to 0 but rounds to 1 instead.

For the cosine, just compute  $\sin(90 - x)$ .

```
860 \def\MFP@Rcos{%
861   \MFP@Rcopy xy\MFP@Rload x1{90}0\MFP@Rsub
862   \MFP@Rcopyzx\MFP@Rsin}%

```

Reduce  $|x|$  by subtracting 180 from the integer part until it is less than 180. Of course,  $\sin x = \operatorname{sgn}(x) \sin(|x|)$  so we only need to compute  $\sin(|x|)$ . The sign will be that of  $x$ ; each reduction by 180 changes the sign, but the reduction code keeps track of that. If  $x$  is 0 after the reduction, return zero.

```
863 \def\MFP@Rsin{%
864   \MFP@tempa\MFP@x@Int
865   \MFP@tempb\MFP@x@Frc
866   \MFP@tempc\MFP@x@Sgn\relax
867   \MFP@reduce@angle
868   \ifnum\MFP@tempa>0 \MFP@@Rsin
869   \else\ifnum\MFP@tempb>0 \MFP@@Rsin
870   \else \MFP@Rzero
871   \fi\fi}%

```

This following reduces  $|x|$  to the case  $0 \leq |x| < 180$ . It assumes the integer part is in count register `\MFP@tempa`, the sign in `\MFP@tempc`.

```
872 \def\MFP@reduce@angle{%
873   \ifnum\MFP@tempa<180
874   \else
875     \advance\MFP@tempa-180
876     \MFP@tempc-\MFP@tempc
877     \@xp\MFP@reduce@angle
878   \fi}%

```

At this point,  $|x|$  is represented by `\MFP@tempa` (integer part) and `\MFP@tempb` (fractional part). Also, we already know the sign stored in `\MFP@tempc`. Moreover  $0 < \MFP@tempa < 180$ . We now reduce to  $0 < |x| \leq 90$  using  $\sin(x) = \sin(180 - |x|)$ , and return 1 if equal to 90.

The calculation of  $180 - x$  is optimized, taking advantage of the fact that both  $x$  and the result are known to be positive. If the fractional part is positive, we borrow 1 by reducing 180 to 179.

```
879 \def\MFP@@Rsin{%
880   \ifnum\MFP@tempa<90
881   \else
882     \MFP@tempa -\MFP@tempa
883     \ifnum\MFP@tempb>0
884       \MFP@tempb -\MFP@tempb

```

```

885     \advance\MFP@tempb \MFP@ttteight\relax
886     \advance\MFP@tempa 179
887     \else \advance\MFP@tempa 180
888     \fi
889 \fi
890 \ifnum\MFP@tempa=90
891     \MFP@Rloadz \MFP@tempc10%
892 \else

```

We would need to convert  $x$  to radians (multiply by  $\pi/180$ ) to use the standard power series, but instead we will incorporate the conversion factor into the power series coefficients.

We will, however, try to increase accuracy by reducing the size of  $x$  and correspondingly increasing the appropriate factors. Since the number of significant figures of a product is limited by the least number of significant figures of the two factors, the bottleneck on accuracy is that of the smaller term: all our numbers have eight digits so if a number is small, the number of nonzero digits is small.

Dividing by 100 seems a good choice (so our units are “hectodegrees”). This makes  $0 < x < .9$  and the integer part ( $\backslash\text{MFP@tempa}$ ) will be henceforth ignored.

The addition of 50 is for rounding purposes. After that, our computations amount to concatenating the top six digits of  $\backslash\text{MFP@tempb}$  to the digits of  $\backslash\text{MFP@tempa}$ . This will produce the integer form of the fractional part of  $x/100$  (the integer part of  $x/100$  is zero).

Division by 100 can turn a number into 0. This is one place we can lose accuracy (up to  $\pm 1$  in the last digit of the result). In compensation, the rest of the calculations become very much more accurate.

```

893     \advance\MFP@tempb 50 \divide\MFP@tempb 100
894     \multiply\MFP@tempa 1000000 \advance\MFP@tempb\MFP@tempa
895     \ifnum\MFP@tempb=0
896         \MFP@Rzero
897     \else

```

We save some multiplications by working with  $t = x^2$ . As we don’t need the original  $x$  anymore, we simply replace it with the newly reduced value. We also save this reduced  $x$  in another register,  $s$ , as we will need it again at the end, and our intermediate calculations do not preserve the  $x$  register. Then we square  $x$  and, if that square is 0 we can skip all the power series and simply return  $x$  converted to radians. If  $x^2$  is not zero, we save it in temporary register  $t$  and call our power series. When this program is finished, all that remains is the final multiplication by a conversion factor ( $\backslash\text{MFP@DPmul}$ ).

```

898     \MFP@Rload s\MFP@tempc0\MFP@tempb
899     \MFP@Rcopy sx%
900     \MFP@Rsqr
901     \ifnum \MFP@z@Frc>0
902         \MFP@Rcopyz t\MFP@Rsin@prog
903     \else
904         \MFP@Rcopy sx%
905     \fi

```

```

906      \MFP@DPmul 1{74532925}{19943296}%
907      \fi
908      \fi}%

```

\MFP@Rsin@prog is the power series computation. The power series need only go to the  $x^{13}$  term as the next is less than  $10^{-9}$  and in our 8-place computations is indistinguishable from 0. Our series is:

$$rx(1 - r^2t/3! + r^4t^2/5! - r^6t^3/7! + r^8t^4/9! - r^{10}t^5/11! + r^{12}t^6/13!)$$

where  $r$  is the factor that converts  $x$  to radian measure (hctodegrees to radians). When  $x$  is so small as to produce  $t = 0$  we have skipped all this.

We minimize any multiplications of tiny numbers by computing this as

$$rx(1 - ft(1 - et(1 - dt(1 - ct(1 - bt(1 - at)))))).$$

In this format, additional terms might actually make a difference, because  $at$  is not particularly small. However, the more computations we have, the more errors accumulate. Therefore we take the fewest that produce acceptable accuracy.

Now  $r = 1.7453292519943296$  and  $a, b$ , etc., have formulas:

$$a = r^2/13/12, \quad b = r^2/11/10, \quad c = r^2/9/8, \\ d = r^2/7/6, \quad e = r^2/5/4, \quad f = r^2/3/2.$$

An alternative method would be to accumulate a sum, computing each term from the previous one (e.g., if  $u = t^3/7!$  is the fourth term, the next one is  $u * t * (1/(8 * 9))$ ). This is a bit more complicated to code and requires moving values around more. It would have the advantage that we can stop whenever a term evaluates to zero, making computation faster for small values of  $x$ . I have not determined whether it would compromise accuracy.

We avoid divisions by precomputing the coefficients  $a, b, c$ , etc. Note that without the reduction in  $x$ , the value of  $a$  for example would be 0.00000195, with only three significant figures of accuracy. Now we can have seven, and the accuracy is more-or-less determined by that of the reduced  $x$ .

$$a = 0.01952676, \quad b = 0.02769249, \quad c = 0.04230797, \\ d = 0.07252796, \quad e = 0.15230871, \quad f = 0.50769570.$$

It is important to note that the following operations step all over the \MFP@temp $x$  \count registers, so we have made sure that we no longer need them.

The \MFP@flipz computes  $1 - z$ , where  $z$  is the result of the previous operation. Instead of simply subtracting, we optimize based on the fact that  $z$  is known to be nonnegative and not larger than 1.

The macro \MFP@com@iter ‘flipz’ the previous result then multiplies by  $t$  and the indicated coefficient. (The name of this macro stands for “common iterated” code; it is reused for some other power series.)

For extra efficiency, the power series uses a “small” version of multiplication \MFP@Rsmul, used only when the factors are sure to lie in  $[0, 1]$ . This does not take into account the sign of  $x$ , whence the ending \edef.

```

909 \def\MFP@Rsin@prog{%
910   \MFP@Rcopy tx\MFP@Rload y10{01952676}\MFP@Rsmul%
911   \MFP@com@iter{02769249}\MFP@com@iter{04230797}\MFP@com@iter{07252796}%
912   \MFP@com@iter{15230871}\MFP@com@iter{50769570}\MFP@flipz \MFP@Rcopyzx
913   \MFP@Rcopy sy\MFP@Rsmul\MFP@Rcopyzx\edef\MFP@x@Sgn{\MFP@s@Sgn}}%
914 \def\MFP@flipz{%
915   \ifnum\MFP@z@Sgn=0
916     \MFP@Rloadz 110%
917   \else
918     \MFP@tempa\MFP@ttteight
919     \advance\MFP@tempa-\MFP@z@Frc\relax
920     \MFP@Rloadz{\ifcase\MFP@tempa 0\else1\fi}0\MFP@tempa
921   \fi}%
922 \def\MFP@com@iter#1{\MFP@flipz
923   \MFP@Rcopyzx\MFP@Rcopy ty\MFP@Rsmul
924   \MFP@Rcopyzx\MFP@Rload y10{#1}\MFP@Rsmul}%

```

As to the accuracy of these computations, we can certainly lose accuracy at each step. In principle, if  $x$  is known to 10 significant figures ( $x \geq 10$  degrees), then even though we lose two figures with division by 100, the accuracy bottleneck is the fact that our coefficients have only seven figures. Now we have 17 multiplications, and while products are said to have the same number of significant figures as the factors, in the worse case we can accumulate inaccuracy of about  $.5 \times 10^{-8}$  per multiplication. So we are not guaranteed an accuracy of more than about  $\pm 10^{-7}$ . Numerical tests, however, show that it isn't that bad, probably because the direction of inaccuracies usually varies randomly, and inaccuracies in one direction compensate for those going the other way. I have not seen a case where the result is off by more than 1 in the last decimal place (i.e.,  $\pm 1.5 \times 10^{-8}$ ). In the case where we can know the result exactly,  $x = 30$ , we get an exact answer, even though we don't single it out (as we do 0, 90 and 180).

The following is the "small" version of `\MFP@Rmul`. Limited to non-negative numbers less than or equal to 1. Theoretically all the numbers are strictly between 0 and 1, but in practice a multiplication could round to 0 and then, after subtraction, a 1 could occur. We handle those easy cases separately, so that in `\MFP@@Rsmul` we don't have to worry about the integer parts at all.

Also, since these are completely internal, we don't even define the overflow and underflow macros.

```

925 \def\MFP@Rsmul{%
926   \ifnum \MFP@x@Sgn=0 \MFP@Rzero
927   \else\ifnum \MFP@y@Sgn=0 \MFP@Rzero
928   \else\ifnum\MFP@x@Int>0 \MFP@Rcopy yz%
929   \else\ifnum\MFP@y@Int>0 \MFP@Rcopy xz%
930   \else \MFP@@Rsmul
931   \fi\fi\fi\fi}%
932 \def\MFP@@Rsmul{%
933   \MFP@split\MFP@x@Frc\MFP@x@Frci\MFP@x@Frcii
934   \MFP@split\MFP@y@Frc\MFP@y@Frci\MFP@y@Frcii
935   \def\MFP@z@Frci {0}\def\MFP@z@Frcii {0}%

```

```

936 \def\MFP@z@Frc@iii{0}\def\MFP@z@Frc@iv {0}%
937 \MFP@tempb\MFP@y@Frc@ii\relax
938 \MFP@multiplyone\MFP@x@Frc@ii\MFP@z@Frc@iv
939 \MFP@multiplyone\MFP@x@Frc@i\MFP@z@Frc@iii
940 \MFP@tempb\MFP@y@Frc@i\relax
941 \MFP@multiplyone\MFP@x@Frc@ii\MFP@z@Frc@iii
942 \MFP@multiplyone\MFP@x@Frc@i\MFP@z@Frc@ii
943 \MFP@carrym\MFP@z@Frc@iv\MFP@z@Frc@iii
944 \MFP@carrym\MFP@z@Frc@iii\MFP@z@Frc@ii
945 \ifnum\MFP@z@Frc@iii<5000 \else
946   \MFP@tempb\MFP@z@Frc@ii
947   \advance\MFP@tempb1
948   \edef\MFP@z@Frc@ii{\number\MFP@tempb}\fi
949 \MFP@carrym\MFP@z@Frc@ii\MFP@z@Frc@i
950 \makeMFP@fourdigits\MFP@z@Frc@ii
951 \makeMFP@fourdigits\MFP@z@Frc@i
952 \def\MFP@z@Int{0}%
953 \edef\MFP@z@Frc{\MFP@z@Frc@i\MFP@z@Frc@ii}%
954 \edef\MFP@z@Sgn{\ifnum\MFP@z@Frc=0 0\else 1\fi}%

```

#### 4.4 Polar angle

Instead of supplying the arcsine and arccosine functions, we supply the more general angle function. This is a binary operation that accepts the two coordinates of a point and computes its angle in polar coordinates. One then has, for example,  $\arctan x = \text{angle}(1, x)$  and  $\arccos x = \text{angle}(x, \sqrt{1 - x^2})$ .

We start, as usual, with a few reductions. When the  $y$ -part is 0, we immediately return 0 or 180. If the  $y$ -part is negative, we compute the angle for  $(x, |y|)$  and negate it. If the  $x$ -part is negative, we compute the angle for  $|x|$  and subtract it from 180. Finally, reduced to both coordinates positive, if  $y > x$  we compute the angle of  $(y, x)$  and subtract that from 90. Ultimately, we apply a power series formula for  $\mathit{mathopangle}(1, y/x)$  and get convergence when the argument is less than 1, but convergence is poor unless the argument is less than 1/2. When that is not the case, conceptually, we rotate the picture clockwise by the arctangent of 1/2, compute the angle of the new point and then add a precomputed value of  $\arctan(1/2)$ .

```

955 \def\MFP@Rangle{%
956   \ifcase\MFP@y@Sgn\relax
957     \ifcase\MFP@x@Sgn\relax
958       \MFP@warn{Point (0,0) has no angle. Returning 0 anyway}%
959       \MFP@Rzero
960     \or
961       \MFP@Rzero
962     \else
963       \MFP@Rloadz 1{180}0%
964     \fi
965     \@xp\@gobble
966   \or

```



```

967     \def\MFP@angle@Sgn{1}\@xp\@firstofone
968   \else
969     \def\MFP@y@Sgn{1}%
970     \def\MFP@angle@Sgn{-1}\@xp\@firstofone
971   \fi
972   {\ifcase\MFP@x@Sgn\relax
973     \MFP@Rloadz1{90}0%
974   \or \MFP@@Rangle
975   \else
976     \def\MFP@x@Sgn{1}\MFP@@Rangle
977     \MFP@Rcopyzy\MFP@Rload x1{180}0\MFP@Rsub
978   \fi
979   \let\MFP@z@Sgn\MFP@angle@Sgn}%
980 \def\MFP@@Rangle{%
981   \MFP@Rcmp
982   \ifMFP@neg
983     \MFP@Rcopy xw\MFP@Rcopy yx\MFP@Rcopy wy%
984     \MFP@@@Rangle
985     \MFP@Rload x1{90}0\MFP@Rcopyzy\MFP@Rsub
986   \else
987     \MFP@@@Rangle
988   \fi}%

```

Precisely what we do when we are finally in the case  $0 < y < x$  is perform a couple of reductions. Ultimately we want to compute the arctan of  $z = y/x$ . We once again use a power series but, for fast convergence, we require  $z$  to be considerably less than 1. For reasons we discuss later, we won't be able to use the more efficient `\MFP@Rsmul` so we want to keep the number of iterations of our power series calculations low.

So we start with two iterations of the algorithm used by Knuth: if  $y/x > 1/2$  we transform the pair  $(x, y)$  to a new one whose angle has been reduced by  $\arctan(1/2)$ . The new pair is  $(x', y') = (2x + y, 2y - x)$ . If we still have  $y/x > 1/4$ , we perform  $(x'', y'') = (4x + y, 4y - x)$ , which then satisfies  $y''/x'' \leq 1/4$ . When either of these transformations is performed, we add the corresponding angle to the “angle-so-far” in register  $a$ .

We could continue this iteration 32 times to get (theoretically) the angle in degrees to  $\pm 10^{-8}$ . That seems a bit long, plus the accumulation of errors over 32 iterations could (in the worst case) produce less than  $\pm 10^{-7}$  accuracy.

To get the accuracy we need, we work in “scaled reals”. That is, we get 10 effective decimal places of accuracy by letting an  $x$  in the range  $0 < x < 100$  stand for  $0 < x/100 < 1$ .

Our initial reductions can increase  $x$  by a factor of about 13. Moreover, we ultimately need to scale  $y$  by 100 when we convert to scaled computations. Thus, if we make sure  $x$  is less than 1 000 000, we will prevent overflow in both cases.

```

989 \def\MFP@Rquad{\MFP@Rdbl\MFP@Rcopyzx\MFP@Rdbl}%
990 \def\MFP@@@Rangle{%
991   \MFP@Rcopy xs\MFP@Rcopy yt%
992   \ifnum\MFP@x@Int<1000000

```

```

993 \else
994   \MFP@RdivC \MFP@Rcopyz s%
995   \MFP@Rcopy tx\MFP@RdivC \MFP@Rcopyz t%
996 \fi
997 \ifnum\MFP@t@Sgn=0 \MFP@Rzero
998 \else
999   \MFP@Rcopy tx\MFP@Rdbl\MFP@Rcopyzx\MFP@Rcopy sy\MFP@Rcmp
1000   \ifMFP@pos
1001     \MFP@Rsub\MFP@Rcopyz u\MFP@Rcopy sx\MFP@Rdbl
1002     \MFP@Rcopyzx\MFP@Rcopy ty\MFP@Radd
1003     \MFP@Rcopyz s\MFP@Rcopy ut%
1004     \MFP@Rload a1{2656}{50511771}%
1005   \else
1006     \MFP@Rload a000%
1007   \fi
1008   \MFP@Rcopy tx\MFP@Rquad\MFP@Rcopyzx\MFP@Rcopy sy\MFP@Rcmp
1009   \ifMFP@pos
1010     \MFP@Rsub\MFP@Rcopyz u\MFP@Rcopy sx\MFP@Rquad
1011     \MFP@Rcopyzx\MFP@Rcopy ty\MFP@Radd
1012     \MFP@Rcopyz s\MFP@Rcopy ut%
1013     \MFP@Rcopy ax\MFP@Rload y1{1403}{62434679}%
1014     \MFP@Radd\MFP@Rcopy za%
1015   \fi
1016   \MFP@Rcopy tx\MFP@RmulC
1017   \MFP@Rcopyzx\MFP@Rcopy sy\MFP@Rdiv
1018   \MFP@Rcopyzx\MFP@Ratanc
1019   \MFP@Rcopyzx\MFP@Rdeg
1020   \MFP@Rcopyzx\MFP@Rcopy ay\MFP@Radd
1021   \MFP@Rcopyzx\MFP@RdivC
1022 \fi}%

```

Here are fast multiplication and division by 100. We need these because we are going to compute the arctangent in radians to ten decimal places. We do this by computing with scaled reals in which, for example, 0.5 is represented by 50.0. When we do this, multiplication requires a division by 100:  $.5 \times .5 = .25$  would be computed as  $(50 \times 50)/100 = 25$ .

```

1023 \def\MFP@twoofmany#1#2#3\MFP@end{#1#2}%
1024 \def\MFP@gobbletwo#1#2{}%
1025 \def\MFP@RmulC{%
1026   \edef\MFP@z@Int{\MFP@x@Int\@xp\MFP@twoofmany\MFP@x@Frc\MFP@end}%
1027   \edef\MFP@z@Frc{\@xp\MFP@gobbletwo\MFP@x@Frc00}%
1028   \edef\MFP@z@Sgn{\MFP@x@Sgn}}%
1029 \def\MFP@RdivC{%
1030   \makeMFP@eightdigits\MFP@x@Int
1031   \makeMFP@eightdigits\MFP@x@Frc
1032   \@XP\MFP@RdivC\@xp\MFP@x@Int\MFP@x@Frc\MFP@end}%
1033 \def\MFP@RdivC#1#2#3#4#5#6{%
1034   \edef\MFP@z@Int{\number#1#2#3#4#5#6}%
1035   \MFP@RdivC}%
1036 \def\MFP@RdivC#1#2#3#4#5#6#7#8#9\MFP@end{%

```

```

1037 \MFP@tempa#1#2#3#4#5#6#7#8\relax
1038 \ifnum#9>49 \advance\MFP@tempa1 \fi
1039 \edef\MFP@z@Frc{\number\MFP@tempa}%
1040 \makeMFP@eightdigits\MFP@z@Frc
1041 \edef\MFP@z@Sgn{\MFP@x@Sgn}%
1042 \ifnum\MFP@tempa=0
1043   \ifnum\MFP@z@Int=0 \def\MFP@z@Sgn{0}\fi
1044 \fi}%

```

Finally, we compute the arctan of a scaled real producing a result as a scaled number (i.e., as “centiradians”—100 times the number of radians) using a power series. Since that number could be around 0.25 (represented by 25.0), we have to sum to at least its 15th power ( $4^{-15}/15 \approx .6 \times 10^{-10}$  and the next term in the series is effectively 0). Fortunately, the power series has only odd terms, so there are only eight terms we actually need to calculate. The calculation proceeds much like the one for the sine, starting with the sum

$$x \left( 1 - \frac{u}{3} + \frac{u^2}{5} - \frac{u^3}{7} + \cdots - \frac{u^7}{15} \right),$$

where  $u = x^2$ .

We start with the common iterated code. It assumes a scaled value in x to be multiplied by the saved (scaled) value of  $x^2$  (in register u) and by a coefficient (supplied in separate integer and fractional parts). It ends with the new value in x.

```

1045 \def\MFP@scaledmul{\MFP@Rmul\MFP@Rcopyzx\MFP@RdivC}%
1046 \def\MFP@atan@iter#1#2{%
1047   \MFP@Rcopy uy\MFP@scaledmul
1048   \MFP@Rcopyzx\MFP@Rload y1{#1}{#2}\MFP@scaledmul
1049   \MFP@Rcopyzy\MFP@Rload x1{100}{00000000}%
1050   \MFP@Rsub\MFP@Rcopyzx}%
1051 \def\MFP@Ratanc{%
1052   \MFP@Rcopy xs\MFP@Rcopy xy\MFP@scaledmul
1053   \ifnum \MFP@z@Sgn=0
1054     \MFP@Rcopy sz%
1055   \else
1056     \MFP@Rcopyz u\MFP@Rcopyzx
1057     \MFP@Rload y1{86}{66666667}\MFP@scaledmul
1058     \MFP@Rcopyzy\MFP@Rload x1{100}{00000000}\MFP@Rsub\MFP@Rcopyzx
1059     \MFP@atan@iter{84}{61538462}\MFP@atan@iter{81}{81818182}%
1060     \MFP@atan@iter{77}{77777778}\MFP@atan@iter{71}{42857143}%
1061     \MFP@atan@iter{60}{00000000}\MFP@atan@iter{33}{33333333}%
1062     \MFP@Rcopy sy\MFP@scaledmul
1063 \fi}%

```

## 4.5 Logarithms

Now for logarithms. We are going to compute both common logarithms (base 10) and natural logarithms (base  $e$ ). The first step of the calculation is be essentially

trivial and works with base 10: to get the integer part of the log for numbers with positive integer part, count the digits in the integer part and subtract 1. For numbers less than one, count the number of zeros at the beginning of the fractional part and add 1 (subtract this from the result of the second part). This reduces the problem to numbers  $1 \leq x < 10$ . A few divisions (when necessary) reduce to the case where  $x = 1 + u$  with  $u$  small enough that the power series for  $\log(1 + u)$  can be computed accurately in an acceptable number of terms. Then we proceed as in the code for sine.

The power series produces a logarithm in base  $e$  so we ultimately get the answer in two parts, with the parts calculated for different bases. The last step is to multiply the second part by a conversion factor and add the first to it. For natural log, convert the first and add the second. Which one is to be returned is passed as a boolean

We keep the value-so-far in register  $s$  and the modified  $x$ -value in register  $t$ .

```

1064 \newif\ifMFP@natural
1065 \def\MFP@Rlog{\MFP@naturalfalse\MFP@Rlog@}%
1066 \def\MFP@Rln{\MFP@naturaltrue\MFP@Rlog@}%
1067 \def\MFP@Rlog@{%
1068   \ifnum\MFP@x@Sgn=0
1069     \MFP@logofzero@err
1070     \MFP@Rloadz{-1}\LogOfZeroInt\LogOfZeroFrac
1071   \else
1072     \ifnum \MFP@x@Sgn<0
1073       \MFP@warn{The logarithm of a negative number is complex.
1074       \MFP@msgbreak Only the real part will be computed}%
1075     \def\MFP@x@Sgn{1}%
1076   \fi
1077   \MFP@Rload s000%

```

If the integer part is 0, the fractional part is not. Save the number of places that will be shifted in  $\text{\MFP@tempa}$ . We use  $\text{\number}$  to strip the leading zeros and (essentially) we count the number of digits that remain. Then we shift left, putting the first digit into the integer part of  $s$  and the rest into the fractional part.

```

1078   \ifnum \MFP@x@Int=0
1079     \edef\MFP@x@Tmp{\number\MFP@x@Frc}%
1080     \MFP@tempa=\MFP@numshiftL\MFP@x@Tmp\relax
1081     \def\MFP@s@Sgn{-1}%
1082     \edef\MFP@t@Int{\@xp\MFP@oneofmany\MFP@x@Tmp\MFP@end}%
1083     \edef\MFP@t@Frc{\@xp@gobble\MFP@x@Tmp0}%
1084     \MFP@padtoeight\MFP@t@Frc
1085   \else

```

When the integer part is not 0, we get the number of digits to shift again in  $\text{\MFP@tempa}$ . It will be one less than the number of integer digits.

```

1086     \MFP@tempa\MFP@numshiftR\MFP@x@Int
1087     \edef\MFP@x@Tmp{\MFP@x@Int\MFP@x@Frc}%
1088     \ifnum\MFP@tempa>0 \def\MFP@s@Sgn{1}\fi
1089     \edef\MFP@t@Int{\@xp\MFP@oneofmany\MFP@x@Tmp\MFP@end}%

```

```

1090 \edef\MFP@x@Tmp{\@xp\@gobble\MFP@x@Tmp}%
1091 \edef\MFP@t@Frc{\@xp\MFP@eightofmany\MFP@x@Tmp\MFP@end}%
1092 \fi

```

Now the integer part of  $\log_{10} x$  is known. We save it in  $s$ . Also, set the sign of the reduced argument (positive). Then call `\MFP@Rlog@reduce`, which reduces  $x$  to less than 1.161 while possibly increasing  $s$ . For the natural log, we convert the value in  $s$ .

If the reduced  $x$  is 1, return the value in  $s$ , otherwise call the power series program (discarding the integer part of  $t$ , which should be a 1). Finally, convert the returned result if necessary and add register  $s$  to it.

```

1093 \edef\MFP@s@Int{\number\MFP@tempa}%
1094 \def\MFP@t@Sgn{1}%
1095 \MFP@Rlog@reduce
1096 \ifMFP@natural \MFP@Rcopy sx\MFP@RbaseE \MFP@Rcopy zs\fi
1097 \ifnum\MFP@t@Frc=0
1098 \MFP@Rcopy sz%
1099 \else
1100 \def\MFP@t@Int{0}\MFP@Rlog@prog
1101 \ifMFP@natural\else \MFP@Rcopyzx \MFP@RbaseX \fi
1102 \MFP@Rcopy sy\MFP@Rcopyzx\MFP@Radd
1103 \fi
1104 \fi}%

```

We determine the size of a right shift by lining up the digits in the integer part, followed by the possible numbers, and picking out the ninth argument. Similarly, to get a left shift we line up the digits of the fractional part (minus the leading zeros) followed by the possible numbers, and again picking the ninth.

```

1105 \def\MFP@numshiftR#1{\@xp\MFP@ninthofmany#176543210\MFP@end}%
1106 \def\MFP@numshiftL#1{\@xp\MFP@ninthofmany#112345678\MFP@end}%

```

In `\MFP@Rlog@reduce` we divide by the square root of 10 if the number is significantly larger than that (adding .5 to value-so-far). We repeat with the 4th, 8th and 16th roots. It seems that this could be where errors can accumulate, so the divisions are done with double precision multiplication and  $x$  is scaled by 100. Our check whether  $x > \sqrt{10}$  is rather rough: comparing the first three digits only, but even in the worst case, the final  $x$  is less than 1.1605, so at most 0.161 is fed to the power series.

```

1107 \def\MFP@Rlog@reduce{%
1108 \MFP@Rcopy tx\MFP@RmulC\MFP@Rcopyz t%
1109 \MFP@reduceonce {316}{31622776}{60168379}{50000000}%
1110 \MFP@reduceonce {177}{56234132}{51903491}{25000000}%
1111 \MFP@reduceonce {133}{74989420}{93324558}{12500000}%
1112 \MFP@reduceonce {115}{86596432}{33600654}{06250000}%
1113 \MFP@Rcopy tx\MFP@RdivC\MFP@Rcopyz t}%
1114 \def\MFP@reduceonce#1#2#3#4{%
1115 \ifnum\MFP@t@Int>#1\relax
1116 \MFP@Rcopy tx%
1117 \MFP@DPmul 0{#2}{#3}\MFP@Rcopyz t%
1118 \MFP@Rcopy sx\MFP@Rload y10{#4}\MFP@Radd

```

```

1119 \MFP@Rcopyz s%
1120 \fi}%

```

Now we have a value for  $t$  of the form  $1 + u$  with  $0 \leq u < 0.161$ . We will use the formula

$$\ln(1 + u) = \sum_{n=1}^{\infty} (-1)^{n-1} \frac{u^n}{n}.$$

We only need to carry it far enough to assure that the remaining terms would be zero in our finite resolution arithmetic, that is  $(.161)^k/k < .5 \times 10^{-8}$ . This is satisfied by  $k = 10$ . So we carry the sum to 9 places.

Again, we compute this by

$$u(1 - au(1 - bu(1 - cu(1 - du(1 - eu(1 - fu(1 - gu(1 - hu))))))))$$

where  $a = 1/2$ ,  $b = 2/3, \dots$ ,  $g = 7/8$ , and  $h = 8/9$ . This arrangement allows us to reuse \MFP@com@iter.

```

1121 \def\MFP@Rlog@prog{%
1122 \MFP@Rcopy tx\MFP@Rload y10{88888889}\MFP@Rsmul
1123 \MFP@com@iter{87500000}\MFP@com@iter{85714286}\MFP@com@iter{83333333}%
1124 \MFP@com@iter{80000000}\MFP@com@iter{75000000}\MFP@com@iter{66666667}%
1125 \MFP@com@iter{50000000}\MFP@flipz\MFP@Rcopyzx\MFP@Rcopy ty\MFP@Rsmul}%

```

## 4.6 Powers

With the exponential function we immediately return 1 if  $x = 0$ . We call two separate handlers for positive and negative  $x$ . This is because the issues are different between positive and negative exponents.

```

1126 \def\MFP@Rexp{%
1127 \ifcase\MFP@x@Sgn\relax
1128 \MFP@Rloadz 110%
1129 \or
1130 \MFP@Rexp@pos
1131 \else
1132 \def\MFP@x@Sgn{1}%
1133 \MFP@Rexp@neg
1134 \fi}%

```

One issue for positive exponents is overflow, so we issue an error message for that case. The largest number that will not produce overflow is 18.42068074 so we first compare to that; if larger, issue the error message and return 99999999.99999999.

We compute the integer power first, using an \ifcase. Because there are only 19 cases to consider a table lookup is faster than multiplications.

Then, we examine the first digit  $d$  after the decimal and compute  $e^{0.d}$ , again by cases. This is multiplied by the integer power previously found. What remains is the rest of the fractional part of  $x$ , which is strictly less than 0.1. The exponential of this is computed using the first several terms of the power series for  $e^x$ .

```

1135 \def\MFP@Rexp@pos{%

```

```

1136 \MFP@Rload y1{18}{42068074}\MFP@Rcmp
1137 \ifMFP@pos
1138 \MFP@expoverflow@err
1139 \MFP@Rloadz 1\MaxRealInt\MaxRealFrac
1140 \else
1141 \MFP@tempa\MFP@x@Int
1142 \edef\MFP@powerof@e{%
1143 1\ifcase\MFP@tempa
1144 10\or
1145 2{71828183}\or
1146 7{38905610}\or
1147 {20}{08553692}\or
1148 {54}{59815003}\or
1149 {148}{41315910}\or
1150 {403}{42879349}\or
1151 {1096}{63315843}\or
1152 {2980}{95798704}\or
1153 {8103}{08392758}\or
1154 {22026}{46579481}\or
1155 {59874}{14171520}\or
1156 {162754}{79141900}\or
1157 {442413}{39200892}\or
1158 {1202604}{28416478}\or
1159 {3269017}{37247211}\or
1160 {8886110}{52050787}\or
1161 {24154952}{75357530}\or
1162 {65659969}{13733051}\else
1163 {\MaxRealInt}{\MaxRealFrac}\fi}%
1164 \@xp\MFP@Rloadz\MFP@powerof@e
1165 \ifnum\MFP@x@Frc=0
1166 \else
1167 \MFP@Rcopyz s%
1168 \MFP@tempa=\@xp\MFP@oneofmany\MFP@x@Frc\MFP@end
1169 \edef\MFP@powerof@e{%
1170 y1\ifcase\MFP@tempa
1171 10\or
1172 1{10517092}\or
1173 1{22140276}\or
1174 1{34985881}\or
1175 1{49182470}\or
1176 1{64872127}\or
1177 1{82211880}\or
1178 2{01375271}\or
1179 2{22554093}\or
1180 2{45960311}\else
1181 10\fi}%
1182 \edef\MFP@t@Frc{0\@xp\@gobble\MFP@x@Frc}%
1183 \MFP@Rcopy sx\@xp\MFP@Rload\MFP@powerof@e\MFP@Rmul
1184 \ifnum\MFP@t@Frc=0
1185 \else

```

```

1186      \MFP@Rcopyz s\MFP@Rload t10\MFP@t@Frc
1187      \MFP@Rexp@pos@prog
1188      \MFP@Rcopy sx\MFP@Rcopyzy\MFP@Rmul
1189      \fi
1190      \fi
1191      \fi}%

```

Since the  $x$  value is now less than 0.1, we can get eight places of accuracy with only six terms of the power series. We can also arrange to use the more efficient \MFP@Rsmul for multiplication.

We organize the computation thusly

$$1 + (x + x/2(x + x/3(x + x/4(x + x/5(x + x/6))))))$$

We start by loading  $x$  (now in register  $t$ ) into register  $z$ , then repeatedly run \MFP@Rexp@iter feeding it the successive values of  $1/n$ . This iterator first multiplies the most recent result (the  $z$  register) by  $1/n$ , then that by  $x$  and then adds  $x$  to that. The final step is to add 1.

```

1192 \def\MFP@Rexp@pos@prog{%
1193   \MFP@Rcopy tz\MFP@Rexp@iter{14285714}\MFP@Rexp@iter{16666667}%
1194   \MFP@Rexp@iter{20000000}\MFP@Rexp@iter{25000000}%
1195   \MFP@Rexp@iter{33333333}\MFP@Rexp@iter{50000000}\MFP@Rcopyzx
1196   \MFP@Rincr}%
1197 \def\MFP@Rexp@iter#1{%
1198   \MFP@Rcopyzx\MFP@Rload y10{#1}\MFP@Rsmul
1199   \MFP@Rcopyzx\MFP@Rcopy ty\MFP@Rsmul
1200   \MFP@Rcopyzx\MFP@Rcopy ty\MFP@Radd}%

```

It is impossible to get accuracy to the last digit when  $e^x$  is large. This is because an absolute error in  $x$  converts to a relative error in  $e^x$ . That is, knowing  $x$  only to  $10^{-8}$  means  $e^x$  is off by (about)  $e^x \cdot 10^{-8}$ . Roughly speaking, this means only about 8 places of  $e^x$  are accurate, so if (for example) the integer part of  $e^x$  has six places then only two places after the decimal are significant.

The first issue with negative exponents is that it doesn't take much to produce a value of  $e^{-x}$  that rounds to 0. Any  $x > 19.11382792$ . So we start by comparing to that value and simply return zero if  $x$  is larger.

We perform exactly the same reductions as for positive exponents, handling the integer part and the first decimal separately. Then we call the power series program (not the same).

```

1201 \def\MFP@Rexp@neg{%
1202   \MFP@Rload y1{19}{11382792}%
1203   \MFP@Rcmp
1204   \ifMFP@pos
1205     \MFP@Rloadz 000%
1206   \else
1207     \MFP@tempa\MFP@x@Int
1208     \edef\MFP@powerof@e{%
1209       \ifcase\MFP@tempa

```



```

1210      11{0}\or
1211      10{36787944}\or
1212      10{13533528}\or
1213      10{04978707}\or
1214      10{01831564}\or
1215      10{00673795}\or
1216      10{00247875}\or
1217      10{00091188}\or
1218      10{00033546}\or
1219      10{00012341}\or
1220      10{00004540}\or
1221      10{00001670}\or
1222      10{00000614}\or
1223      10{00000226}\or
1224      10{00000083}\or
1225      10{00000031}\or
1226      10{00000011}\or
1227      10{00000004}\or
1228      10{00000002}\or
1229      10{00000001}\else
1230      000\fi}%
1231      \@xp\MFP@Rloadz\MFP@powerof@e
1232      \ifnum\MFP@x@Frc=0
1233      \else
1234      \MFP@Rcopyz s%
1235      \MFP@tempa=\@xp\MFP@oneofmany\MFP@x@Frc\MFP@end
1236      \edef\MFP@powerof@e{%
1237      y1\ifcase\MFP@tempa
1238      10\or
1239      0{90483742}\or
1240      0{81873075}\or
1241      0{74081822}\or
1242      0{67032005}\or
1243      0{60653066}\or
1244      0{54881164}\or
1245      0{49658530}\or
1246      0{44932896}\or
1247      0{40656966}\else
1248      10\fi}%
1249      \edef\MFP@t@Frc{0\@xp\@gobble\MFP@x@Frc}%
1250      \MFP@Rcopy sx\@xp\MFP@Rload\MFP@powerof@e\MFP@Rmul
1251      \ifnum\MFP@t@Frc=0
1252      \else
1253      \MFP@Rcopyz s\MFP@Rload t10\MFP@t@Frc
1254      \MFP@Rexp@neg@prog
1255      \MFP@Rcopyzx\MFP@Rcopy sy\MFP@Rmul
1256      \fi
1257      \fi
1258      \fi}%

```

Since  $x$  is now positive we calculate  $e^{-x}$ . Again we need only up to the 6th power, organized as follows

$$1 - x(1 - x/2(1 - x/3(1 - x/4(1 - x/5(1 - x/6))))))$$

Since this has exactly the same form as the the power series calculation for log and sin, we can reuse the code in `\MFP@com@iter`. We end with the final multiplication by  $x$  and the subtraction from 1 rather than call `\MFP@com@iter` with a useless multiplication by 1.

```

1259 \def\MFP@Rexp@neg@prog{%
1260   \MFP@Rcopy tx\MFP@Rload y10{14285712}\MFP@Rsmul
1261   \MFP@com@iter{16666667}\MFP@com@iter{20000000}%
1262   \MFP@com@iter{25000000}\MFP@com@iter{33333333}%
1263   \MFP@com@iter{50000000}\MFP@flipz\MFP@Rcopyzx
1264   \MFP@Rcopy ty\MFP@Rsmul\MFP@flipz}%

```

The most efficient way to take an integer power of a number  $x$  is to scan the binary code for the exponent. Each digit 1 in this code corresponds to a  $2^k$  power of  $x$ , which can be computed by repeatedly squaring  $x$ . These *dyadic* powers are multiplied together. We can convert this idea to a simple loop illustrated by this example of finding  $x^{13}$  ( $13 = 1101$  in base 2). Here  $p$  holds the current product and  $q$  holds the current dyadic power of  $x$ , initialized with  $p = 1$  and  $q = x$ :

1. Rightmost digit 1: update  $p \leftarrow pq = x$  and  $q \leftarrow q^2 = x^2$ .
2. Next digit 0: Just update  $q \leftarrow q^2 = x^4$ .
3. Next digit 1: update  $p \leftarrow pq = x^5$  and  $q \leftarrow q^2 = x^8$ .
4. Next digit 1: update  $p \leftarrow pq = x^{13}$ , detect that we are at the end and skip the update of  $q$ . Return  $p$ .

Of course, this requires the binary digits of the exponent  $n$ . But the rightmost digit of  $n$  is 1 if and only if  $n$  is odd, and we can examine each digit in turn if we divide  $n$  by 2 (discarding the remainder) at each stage. We detect the end when  $n$  is reduced to 1.

Accuracy is partly a function of the number of multiplications. The above scheme requires at most  $\lfloor \log_2 n \rfloor$  squarings and at most  $\lceil \log_2 n \rceil$  multiplications for  $x^n$ , while directly multiplying  $x \cdot x \cdots x$  would require  $n - 1$  multiplications. I have tested with an exponent equal to 8000 and it takes only about 25 times as long as a single multiplication (rather than 7999 times).

For negative powers we can either find the positive power of  $x$  and take its reciprocal or take the reciprocal of  $x$  and find its positive power. We do the second so that overflow can be detected in `\MFP@@Rpow`.

```

1265 \def\MFP@Rpow{%
1266   \ifnum\MFP@y@Frc>0
1267     \MFP@warn{The "pow" function requires an integer power.
1268     \MFP@msgbreak The fractional part will be ignored}%
1269   \fi
1270   \MFP@loopctr=\MFP@y@Int\relax
1271   \ifnum\MFP@loopctr=0
1272     \MFP@Rloadz 110%

```

```

1273 \else
1274   \ifnum\MFP@x@Sgn=0
1275     \ifnum\MFP@y@Sgn>0
1276       \MFP@Rloadz 000%
1277     \else
1278       \MFP@badpower@err
1279       \MFP@Rloadz 1\xOverZeroInt\xOverZeroFrac
1280     \fi
1281   \else
1282     \ifnum\MFP@x@Sgn>0
1283       \def\MFP@power@Sgn{1}%
1284     \else
1285       \edef\MFP@power@Sgn{\ifodd\MFP@loopctr -\fi 1}%
1286     \fi
1287     \ifnum\MFP@y@Sgn<0 \MFP@Rinv \MFP@Rcopyzx\fi
1288     \ifnum\MFP@loopctr=1
1289       \MFP@Rloadz \MFP@power@Sgn\MFP@x@Int\MFP@x@Frc
1290     \else
1291       \MFP@@@Rpow
1292     \fi
1293   \fi
1294 \fi}%

```

This implements the algorithm discussed above. We save  $x$  in register  $q$ , initialize the starting value of 1 in  $p$  and then run the loop. If the binary digit just read is a 1 (i.e., `\ifodd` is true), it multiplies  $p$  and  $q$ . It also saves the last product (copies  $z$  to  $p$ ). This need not be done on the last iteration, but must not be moved out of the `\ifodd` conditional because intervening computations modify  $z$ . If there are more iterations to do (i.e., the `\ifnum` is true), this squares  $q$  and reduces the counter. Note that the exponents 0 and 1 do not occur since we have handled them separately.

In case of overflow (either the multiplication or the squaring) we break the loop and return  $\pm\infty$ .

```

1295 \def\MFP@@@Rpow{%
1296   \MFP@Rcopy xq%
1297   \MFP@Rload p110%
1298   \MFP@Rpow@loop}%
1299 \def\MFP@Rpow@loop{%
1300   \ifodd\MFP@loopctr
1301     \MFP@Rcopy px\MFP@Rcopy qy\MFP@Rmul
1302     \ifnum \MFP@z@Ovr>0 \MFP@handle@expoverflow
1303   \else
1304     \ifnum\MFP@loopctr>1 \MFP@Rcopyz p\fi
1305   \fi
1306 \fi
1307 \ifnum\MFP@loopctr>1
1308   \MFP@Rcopy qx\MFP@Rsq
1309   \ifnum \MFP@z@Ovr>0 \MFP@handle@expoverflow
1310 \else

```

```

1311      \MFP@Rcopyz q%
1312      \divide\MFP@loopctr 2
1313      \@XP\MFP@Rpow@loop
1314      \fi
1315      \fi}%
1316 \def\MFP@handle@expoverflow{%
1317   \MFP@expoverflow@err
1318   \MFP@loopctr=0
1319   \MFP@Rloadz\MFP@power@Sgn\MaxRealInt\MaxRealFrac}%

```

## 4.7 The square root

One can combine logarithms and exponentials to get any power: to get  $x^y$ , compute  $e^{y \ln x}$ . This has the disadvantage that it doesn't work if  $x$  is negative. Most powers of negative numbers are not defined, but certainly integer powers are. Thus we have defined `\MFPpow` and `\Rpow` for that case.

If we enforce a positive  $x$ , then  $y$  can have any value. However, the computation of  $e^{.5 \ln x}$  cannot give a result as good as one can get from a special purpose algorithm for the square root. For example, the inaccuracies in computing  $\ln x$  will make  $e^{.5 \ln 9}$  inexact, while the square root function we implement below will produce exactly  $\sqrt{9} = 3$ . In fact, if a square root can be expressed exactly within our 8-digit precision, our code will find it.

For the square root we return zero if  $x$  is not positive. If the integer part of  $x$  is 0, we copy the fractional part to the integer part (that is, we multiply by  $10^8$ , remembering to multiply by  $10^{-4}$  later). This makes the square root of such numbers rather more accurate. (To get around some other rare but annoying inaccuracies, we go through a similar process when the integer part of  $x$  is at most 4 digits, multiplying by  $10^4$  before and by  $10^{-2}$  after.)

We then compute the square root using an algorithm that will be exact whenever possible. We perform one additional processing step. To explain it, note that our algorithm actually produces the largest number  $s$  with four digits right of the decimal place that satisfies  $s^2 \leq x$ . That is

$$s^2 \leq x < (s + 10^{-4})^2$$

From this it follows that  $x = (s + \epsilon)^2 = s^2 + 2s\epsilon + \epsilon^2$  with  $\epsilon < 10^{-4}$  (and so  $\epsilon^2 < 10^{-8}$ ). We estimate this  $\epsilon$  and add that estimate to  $s$ . The estimate we use is obtained by discarding the very small  $\epsilon^2$  and solving for the remaining  $\epsilon$  get

$$\epsilon \approx \bar{\epsilon} = \frac{x - s^2}{2s}$$

With this value,  $s + \bar{\epsilon}$  misses the exact square root by at most  $\epsilon^2/(2s) < .5 \cdot 10^{-8}$ , because  $s \geq 1$ . The final result  $s + \bar{\epsilon}$  is equivalent to computing the average  $s$  and  $x/s$ . This, possibly divided by  $10^4$  or  $10^2$  is the returned value.

By tests, with rare exceptions, our computations produces a result correct in all eight decimal places. In the rare exception, the last place is within 1 of the correct value.

```

1320 \def\MFP@Rsqr{t{%
1321   \ifcase\MFP@x@Sgn\relax
1322     \MFP@Rzero
1323   \or
1324     \ifnum\MFP@x@Int=0
1325       \def\MFP@sqr{t@reduce{2}%
1326       \edef\MFP@x@Int{\number\MFP@x@Frc}%
1327       \edef\MFP@x@Frc{00000000}%
1328     \else\ifnum\MFP@x@Int<10000
1329       \def\MFP@sqr{t@reduce{1}%
1330       \edef\MFP@x@Int{\MFP@x@Int\@xp\MFP@four{ofmany}\MFP@x@Frc\MFP@end}%
1331       \edef\MFP@x@Frc{\@xp\MFP@gobblefour\MFP@x@Frc0000}%
1332     \else
1333       \def\MFP@sqr{t@reduce{0}%
1334     \fi\fi
1335     \MFP@Rcopy xt%
1336     \MFP@Isqr{t
1337     \MFP@Rcopyz s\MFP@Rcopyzy
1338     \MFP@Rcopy tx\MFP@Rdiv
1339     \MFP@Rcopy sx\MFP@Rcopyzy\MFP@Radd
1340     \MFP@Rcopyzx\MFP@Rhalve
1341     \ifcase \MFP@sqr{t@reduce\relax
1342     \or
1343       \MFP@Rcopyzx\MFP@Rload y10{01000000}\MFP@Rmul
1344     \or
1345       \MFP@Rcopyzx\MFP@Rload y10{00010000}\MFP@Rmul
1346     \fi
1347   \else
1348     \MFP@warn{Square root of a negative number. Zero will be returned.}%
1349     \MFP@Rzero
1350   \fi}%
1351 \def\MFP@four{ofmany#1#2#3#4#5\MFP@end{#1#2#3#4}%
1352 \def\MFP@gobblefour{#1#2#3#4}%

```

There is a rather straightforward pencil and paper algorithm that provides the square root digit by digit, and it produces an exact answer when that is possible. Unfortunately, the decimal version is not easy to code. Fortunately the same algorithm works in any number base and it is rather simple to code the binary version (because we only need to decide at each stage whether the “next digit” is 0 or 1. This produces a square root in binary digits, from which it is easy to compute the number itself. The result is exact if the answer would be a finite number of binary digits. We apply it to the integer  $10^8x$ . While this number is too large for  $\text{\TeX}$  to handle as an integer, it is not that hard to convert it to a string of binary digits stored in a macro.

The process turns out to be simpler if we convert  $10^8x$  to base 4 rather than binary. Also, instead of producing the square root encoded in a string of binary digits, we simply build the numerical result as we discover the binary digits (multiply previous value by two and add the new digit.) Fortunately, the square root of  $10^8x$  (and the temporary scratch registers used in the code) will never exceed

$\TeX$ 's limit for integers.

The macro `\MFP@ItoQ` implements the conversion to base-4 digits. The two arguments are the integer and fractional part of  $x$ . The result is stored in `\MFP@ItoQ@Tmp`, which is so far only used by the square root code.

The test `\ifodd\MFP@tempb` is used to get the binary digits. Combining two of them yields the quadrenary digits. The `\ifodd\MFP@tempa` tests are there to check whether there will be a remainder after division by 2, which should then be inserted at the front of `\MFP@tempb` before division by 2. Two divisions by 2 each iteration amounts to division by 4. This is slightly more efficient than dividing by 4 and determining the remainder.

```

1353 \def\MFP@ItoQ#1#2{%
1354   \MFP@tempa#1\relax\MFP@tempb#2\relax
1355   \def\MFP@ItoQ@Tmp{}\MFP@ItoQ@loop}%
1356 \def\MFP@ItoQ@loop{%
1357   \ifodd\MFP@tempb
1358     \ifodd\MFP@tempa \advance\MFP@tempb \MFP@ttteight\relax\fi
1359     \divide\MFP@tempa2 \divide\MFP@tempb2
1360     \edef\MFP@ItoQ@Tmp{\ifodd\MFP@tempb 3\else 1\fi\MFP@ItoQ@Tmp}%
1361   \else
1362     \ifodd\MFP@tempa \advance\MFP@tempb \MFP@ttteight\relax\fi
1363     \divide\MFP@tempa2 \divide\MFP@tempb2
1364     \edef\MFP@ItoQ@Tmp{\ifodd\MFP@tempb 2\else 0\fi\MFP@ItoQ@Tmp}%
1365   \fi
1366   \ifodd\MFP@tempa \advance\MFP@tempb \MFP@ttteight\relax\fi
1367   \divide\MFP@tempa 2 \divide\MFP@tempb 2
1368   \ifnum\MFP@tempa>0
1369     \xp\MFP@ItoQ@loop
1370   \else\ifnum\MFP@tempb>0
1371     \XP\MFP@ItoQ@loop
1372   \fi\fi}%

```

This integer square root  $n$  is  $10^4$  times the largest number  $y$  satisfying  $y^2 \leq x$  and having at most four decimal places. The rest of the code after the `\MFP@Isqrt@loop` is intended to divide  $n$  (returned in `\MFP@tempc`) by  $10^4$  in order to get the number  $y$  itself.

```

1373 \def\MFP@Isqrt{%
1374   \MFP@ItoQ\MFP@x@Int\MFP@x@Frc
1375   \MFP@tempa=0 \MFP@tempb=0 \MFP@tempc=0
1376   \expandafter\MFP@Isqrt@loop\MFP@ItoQ@Tmp\MFP@end
1377   \MFP@tempa=\MFP@tempc
1378   \divide\MFP@tempc\MFP@tttfour
1379   \edef\MFP@z@Int{\number\MFP@tempc}%
1380   \multiply\MFP@tempc \MFP@tttfour
1381   \advance\MFP@tempa -\MFP@tempc
1382   \edef\MFP@z@Frc{\number\MFP@tempa}%
1383   \makeMFP@fourdigits\MFP@z@Frc
1384   \edef\MFP@z@Frc{\MFP@z@Frc0000}%
1385   \def\MFP@z@Sgn{1}}%

```

The following is a loop that essentially performs a base-2 version of the base-10

algorithm that I learned at age 12 from my father (apparently it was taught in eighth or ninth grade in his day). Seeing it written out, I am surprise at how concise and elegant it is!

```

1386 \def\MFP@Isqrt@loop#1{%
1387   \ifx\MFP@end #1%
1388   \else
1389     \multiply\MFP@tempa 2 \multiply\MFP@tempb 4 \multiply\MFP@tempc 2
1390     \advance \MFP@tempb#1\relax
1391     \ifnum\MFP@tempa<\MFP@tempb
1392       \advance\MFP@tempc 1 \advance\MFP@tempa 1
1393       \advance\MFP@tempb -\MFP@tempa
1394       \advance\MFP@tempa 1
1395     \fi
1396     \expandafter\MFP@Isqrt@loop
1397   \fi}%
1398 \MFP@xfinish
1399 \</extra>

```

For my own benefit: the above code finds the next binary digit and updates the square root (in `\MFP@tempc`) by appending that digit. The new digit is also appended to the end of `\MFP@tempa`. This is subtracted from `\MFP@tempb`, but only if the last digit is a 1. Then the next quadrenary digit is appended to `\MFP@tempb`. Finally, the last binary digit found is added (not appended) to `\MFP@tempa`. The “appending” of a digit means a multiplication by 2 (or 4) and the addition of the digit. We perform such additions only if the digit is a 1, and we determine if the digit is 1 or 0 by the `\ifnum` test.

## Index

Numbers refer to the page(s) where the corresponding entry is described.

<b>E</b>	<b>L</b>	<code>\MFPdeg</code> . . . . . 39, 41
<code>\EndofStack</code> . . . . . 12, 16	<code>\LogOfZeroFrac</code> . . . . . 41	<code>\MFPdiv</code> . . . . . 5, 22
<code>\Export</code> . . . . . 11, 22	<code>\LogOfZeroInt</code> . . . . . 41	<code>\MFPe</code> . . . . . 7, 39
<code>\ExportStack</code> . . . . . 11, 22		<code>\MFPexp</code> . . . . . 39, 41
<b>G</b>	<b>M</b>	<code>\MFPfloor</code> . . . . . 6, 22
<code>\Global</code> . . . . . 11, 22	<code>\MaxRealFrac</code> . . . . . 16, 41	<code>\MFPfrac</code> . . . . . 6, 22
<code>\GlobalStack</code> . . . . . 11, 22	<code>\MaxRealInt</code> . . . . . 16, 41	<code>\MFPhalve</code> . . . . . 6, 22
	<code>\MFPabs</code> . . . . . 6, 22	<code>\MFPincr</code> . . . . . 6, 22
	<code>\MFPadd</code> . . . . . 5, 22	<code>\MFPint</code> . . . . . 6, 22
	<code>\MFPangle</code> . . . . . 39	<code>\MFPinv</code> . . . . . 6, 22
<b>I</b>	<code>\MFPceil</code> . . . . . 6, 22	<code>\MFPln</code> . . . . . 39, 41
<code>\IFeq</code> . . . . . 7, 19	<code>\MFPchk</code> . . . . . 7, 19	<code>\MFPlog</code> . . . . . 39, 41
<code>\IFgt</code> . . . . . 7, 19	<code>\MFPchs</code> . . . . . 6, 22	<code>\MFPmax</code> . . . . . 5, 22
<code>\IFlt</code> . . . . . 7, 19	<code>\MFPcmp</code> . . . . . 7, 19	<code>\MFPmin</code> . . . . . 5, 22
<code>\IFneg</code> . . . . . 7, 19	<code>\MFPcos</code> . . . . . 39, 41	<code>\MFPmpy</code> . . . . . 5, 22
<code>\IFpos</code> . . . . . 7, 19	<code>\MFPdbl</code> . . . . . 6, 22	<code>\MFPmul</code> . . . . . 5, 22
<code>\IFzero</code> . . . . . 7, 19	<code>\MFPdecr</code> . . . . . 6, 22	<code>\MFPnoop</code> . . . . . 6, 23

