

# Floating-point numbers

Version 0.6.0

21/06/2014

**Laurent Bartholdi**

MPFR- and CXSC-based library for GAP

**Laurent Bartholdi** Email: [laurent.bartholdi@gmail.com](mailto:laurent.bartholdi@gmail.com)

Homepage: <http://www.uni-math.gwdg.de/laurent/>

**Address:** Mathematisches Institut  
Bunsenstraße 3-5  
D-37073 Göttingen  
Germany

## **Abstract**

This document describes the package `Float`, which implements in GAP arbitrary-precision floating-point numbers.

For comments or questions on `Float` please contact the author.

## **Copyright**

© 2011-2014 by Laurent Bartholdi

## **Acknowledgements**

Part of this work is supported by the "Swiss National Fund for Scientific Research (SNF)", the "German National Science Foundation (DFG)", and the Courant Research Centre "Higher Order Structures" of the University of Göttingen.

# Contents

<b>1</b>	<b>Licensing</b>	<b>4</b>
<b>2</b>	<b>Float package</b>	<b>5</b>
2.1	A sample run . . . . .	5
<b>3</b>	<b>Polynomials</b>	<b>7</b>
3.1	The Floats pseudo-field . . . . .	7
3.2	Roots of polynomials . . . . .	7
3.3	LLL lattice reduction . . . . .	7
<b>4</b>	<b>Implemented packages</b>	<b>8</b>
4.1	MPFR . . . . .	8
4.2	MPFI . . . . .	8
4.3	MPC . . . . .	8
4.4	CXSC . . . . .	9
4.5	FPLLL . . . . .	9
	<b>Index</b>	<b>10</b>

# Chapter 1

## Licensing

This program is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2 of the License, or any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program, in the file COPYING. If not, see <http://www.gnu.org/licenses/>.

## Chapter 2

# Float package

### 2.1 A sample run

The extended floating-point capabilities of GAP are installed by loading the package via `LoadPackage("float");` and selecting new floating-point handlers via `SetFloats(MPFR)`, `SetFloats(MPFI)`, `SetFloats(MPC)` or `SetFloats(CXSC)`, depending on whether high-precision real, interval or complex arithmetic are desired, or whether a fast package containing all four real/complex element/interval arithmetic is desired:

Example

```
gap> LoadPackage("float");
Loading FLOAT 0.3 ...
true
gap> SetFloats(MPFR); # floating-point
gap> x := 4*Atan(1.0);
.314159e1
gap> Sin(x);
.169569e-30
gap> SetFloats(MPFR,1000); # 1000 bits
gap> x := 4*Atan(1.0);
.314159e1
gap> Sin(x);
.125154e-300
gap> String(x,300);
".3141592653589793238462643383279502884197169399375105820974944592307816406286\
208998628034825342117067982148086513282306647093844609550582231725359408128481\
117450284102701938521105559644622948954930381964428810975665933446128475648233\
78678316527120190914564856692346034861045432664821339360726024914127e1"
gap>
gap> SetFloats(MPFI); # intervals
gap> x := 4*Atan(1.0);
.314159e1(99)
gap> AbsoluteDiameter(x); Sup(x); Inf(x);
.100441e-29
.314159e1
.314159e1
gap> Sin(x);
-.140815e-29(97)
gap> 0.0 in last;
```

```
true
gap> 1.0; # exact representation
.1e1(inf)
gap> IncreaseInterval(last,0.001); # now only 8 significant bits
.1e1(8)
gap> IncreaseInterval(last,-0.002); # now becomes empty
\emptyset
gap> MinimalPolynomial(Rationals,Sqrt(2.0));
-2*x_1^2+1
gap> Cyc(last);
E(8)-E(8)^3
gap>
gap> SetFloats(MPC); # complex numbers
```

## Chapter 3

# Polynomials

### 3.1 The Floats pseudo-field

Polynomials with floating-point coefficients may be manipulated in GAP; though they behave, in subtle ways, quite differently than polynomials over rings.

The "pseudo-field" of floating-point numbers is an object in GAP, called `FLOAT_PSEUDOFIELD`. (It is not really a field, e.g. because addition of floating-point numbers is not associative). It may be used to create indeterminates, for example as

Example

```
gap> x := Indeterminate(FLOAT_PSEUDOFIELD, "x");
x
gap> 2*x^2+3;
2.0*x^2+3.0
gap> Value(last, 10);
203.0
```

### 3.2 Roots of polynomials

The Jenkins-Traub algorithm has been implemented, in arbitrary precision for MPFR and MPC.

Furthermore, CXSC can provide complex enclosures for the roots of a complex polynomial.

### 3.3 LLL lattice reduction

A faster implementation of the LLL lattice reduction algorithm has also been implemented. It is accessible via the commands `FPLLLReducedBasis(m)` and `FPLLLShortestVector(m)`.

## Chapter 4

# Implemented packages

### 4.1 MPFR

#### 4.1.1 IsMPFRFloat

- ▷ IsMPFRFloat (filter)
- ▷ TYPE\_MPFR (global variable)

The category of floating-point numbers.

Note that they are treated as commutative and scalar, but are not necessarily associative.

### 4.2 MPFI

#### 4.2.1 IsMPFIFloat

- ▷ IsMPFIFloat (filter)
- ▷ TYPE\_MPFI (global variable)

The category of intervals of floating-point numbers.

Note that they are treated as commutative and scalar, but are not necessarily associative.

### 4.3 MPC

#### 4.3.1 IsMPCFloat

- ▷ IsMPCFloat (filter)
- ▷ TYPE\_MPC (global variable)

The category of intervals of floating-point numbers.

Note that they are treated as commutative and scalar, but are not necessarily associative.



## 4.4 CXSC

### 4.4.1 IsCXSCReal

▷ IsCXSCReal	(filter)
▷ IsCXSCComplex	(filter)
▷ IsCXSCInterval	(filter)
▷ IsCXSCBox	(filter)
▷ TYPE_CXSC_RP	(global variable)
▷ TYPE_CXSC_CP	(global variable)
▷ TYPE_CXSC_RI	(global variable)
▷ TYPE_CXSC_CI	(global variable)

The category of floating-point numbers.

Note that they are treated as commutative and scalar, but are not necessarily associative.

## 4.5 FPLLL

### 4.5.1 FPLLLReducedBasis

▷ FPLLLReducedBasis( $m$ )	(operation)
----------------------------	-------------

**Returns:** A matrix spanning the same lattice as  $m$ .

This function implements the LLL (Lenstra-Lenstra-Lovász) lattice reduction algorithm via the external library `fpLLL`.

The result is guaranteed to be optimal up to 1%.

### 4.5.2 FPLLLShortestVector

▷ FPLLLShortestVector( $m$ )	(operation)
------------------------------	-------------

**Returns:** A short vector in the lattice spanned by  $m$ .

This function implements the LLL (Lenstra-Lenstra-Lovász) lattice reduction algorithm via the external library `fpLLL`, and then computes a short vector in this lattice.

The result is guaranteed to be optimal up to 1%.

# Index

FPLLLReducedBasis, 9  
FPLLLShortestVector, 9

IsCXSCBox, 9  
IsCXSCComplex, 9  
IsCXSCInterval, 9  
IsCXSCReal, 9  
IsMPCFloat, 8  
IsMPFIFloat, 8  
IsMPFRFloat, 8

TYPE\_CXSC\_CI, 9  
TYPE\_CXSC\_CP, 9  
TYPE\_CXSC\_RI, 9  
TYPE\_CXSC\_RP, 9  
TYPE\_MPC, 8  
TYPE\_MPFI, 8  
TYPE\_MPFR, 8