

Stochastic Arithmetic in Multiprecision

The SAM library

LIP6, Sorbonne Université, CNRS
Paris, France



Contents

1	Introduction	5
1.1	Aim of the SAM library	5
1.2	The DSA (Discrete Stochastic Arithmetic)	7
1.2.1	The CESTAC method	7
1.2.2	The computational zero	9
1.2.3	Discrete stochastic relations	10
2	Reference guide	11
2.1	Stochastic type	11
2.2	Assignment	11
2.3	Intrinsic functions	11
2.3.1	Conversion functions	11
2.3.2	Numerical functions	12
2.3.3	Mathematical functions	12
2.4	Relational operators	12
2.5	SAM specific functions	13
2.5.1	Initializing and closing the library	13
2.5.2	Obtaining a string from a result with its evaluated accuracy	14
2.5.3	Obtaining the number of exact significant digits of a stochastic variable	15
2.5.4	Obtaining the triplet associated with a stochastic variable	16
2.5.5	Testing if a variable is a computational zero	16
2.5.6	Reducing accuracy of initial data	16
3	User's guide	19
3.1	Declaration of the SAM library	19
3.2	Initialization of the SAM library	20

3.3	Declaration of variables	20
3.3.1	Changes in the type of variables	20
3.4	Changes in assignments or arithmetic operations	20
3.4.1	Conversions between usual types and the stochastic type	20
3.4.2	Classical arithmetic operators	21
3.5	Changes in reading statements	21
3.6	Changes in printing statements	22
3.7	Constants passed as function arguments	22
3.8	Termination of the SAM library	23
3.9	An example of numerical code and its modified version	23
3.9.1	Standard C source code	23
3.9.2	Source code using the SAM library	24
3.9.3	Example of execution without SAM	25
3.9.4	Example of execution with SAM	26
3.10	Numerical debugging with SAM	26
4	Test runs	29
4.1	Example 1: a rational fraction function of two variables . . .	29
4.2	Example 2: solving a second order equation	30
4.3	Example 3: computing a determinant	31
4.4	Example 4: computing a second order recurrent sequence . .	33
4.5	Example 5: computing a root of a polynomial	37
4.6	Example 6: solving a linear system	39
4.7	Example 7: when SAM fails	40

Chapter 1

Introduction

1.1 Aim of the SAM library

The arithmetic commonly used on computers for scientific programming is floating point arithmetic. This arithmetic only approximates exact arithmetic. Consequently each arithmetic statement generates a round-off error. So when a correct program with regard to syntax and logical organization is running on a computer, every produced result is unavoidably given with a so called “computing error”. This error is due to all the round-off errors produced along the elementary statements required to obtain the result. Sometimes the error may be such that the final result is really wrong (and not only inaccurate).

The aim of the SAM library presented here is to answer the following question:

What is the computing error due to floating point arithmetic on the results produced by any program running on a computer?

So, we want to estimate the round-off error on each result with a technique which is independent on the program and hence on the algorithm used.

SAM is a library, based on the MPFR library. More precisely, SAM is a set of data types, functions and subroutines that may be used in any program written in C/C++. It implements the CESTAC method in a synchronous way (the Discrete Stochastic Arithmetic DSA). With a few modifications in the source code, this library has for main purpose to estimate the effects of round-off error propagation on every numerical computed result. It also allows to study the effects of the initial data uncertainties upon computed results, as described in 2.5.

This implementation consists in replacing the computer deterministic arithmetic by a stochastic arithmetic (the Discrete Stochastic Arithmetic DSA) and in performing N times ($N = 3$) each elementary operation before executing the next statement.

Thus, it is as N identical programs were simultaneously running on N synchronized computers each of them using random arithmetic. So for each result, we obtain N samples from which we compute the mean value and the standard deviation which characterize the corresponding stochastic number. The value of this number is defined as the mean value of the different samples. The accuracy of this number, *i.e.* its number of exact significant digits, is estimated using the mean value and the standard deviation. If all the samples are equal to zero or if the number of exact significant (decimal) digits is less than one, then the number is defined as a computational zero. This means that a computational zero is either the mathematical zero or a number without any significance.

So round-off error propagation can be analyzed step by step. Numerical instabilities and non significant results are detected. The branchings based on order relations may also be controlled. Therefore, this synchronous implementation of the CESTAC method allows to validate any scientific code during its run.

With the SAM library, one can run any scientific code using random arithmetic, without having to rewrite or notably change the initial code. This tool has been written in C++. This language enables to create new numerical **types** with their operators; furthermore the designating symbol of an operator can be chosen among the primitive symbols in the language (+, *, ...). In other words, this language enables the so called “operator overloading”. Thanks to these new properties, SAM has been developed for C/C++ programs.

Thus a new numerical type has been created, the **stochastic number**; it is nothing else than an N -set ($N = 3$) containing perturbed floating-point values (of type *mpfr_t*). All the arithmetic operators (+, −, *, /) have been overloaded in such a manner that when an operator is used, the operands are N -sets and the returned result is a randomly perturbed N -set. The relational operators (>, ≥, <, ≤, ==, ≠) are overloaded. All standard functions defined in “math.h” (SIN, COS, EXP, ...) have also been overloaded. Likewise, in/out statements have been modified, mainly the printing statement which gives as a result the mean value of the N -set written with only its exact significant (decimal) digits.

Furthermore, in order to enable the evaluation of the weight of uncertainties

on initial data on the results, a function called `data_st` may be used to perturb data as exposed in 2.5.6.

During the run of a program, as soon as a numerical anomaly (for example the product of non-significant numbers, or a relational test involving a non-significant result) is produced, some special counters are updated. At the end of the run, all information about numerical anomalies is printed on the standard output.

If no anomaly has been detected, it means that the program runs without any numerical problem. Results are then given with their accuracy - number of exact significant (decimal) digits.

If some numerical anomalies have been detected, they must be analysed. Helped by the debugger associated with the compiler, the user may retrieve the statements that produced the anomalies and determine if changes in the code are required.

For every source code, the run time and also the memory required are only multiplied by about three, what is quite reasonable according to the major interest of validating numerical results. SAM is the only existing tool that is able of a such performance.

Integrated with MPFR, SAM is able to set the decimal accuracy (precision in bits) exactly to any valid value for each variable (including very small precision).

The stochastic types and the overloaded or newly defined functions of the library are presented in the next sections.

1.2 The DSA (Discrete Stochastic Arithmetic)

1.2.1 The CESTAC method

The CESTAC (Contrôle et Estimation Stochastique des Arrondis de Calcul) method, which has been developed by La Porte and Vignes [29, 18, 21], enables one to estimate the number of exact significant digits of any computed result.

The basic idea of the method is defined in [34, 35] and consists of the following:

- to perform the same code N times with a different round-off error propagation for each run,

- to estimate the common part of these results and to consider that this part is representative of exact result.

In practice, these different round-off error propagations are obtained by using the random rounding mode defined below.

Each result ρ of a floating-point operation (assignment, arithmetical operation) which is not an exact floating-point value, is bounded by two floating-point values, one by default ρ^- and the other by excess ρ^+ .

The random rounding mode consists, at the level of each floating-point operation or assignment, in choosing as a result randomly with an equal probability either ρ^- or ρ^+ .

With this random rounding mode, the same program run several times provides different results, due to different round-off errors.

Let us consider a sequence of computations providing an exact result r . When this sequence is performed with the CESTAC [25, 24, 18] method, N results R_k , $k = 1, \dots, N$ are obtained. From the formalization of the round-off errors of the floating-point arithmetic operations $(+, -, *, /)$ a probabilistic model for estimating the round-off error on the mean value \bar{R} of the R_k , considered as the computed result, has been established. This model is a first order model. It means that the terms in 2^{-2p} (p being the number of bits of the mantissa) which appear in the expression of the round-off error of the floating-point multiplications and divisions have been neglected. Only the terms in 2^{-p} are considered.

This model is based on two hypotheses.

- Hyp1.: The elementary round-off errors α_i of the floating-point arithmetic operations are random independent, centered and uniformly distributed variables.
- Hyp2.: The approximation of the first order in 2^{-p} is legitimate.

It has been proved that if the two hypotheses hold then the R_k , $k = 1, \dots, N$ are samples of the Gaussian distribution, centered on the exact result r . Thus it is possible to use the Student's test which allows to obtain a confident interval of \bar{R} with a $(1 - \beta)$ probability and then to estimate the number of exact significant digits of \bar{R} by the formula

$$C_{\bar{R}} = \log_{10} \left(\frac{\sqrt{N} |\bar{R}|}{\tau_{\beta} \sigma} \right) \quad (1.1)$$

with

$$\overline{R} = \frac{1}{N} \sum_{i=1}^N R_i$$

and

$$\sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (R_i - \overline{R})^2.$$

τ_β is the value of the Student's distribution for $N - 1$ degrees of freedom and a probability level $1 - \beta$. In practice $N = 3$, $\beta = 0.05$ and then $\tau_\beta = 4.303$.

The result provided by eq(1.1) is reliable when the two previous hypotheses hold in practice [24, 36].

- Concerning Hyp1, with the use of random rounding, the α_i are truly independent random variables. However they are not exactly centered, consequently the \overline{R} is biased. But because of the robustness of Student's test, Hyp1 still holds. This hypothesis is not an inconvenience for the reliability of eq(1.1).
- Concerning Hyp2, it holds if the terms in 2^{-2p} are negligible in comparison to the terms in 2^{-p} . It has been proved that this fact is satisfied if
 - the operands of any multiplication are both significant
 - the divisor of any division is significant.

It is then absolutely necessary to control these two points during a run of code.

Indeed if they are not satisfied, this means that the Hyp2 has been violated and then the results obtained with eq(1.1) must be considered as not reliable. This control is done with the concept of the informatical zero also named computational zero or computed zero.

1.2.2 The computational zero

Each result provided by the CESTAC method is an informatical zero also called "computational zero" denoted by @.0 if one of the two following conditions holds:

- $\forall i, i = 1, \dots, N, R_i = 0$
- $C_{\overline{R}} \leq 0$ ($C_{\overline{R}}$ obtained with eq(1.1))

When $C_{\overline{R}} \leq 0$, then \overline{R} is an insignificant value.

From the concept of computational zero, discrete stochastic relations have been defined (equality and order relations).

1.2.3 Discrete stochastic relations

Let X and Y be N -samples provided by CESTAC method.

- Discrete stochastic equality denoted by $s =$ is defined as:
 $Xs = Y$ if $X - Y = @.0$
- Discrete stochastic inequalities denoted by $s >$ and $s \geq$ are defined as:
 $Xs > Y$ if $\overline{X} > \overline{Y}$ and $X - Y \neq @.0$
 $Xs \geq Y$ if $\overline{X} \geq \overline{Y}$ or $X - Y = @.0$

The Discrete Stochastic Arithmetic (DSA) [9, 36, 16] is defined from the CESTAC method, the concept of informatical zero and the discrete stochastic relations. With this DSA, it is possible to control the run of a scientific code, to detect the numerical instabilities and the violation of the hypotheses underlying the method.

Chapter 2

Reference guide

2.1 Stochastic type

SAM provides one new numerical type, the *stochastic type*:

`mp_st` for stochastic variables in multiple-precision
stochastic type associated to `mpfr_t`

2.2 Assignment

The operator “=” is overloaded and accepts stochastic types. It sets a stochastic variable with different types of values such as `float`, `double`, `long`, `int`, `unsigned` and MPFR object.

Also we have the function `sam_set_str(stochastic argument, string)` which sets the stochastic argument to the value of the string in base 10, rounded in the direction `MPFR_RNDN`, and returns the value of the stochastic argument.

2.3 Intrinsic functions

We present here how the intrinsic functions defined in C have been extended for stochastic types.

2.3.1 Conversion functions

The `float`, `double`, `long`, `unsigned` and `int` cast operators:

They act on variables of stochastic type and work like for numerical predefined types. Thus the result is of classical type and the knowledge of the

accuracy is lost.

If X is a stochastic variable consisting in N samples X_i , for instance

- (int) X is computed as $(\text{int})(\frac{\sum_{i=1}^N X_i}{N})$.

2.3.2 Numerical functions

The **fabs** function:

Given a `mp_st` argument, this function returns a positive `mp_st` value.

The **floor**, **ceil** and **rint** functions:

These functions accept stochastic arguments and work like on the classical types.

The **pow** function:

This function accepts both classical or stochastic types. If at least one argument is a stochastic variable, the output value is of stochastic type.

2.3.3 Mathematical functions

These are the following functions: `sqrt`, `exp`, `log`, `log10`, `sin`, `cos`, `tan`, `asin`, `acos`, `atan`, `atan2`, `sinh`, `cosh`, `tanh`, `hypot`. They accept arguments of `mp_st` stochastic type. The output value has the same type as the argument. If the function has two arguments, they must be of the stochastic type.

2.4 Relational operators

Comparison operators are overloaded and accept stochastic types and a mixture of classical types and stochastic types. They take into account the accuracy of the operands.

Thus when the expression `a == 0.0` is true, it means that `a` is a *computational zero*, i.e.

- `a` is a mathematical zero or
- `a` has no exact significant digit.

Similarly, when the expression `a >= b` is true, it means that

- `a-b` is a computational zero or
- $\frac{\sum_{i=1}^N a_i}{N} > \frac{\sum_{i=1}^N b_i}{N}$,

and, when the expression $a > b$ is true, it means that

- $a-b$ is NOT a computational zero, i.e. has at least one exact significant digit, and
- $\frac{\sum_{i=1}^N a_i}{N} > \frac{\sum_{i=1}^N b_i}{N}$.

2.5 SAM specific functions

The previous part described how some classical C statements are slightly affected when using the SAM tool. Now we present functions that are specific to the library. Note that the functions `SAM_init` and `SAM_end` have to appear, respectively to initialize and to close the library. The other functions `SAM_enable`, `SAM_disable`, `self_validation_only`, `computedzero`, `data_st`, `nb_significant_digit`, `str` and `strp` will appear in some applications.

2.5.1 Initializing and closing the library

The `SAM_init` function has to be called once, early in the main program, before any kind of declaration.

This function has five integer arguments:

`SAM_init(numb_instability, SAM_instability, cancel_level, init_random, precision)`.

The first and the last argument must always be present.

The user chooses the maximum number of numerical instabilities that will be detected.

- if `numb_instability` = -1, all the instabilities will be detected
- if `numb_instability` = 0, no instability will be detected
- if `numb_instability` = M (strictly positive M), the first M instabilities will be detected.

The last argument is an integer which represents the precision in bits.

For example, a number presented using 53 bits has 15 decimal significant digits.

The other arguments are optional.

The second argument allows the user to determine what kind of instabilities will be enabled or disabled.

There are 7 integer parameters in the library:

SAM_BRANCHING,
SAM_CANCEL,
SAM_DIV,
SAM_INTRINSIC,
SAM_MATH,
SAM_MUL,
SAM_POWER.

By default, the detection of all types of instability is enabled. The user has only to specify what kind of instability is to be **disabled** by passing, as the second argument, the addition of the chosen parameters.

The third argument is an integer which is used to initialize some internal variables for random arithmetic. The default value for this argument is 51.

The fourth argument corresponds to the following. An unstable cancellation is pointed out when the difference between the number of exact significant digits (i.e. digits which are not affected by round-off errors) of the result of an addition or a subtraction and the minimum of the number of exact significant digits of the two operands is greater than the `cancel_level` argument. The default value of this argument is 4. In other words, when one loses more than `cancel_level` significant digits in one addition or subtraction, SAM considers that a catastrophic cancellation has been detected (if the detection of this kind of instability is enabled).

The `SAM_end` function “closes” the library and prints to the standard output the result of the detection of numerical instabilities.

2.5.2 Obtaining a string from a result with its evaluated accuracy

The `str` function has a string argument and a stochastic argument. It returns a pointer to the first argument. This output string contains the scientific notation of the stochastic argument; only the exact significant digits appear in the string. Thus accuracy is easy to read.

When the argument has no exact significant digit, the string that is returned is `@.0`.

To avoid the use of the string parameter, a special implementation of the `str` function has been written. It must be used only with the family of `printf`

functions. The name of this new function is `strp`. Using this function, the allocation of the string is dynamically managed by the function itself. The only restriction is that it is not possible to have more than 256 calls to the `strp` function in one call to the `printf` function.

For C++ programmers, the classical `c.out` and `>>` notations have been overloaded for the stochastic types. No modification is needed.

Example:

Let us consider the following instructions:

```
int i;
double a;
...
printf("iteration %d a=%lf\n",i,a);
...
```

Using SAM, the corresponding instructions in C can be:

```
int i;
mp_st a;
char s[25];
...
printf("iteration %d a=%s\n",i,str(s,a));
...
```

Or using the `strp` function:

```
int i;
mp_st a;
...
printf("iteration %d a=%s \n",i,strp(a));
...
```

In C++, it is simpler:

```
int i;
mp_st a;
...
cout << "iteration " << i << "a=" << a << endl;
...
```

2.5.3 Obtaining the number of exact significant digits of a stochastic variable

The `nb_significant_digit` method returns an integer giving the number of exact significant decimal digits of a stochastic variable when the method is called. At some point `x.nb_significant_digit()` may return 7; later during the run

it may return 5. If x becomes non-significant then `x.nb_significant_digit()` returns 0.

2.5.4 Obtaining the triplet associated with a stochastic variable

The `display` method prints the triplet associated with a stochastic variable. For instance, let `d` be a multiple-precision stochastic variable. The following instructions

```
printf("%s\n",strp(d));  
d.display();
```

may provide

```
0.30E-64  
3.0217133019536030e-65 -- 3.0133146666062181e-65 -- 3.0248565827034563e-65
```

The three multiple-precision values associated to `d` have 2 common significant digits.

2.5.5 Testing if a variable is a computational zero

The `computedzero` method acts on a stochastic variable and returns 0 or 1. The `computedzero` method returns 1 if this stochastic variable is a computational zero, *i.e.* it is a mathematical zero or it has no exact significant digit.

2.5.6 Reducing accuracy of initial data

Initial data are often known with less significant digits than provided by their internal representation. The `data_st` method allows the user to introduce some effective uncertainties on these data, reducing their initial accuracy. So the accuracy of results depends in some way on the accuracy of initial data.

The `data_st` method acts on a stochastic variable `X` and has two optional arguments: `X.data_st(ERX,IER)`;

The first argument is an optional `double` argument that contains the relative or absolute uncertainty of the stochastic variable `X`. The second argument determines the kind of the uncertainty: relative or absolute. If `X` is a stochastic variable and `ERX` is a `double` value strictly less than 1, the `X.data_st(ERX,IER)`; instruction modifies the values of the N samples in `X` according to the following formula:

$X_i = X_i * (1 + ERX * ALEA)$ for $i = 1$ to N if $IER = 0$

$X_i = X_i + ERX * ALEA$ for $i = 1$ to N if $IER = 1$

ALEA is a random variable uniformly distributed between -1 and 1.
If *ERX* is 0, no perturbation takes place as if the statement was suppressed.
If *ERX* is absent, perturbation will concern only the last bit of the mantissa.
If *IER* is absent, it is like *IER* = 0. The `data_st` method without `ERX` must be used when data are considered as exact but cannot be exactly coded in the memory.

Chapter 3

User's guide

The use of the SAM library involves seven steps:

- declaration of the SAM library for the compiler,
- initialization of the SAM library,
- substitution of the type `float` or `double` by stochastic type `mp_st` in variable declarations,
- possible changes in the input data if perturbation is desired, to take into account uncertainty in initial values,
- change of output statements to print stochastic results with their accuracy,
- possible use of SAM functions to evaluate the number of exact significant digits,
- termination of the SAM library.

The reader may refer to the sample program given in 3.9 with two versions, *i.e.* the initial C code and the code modified to be compiled with the SAM library.

3.1 Declaration of the SAM library

The following pseudo-statement

```
#include <SAM.h>
```

must take place in any file which contains declarations of stochastic variables or SAM functions to be found by the compiler.

3.2 Initialization of the SAM library

The `SAM_init` function has to be called once, early in the main program, before any kind of declaration, to initialize random arithmetic and to set the default precision. The choice of the default precision is global.

Note that `double` variables are 53-bit mantissa length numbers and that `float` variables are 24-bit mantissa length numbers in IEEE standard floating-point arithmetic [39].

For more information about the arguments of the `SAM_init` function, see 2.5.1.

3.3 Declaration of variables

3.3.1 Changes in the type of variables

To control the numerical quality of a variable, just replace its standard type by the stochastic type.

Example:

standard declarations	SAM declarations
<code>float a, b;</code>	<code>mp_st a, b;</code>
<code>double c;</code>	<code>mp_st c;</code>
<code>float d[6], e, f;</code>	<code>mp_st d[6], e, f;</code>

3.4 Changes in assignments or arithmetic operations

3.4.1 Conversions between usual types and the stochastic type

In assignment statements, conversions are implicit from C `float`, `int` or `double` types *to* and *from* the stochastic type `mp_st` (because the `=` operator has been overloaded), **but for conversions from the stochastic type `mp_st` to standard types, the knowledge of accuracy is lost.**

With the following declarations:

```
mp_st a, b;  
float r;
```

the assignments `a = r;`, `b = 2;` and `r = a;` are correct but there is, of course, no information on the accuracy of `r`.

When a variable is set to a value which cannot be exactly coded on a computer, the `sam_set_str` function should be used.

Example:

Initial C statements	Modified C statements for SAM
float x, y; x=1.234; y=-3.0;	#include <SAM.h> mp_st x, y; sam_set_str(x, "1.234"); y=-3.0;

3.4.2 Classical arithmetic operators

As previously described, all arithmetic operators on floating-point variables are overloaded and arithmetic expressions without functions do not have to be modified. Expressions may contain a mixture of the stochastic type, classical types and integer types.

With the following declarations:

`mp_st a, b;`

`double c;`

the statement `a = b * c + 3;` needs no change.

The result of expressions containing stochastic terms will be of stochastic type.

3.5 Changes in reading statements

The family of `scanf` functions is adapted to classical floating-point variables, which must be transformed into stochastic variables.

Example:

Initial C statements	Modified C statements for SAM
float x; scanf("x = %14.7e \n", &x);	#include <SAM.h> chariaux[100]; mp_st x; scanf("%s", iaux); sam_set_str(x, iaux);

Note that initial data read from a file or from keyboard may have sometimes to be duplicated in some way, because they are read as classical variables which are then assigned to stochastic variables with controlled uncertainty.

3.6 Changes in printing statements

Before printing each stochastic variable, it must be transformed into a string by the `str` or `strp` function. The required length is 15 for a `mp_st` variable of 24 bits precision and 25 for a `mp_st` variable of 53 bits precision. Therefore formats should be modified.

For example, if a `float` variable `x` becomes a `mp_st` variable, the printing instruction can be modified as follows:

Initial C statements	Modified C statements for SAM
float x; ... printf("x = %14.7e", x);	#include <SAM.h> mp_st x; ... printf("x = %s", strp(x));

3.7 Constants passed as function arguments

Function definitions and function calls must sometimes be adapted because stochastic parameters of functions must not be passed by value.

Example:

Initial C statements	Modified C statements for SAM
float a; a=3.14*f(2.0); ... float f(x) { float x; ... }	#include <SAM.h> mp_st aux, a; aux=2.0; a=3.14*f(aux); ... mp_st f(x) { mp_st x; ... }

3.8 Termination of the SAM library

The call to the `SAM_end` function must be the last program statement.

3.9 An example of numerical code and its modified version

The following source codes use the Gauss-Jordan method to invert a matrix.

3.9.1 Standard C source code

```
#include <stdio.h>
#define N 4

// Initialization:
void InitMat(float M[N][N]){
    int i,j;
    for(i=0;i<N;i++)
        for(j=0;j<N;j++) scanf("%e",&M[i][j]);
}

// Inversion using the Gauss-Jordan method:
void InvertMat(float M[N][N]){
    int i,j,k;
    float temp;
    for(k=0;k<N;k++)
        {temp = M[k][k];
        M[k][k] = 1.0;
        for(j=0;j<N;j++) M[k][j]/=temp;
        for(i=0;i<N;i++)
            if(i!=k)
                {temp=M[i][k];
                M[i][k] = 0.0;
                for(j=0;j<N;j++) M[i][j] -= temp*M[k][j];
                };
        };
}

// Display of a matrix:
```

```

void DisplayMat(float M[N][N]){
    int i,j;
    for(i=0;i<N;i++)
        {for(j=0;j<N;j++)
            printf("%14.7e  ",M[i][j]);
            printf("\n");
        }
}

```

```

void main(){
    float M[N][N];
    printf("Initial matrix:\n");
    InitMat(M);
    DisplayMat(M);
    InvertMat(M);
    printf("Inverted matrix:\n");
    DisplayMat(M);
}

```

3.9.2 Source code using the SAM library

```

#include "SAM.h"
#include <stdio.h>
#define N 4

// Initialization:
void InitMat(mp_st M[N][N]){
    char aux[100];
    int i,j;
    for(i=0;i<N;i++)
        for(j=0;j<N;j++) {scanf("%s",aux); sam_set_str(M[i][j],aux);};
}

// Inversion using the Gauss-Jordan method:
void InvertMat(mp_st M[N][N]){
    int i,j,k;
    mp_st temp;
    for(k=0;k<N;k++)
        {temp = M[k][k];
        M[k][k] = 1.0;

```

```

        for(j=0;j<N;j++) M[k][j]/=temp;
    for(i=0;i<N;i++)
        if(i!=k)
            {temp=M[i][k];
             M[i][k]=0.0;
             for(j=0;j<N;j++) M[i][j] -=temp*M[k][j];
            };
    };
}

// Display of a matrix:
void DisplayMat(mp_st M[N][N]){
    int i,j;
    for(i=0;i<N;i++)
        {for(j=0;j<N;j++)
            printf("%s ",strp(M[i][j]));
          printf("\n");
        }
}

main(){
    SAM_init(-1,24);
    mp_st M[N][N];
    printf("Initial matrix:\n");
    InitMat(M);
    DisplayMat(M);
    InvertMat(M);
    printf("Inverted matrix:\n");
    DisplayMat(M);
    SAM_end();
}

```

3.9.3 Example of execution without SAM

Initial matrix:

1.0000000e+00	2.0000000e+03	5.0000000e-01	4.0000000e+00
2.9999999e-05	1.0000000e+00	2.0000000e+00	8.0000000e+00
4.0000000e+00	5.0000000e-01	2.9999999e-08	2.0000000e+00
2.0000000e+00	3.0000000e+00	5.0000000e-01	5.0000000e+09

Inverted matrix:

```

-6.2576764e-05  -8.1558341e-05   2.5001565e-01  -9.9964599e-11
 5.0009380e-04  -1.2504448e-04  -1.2502252e-04  -1.4995290e-13
-2.5004597e-04   5.0006253e-01   5.8761423e-05  -7.9992352e-10
-2.4999515e-13  -4.9991469e-11  -9.9937121e-11   2.0000000e-10

```

3.9.4 Example of execution with SAM

```

Initial matrix:
0.1000000E+1   0.1999999E+4   0.5000000   0.4000000E+1
0.300000E-4    0.1000000E+1   0.2000000E+1  0.8000000E+1
0.4000000E+1   0.5000000   0.2999999E-7  0.2000000E+1
0.2000000E+1   0.3000000E+1   0.5000000   0.500000E+10

Inverted matrix:
-0.63E-4      @.0          0.250016   -0.10E-9
 0.500094E-3  -0.124812E-3  -0.125023E-3  -0.15E-12
-0.2500459E-3  0.500063     0.587614E-4  -0.799924E-9
-0.249E-12    -0.49E-10    -0.999371E-10  0.200000E-9

```

3.10 Numerical debugging with SAM

One can enable the detection of the following instabilities:

```

UNSTABLE DIVISION(S),
UNSTABLE POWER FUNCTION(S),
UNSTABLE MULTIPLICATION(S),
UNSTABLE BRANCHING(S),
UNSTABLE MATHEMATICAL FUNCTION(S),
UNSTABLE INTRINSIC FUNCTION(S),
LOSS OF ACCURACY DUE TO CANCELLATION(S).

```

The library counts the number of detections for each instability. The global information for these detections is printed out with the `SAM_end` function, see 2.5.1.

The accuracy estimated by SAM is valid if there is no deep numerical anomaly during the computation, i.e. no `UNSTABLE DIVISION`, `UNSTABLE POWER FUNCTION` and `UNSTABLE MULTIPLICATION`, see [23, 15, 9].

The meaning of the message is:

- **unstable division:** the divisor is non-significant

- **unstable power function:** one operand of the power function is non-significant
- **unstable multiplication:** both operands are non-significant
- **unstable branching:** the difference between the two operands is non-significant (a computational zero).
The chosen branching statement is associated with the equality.
- **unstable mathematical function:**
in the `log`, `sqrt`, `exp` or `log10` function, the argument is non-significant.
- **unstable intrinsic function:**
 - when using integer cast functions, the integral part of the argument can not be exactly determined due to the round-off error propagation;
 - in the `fabs` function: the argument is non-significant;
 - the `floor`, `ceil` or `trunc` function returns different values for each component of the stochastic argument.
- **loss of accuracy due to cancellation:** as explained in 2.5.1, an unstable cancellation is pointed out when the difference between the number of exact significant digits (i.e. digits which are not affected by round-off errors) of the result of an addition or a subtraction and the minimum of the number of exact significant digits of the two operands is greater than the `cancel_level` argument. The default value of this argument is 4. In other words, when one loses more than `cancel_level` significant digits in one addition or subtraction, SAM considers that a catastrophic cancellation has been detected (if the detection of this kind of instability is enabled).

To perform actual numerical debugging, it is necessary, for each instability, to identify the statement in the code that generates this instability. This can be performed directly using a symbolic debugger like **`gdb`** with Linux or as a background task using special input and output files.

In both cases, one has to put a breakpoint at the entry of the **`instability`** internal function of the SAM library. This function is called each time a numerical instability is detected. To get the right label for this system and compiler dependent function, one can use the following statement:

```
nm name_of.the.binary.code | grep instability
```

For instance, using **gdb** with Linux, the general statement which enables the detection of all the instabilities in a single run is

```
nohup gdb name_of_the_binary_code < gdb.in >! gdb.out &
```

The *gdb.in* file may contain

```
break instability
run
while 1
where
cont
end
```

where prints out the complete trace of the instability which has stopped the run and **cont** makes the execution going on.

P.S.: **nohup** allows to keep the process alive even when logging off.

The *gdb.out* file will contain all the traces of instabilities.

Chapter 4

Test runs

You first need to install GMP (GNU Multiprecision Library), MPFR and MPFI for running all the tests.

We present, with the examples included in the distribution, an illustration of the use of the SAM library and the benefits of the DSA. For each example, we describe the results obtained using the standard floating-point arithmetic and then the results provided by the SAM library.

The results reported in this section have been obtained using the darwin g++ compiler on an Macbook pro running MAC OS. Different results may be obtained with another processor or another compiler, especially when the digits printed out using the standard floating-point arithmetic are affected by round-off errors. With SAM, only the exact significant digits appear in the results.

4.1 Example 1: a rational fraction function of two variables

In the following example [26], the rational fraction

$$F(x, y) = 333.75y^6 + x^2(11x^2y^2 - y^6 - 121y^4 - 2) + 5.5y^8 + \frac{x}{2y}$$

is computed with $x = 77617$, $y = 33096$. The first 30 digits of the exact result are -0.827396059946821368141165095479.

Using IEEE double precision arithmetic with rounding to the nearest, one obtains: res = 5.764607523034235E+17.

With SAM, we obtain the exact result when the default precision is set to more than 121 bits. So using SAM in 122 bits, one obtains:

```

-----
SAM software --- University P. et M. Curie --- LIP6
Self-validation detection: ON
Mathematical instabilities detection: ON
Branching instabilities detection: ON
Intrinsic instabilities detection: ON
Cancellation instabilities detection: ON
-----

```

```

-----
| Polynomial function of two variables |
| with SAM                            |
-----

```

```

res=-0.827396059946821368141165095479816292
-----

```

```

SAM software --- University P. et M. Curie --- LIP6
No instability detected
-----

```

4.2 Example 2: solving a second order equation

The roots of the following second order equation are computed:

$$0.3x^2 - 2.1x + 3.675 = 0.$$

The exact results are: Discriminant $d=0$, $x_1=x_2=3.5$.

Using IEEE single precision arithmetic with rounding to the nearest, one obtains:

```

-----
| Second order equation          |
| without SAM                    |
-----

```

```

d = -2.861023e-06

```

There are two complex solutions.

```

z1 = +3.500000e+00 + i * +8.457279e-04

```

```

z2 = +3.500000e+00 + i * -8.457279e-04

```

and using SAM in 100 bits, one obtains:

```

-----
SAM software --- University P. et M. Curie --- LIP6

```


SAM software --- University P. et M. Curie --- LIP6
 No instability detected

4.4 Example 4: computing a second order recurrent sequence

This sequence was proposed by J.-M. Muller [22]. The first 30 iterations of the following recurrent sequence are computed:

$$U_{n+1} = 111 - \frac{1130}{U_n} + \frac{3000}{U_n U_{n-1}}$$

with $U_0 = 5.5$ and $U_1 = \frac{61}{11}$. The exact limit is 6.

Using IEEE double precision arithmetic with rounding to the nearest, one obtains:

 | A second order recurrent sequence |
without SAM

U(3) = +5.590163934426237e+00
 U(4) = +5.633431085044127e+00
 U(5) = +5.674648620512615e+00
 U(6) = +5.713329052423919e+00
 U(7) = +5.749120920462043e+00
 U(8) = +5.781810933690098e+00
 U(9) = +5.811314466602178e+00
 U(10) = +5.837660476543959e+00
 U(11) = +5.861018785996283e+00
 U(12) = +5.882524608269310e+00
 U(13) = +5.918655323805488e+00
 U(14) = +6.243961815306110e+00
 U(15) = +1.120308737284091e+01
 U(16) = +5.302171264499677e+01
 U(17) = +9.473842279276452e+01
 U(18) = +9.966965087355071e+01
 U(19) = +9.998025776093678e+01
 U(20) = +9.999882245337588e+01
 U(21) = +9.999992970745579e+01

```

U(22) = +9.999999580049865e+01
U(23) = +9.999999974893262e+01
U(24) = +9.99999998498109e+01
U(25) = +9.99999999910112e+01
U(26) = +9.9999999994618e+01
U(27) = +9.999999999677e+01
U(28) = +9.999999999980e+01
U(29) = +9.999999999999e+01
U(30) = +1.00000000000000e+02
The exact limit is 6.

```

and using SAM in double precision (53 bits), one obtains:

```

-----
SAM software --- University P. et M. Curie --- LIP6
Self-validation detection: ON
Mathematical instabilities detection: ON
Branching instabilities detection: ON
Intrinsic instabilities detection: ON
Cancellation instabilities detection: ON
-----

```

```

-----
| A second order recurrent sequence |
| with SAM                           |
-----

```

```

U(3) = 0.5590163934426E+1
U(4) = 0.563343108504E+1
U(5) = 0.56746486205E+1
U(6) = 0.571332905E+1
U(7) = 0.57491209E+1
U(8) = 0.5781811E+1
U(9) = 0.581131E+1
U(10) = 0.58376E+1
U(11) = 0.586E+1
U(12) = 0.59E+1
U(13) = 0.6E+1
U(14) = @.0
U(15) = @.0
U(16) = @.0
U(17) = @.0
U(18) = 0.9E+2

```

```

U(19) = 0.99E+2
U(20) = 0.999E+2
U(21) = 0.99999E+2
U(22) = 0.999999E+2
U(23) = 0.9999999E+2
U(24) = 0.99999999E+2
U(25) = 0.999999999E+2
U(26) = 0.9999999999E+2
U(27) = 0.99999999999E+2
U(28) = 0.999999999999E+2
U(29) = 0.9999999999999E+2
U(30) = 0.100000000000000E+3
The exact limit is 6.

```

SAM software --- University P. et M. Curie --- LIP6

CRITICAL WARNING: the self-validation detects major problem(s).
The results are NOT guaranteed.

There are 12 numerical instabilities
9 UNSTABLE DIVISION(S)
3 UNSTABLE MULTIPLICATION(S)

and using SAM with 100 bits, one obtains:

SAM software --- University P. et M. Curie --- LIP6
Self-validation detection: ON
Mathematical instabilities detection: ON
Branching instabilities detection: ON
Intrinsic instabilities detection: ON
Cancellation instabilities detection: ON

```

| A second order recurrent sequence |
| with SAM                          |

```

```

U(3) = 0.559016393442622950819672131E+1
U(4) = 0.56334310850439882697947214E+1
U(5) = 0.5674648620510150963040083E+1

```

U(6) = 0.571332905238051554903219E+1
 U(7) = 0.5749120919702638043705E+1
 U(8) = 0.578181092048561557947E+1
 U(9) = 0.58113142382939957232E+1
 U(10) = 0.5837656548958711962E+1
 U(11) = 0.586095152251613197E+1
 U(12) = 0.5881377215841419E+1
 U(13) = 0.589915390579007E+1
 U(14) = 0.59145249506789E+1
 U(15) = 0.5927741407777E+1
 U(16) = 0.59390504855E+1
 U(17) = 0.5948687492E+1
 U(18) = 0.595687073E+1
 U(19) = 0.59637987E+1
 U(20) = 0.596965E+1
 U(21) = 0.59746E+1
 U(22) = 0.5979E+1
 U(23) = 0.598E+1
 U(24) = 0.59E+1
 U(25) = @.0
 U(26) = @.0
 U(27) = @.0
 U(28) = @.0
 U(29) = 0.9E+2
 U(30) = 0.99E+2
 U(31) = 0.999E+2
 U(32) = 0.99999E+2
 U(33) = 0.999999E+2
 U(34) = 0.9999999E+2
 U(35) = 0.99999999E+2
 U(36) = 0.999999999E+2
 U(37) = 0.9999999999E+2
 U(38) = 0.99999999999E+2
 U(39) = 0.999999999999E+2
 U(40) = 0.9999999999999E+2
 U(41) = 0.99999999999999E+2
 U(42) = 0.999999999999999E+2
 U(43) = 0.9999999999999999E+2
 U(44) = 0.9999999999999999E+2
 U(45) = 0.99999999999999999E+2

Using IEEE double precision arithmetic with rounding to the nearest, one obtains:

```
-----
| Computation of a root of a polynomial by Newton's method |
| without SAM                                              |
-----
```

```
x( 66) = +4.285714325273430e-01
x( 67) = +4.285714325273430e-01
```

and using SAM in 100 bits (30 decimal digits), one obtains:

```
-----
SAM software --- University P. et M. Curie --- LIP6
Self-validation detection: ON
Mathematical instabilities detection: ON
Branching instabilities detection: ON
Intrinsic instabilities detection: ON
Cancellation instabilities detection: ON
-----
```

```
-----
| Computation of a root of a polynomial by Newton's method |
| with SAM                                                  |
-----
```

```
x( 36) =  0.428571428572503043
x( 37) =  0.428571428571965807
```

```
-----
SAM software --- University P. et M. Curie --- LIP6
There are 87 numerical instabilities
87 LOSS OF ACCURACY DUE TO CANCELLATION(S)
-----
```

With SAM, one can see that 12 significant decimal digits were lost. SAM allows to stop the algorithm when the subtraction $x_n - x_{n-1}$ is insignificant (there is no more information to compute at the next iteration). In Newton's method, a division by a computational zero may suggest a double root. One can simplify the fraction. When these two transformations are performed, the code is stabilized and the results are obtained with the best accuracy of the computer. The exact value of the root is $x_{sol} = 3/7 = 0.428571428571428571...$ Now, we obtain:

SAM software --- University P. et M. Curie --- LIP6

Self-validation detection: ON

Mathematical instabilities detection: ON

Branching instabilities detection: ON

Intrinsic instabilities detection: ON

Cancellation instabilities detection: ON

```
-----
| Computation of a root of a polynomial by Newton's method |
| with optimisation and SAM                               |
-----
```

x(95) = 0.42857142857142857142857142857

x(96) = 0.42857142857142857142857142857

SAM software --- University P. et M. Curie --- LIP6

No instability detected

4.6 Example 6: solving a linear system

In this example, SAM is able to provide correct results which were impossible to be obtained with the standard floating-point arithmetic. The following linear system $Ax = b$ is solved using Gaussian elimination with partial pivoting.

$$\begin{pmatrix} 21 & 130 & 0 & 2.1 \\ 13 & 80 & 4.74 \cdot 10^8 & 752 \\ 0 & -0.4 & 3.9816 \cdot 10^8 & 4.2 \\ 0 & 0 & 1.7 & 9 \cdot 10^{-9} \end{pmatrix} \cdot x = \begin{pmatrix} 153.1 \\ 849.74 \\ 7.7816 \\ 2.6 \cdot 10^{-8} \end{pmatrix}$$

The exact solution is $x_{sol}^t = (1, 1, 10^{-8}, 1)$.

Using IEEE single precision arithmetic with rounding to the nearest, one obtains:

```
-----
| Solving a linear system using Gaussian elimination |
| by partial pivoting without SAM                     |
-----
```

x_sol(0) = +6.261988e+01 (exact solution: xsol(0)= +1.000000e+00)

x_sol(1) = -8.953979e+00 (exact solution: xsol(1)= +1.000000e+00)

```
x_sol(2) = +0.000000e+00 (exact solution: xsol(2)= +1.000000e-08)
x_sol(3) = +1.000000e+00 (exact solution: xsol(3)= +1.000000e+00)
```

and using SAM in 53 bits, one obtains:

```
-----
SAM software --- University P. et M. Curie --- LIP6
Self-validation detection: ON
Mathematical instabilities detection: ON
Branching instabilities detection: ON
Intrinsic instabilities detection: ON
Cancellation instabilities detection: ON
-----
```

```
-----
| Solving a linear system using Gaussian elimination |
| with partial pivoting                               |
-----
```

```
x_sol(0) = 0.100000000000E+1 (exact solution: xsol(0)= 0.1000000000000000E+1)
x_sol(1) = 0.999999999999 (exact solution: xsol(1)= 0.1000000000000000E+1)
x_sol(2) = 0.1000000000000000E-7 (exact solution: xsol(2)= 0.1000000000000000E-7)
x_sol(3) = 0.9999999999999999 (exact solution: xsol(3)= 0.1000000000000000E+1)
-----
```

```
SAM software --- University P. et M. Curie --- LIP6
There is 1 numerical instability
1 LOSS OF ACCURACY DUE TO CANCELLATION(S)
-----
```

Using standard floating-point arithmetic, during the reduction of the third column, the matrix element $A(3,3)$ is equal to 4864. But the exact value of $A(3,3)$ is zero. The standard floating-point arithmetic cannot detect that $a(3,3)$ is insignificant. This value is chosen as pivot. That leads to erroneous results. SAM detects the non-significant value of $A(3,3)$. This value is eliminated as pivot. That leads to satisfactory results.

4.7 Example 7: when SAM fails

SAM is based on a probabilistic model. It should never be forgotten that all the estimations computed by SAM are probabilistic, even if the probability is close to 1. Moreover, the theoretical model shows that SAM is able to estimate the round-off errors to the first order. If they represent the global

round-off errors, SAM works well but, if they are dominated by terms of greater order, SAM may fail.

In the present example, we have the same behaviour but only with additions and subtractions, so without any warning of numerical instability. Let us perform the following computation:

```
x=6.83561e+5;
y=6.83560e+5;
z=1.00000000007;
r = z - x;
r1 = z - y;
r = r + y;
r1 = r1 + x;
r1 = r1 - 2;
r = r + r1;
//      r = ((z-x)+y) + ((z-y)+x-2)
```

The exact result is $1.4 \cdot 10^{-10}$. The result obtained using IEEE double precision arithmetic with rounding to the nearest is 2.32830643653870E-10.

With SAM in double precision (53 bits), because we essentially performed the same computation, $((z - x) + y)$ and $((z - y) + x - 2)$, we find that if the same rounding mode is chosen for both parts, the final result appears as exact but it is wrong. It happens in one case in three and the result provided by SAM is then 0.116415321826935E-009 with 15 exact significant digits. If computations are performed 100000 times using SAM, one may obtain:

```
-----
SAM software --- University P. et M. Curie --- LIP6
Self-validation detection: ON
Mathematical instabilities detection: ON
Branching instabilities detection: ON
Intrinsic instabilities detection: ON
Cancellation instabilities detection: ON
-----
```

```
Enter the number of iterations: 100000
r = @.0 ierr = 37557
-----
```

```
SAM software --- University P. et M. Curie --- LIP6
There are 150790 numerical instabilities
150790 LOSS OF ACCURACY DUE TO CANCELLATION(S)
```

The last value of `r` is printed out, and also `ierr` the number of times when the result was wrong.

Bibliography

- [1] J.-M. Chesneaux, F. Jézéquel, J.-L. Lamotte, Stochastic arithmetic and verification of mathematical models In *Uncertainties in environmental modelling and consequences for policy making*, P. Baveye, J. Mysiak, M. Laba Eds., NATO Science for Peace and Security Series - C: Environmental Security, Springer, pages 101-125, 2009.
- [2] J.-M. Chesneaux, S. Graillat, F. Jézéquel, Numerical validation and assessment of numerical accuracy, Oxford e-Research Centre, overview article, 44 pages, march 2009.
http://cpc.cs.qub.ac.uk/oerc_numerical_accuracy.pdf
- [3] J.-M. Chesneaux, S. Graillat, F. Jézéquel, Rounding errors, invited paper, In *Wiley Encyclopedia of Computer Science and Engineering* (Benjamin Wah, ed.) Hoboken: John Wiley & Sons, vol. 4, pages 2480-2494, 2009.
- [4] F. Jézéquel, J.-M. Chesneaux, CADNA: a library for estimating round-off error propagation, *Computer Physics Communications*, 178(12), pages 933-955, 2008.
- [5] N.S. Scott, F. Jézéquel, C. Denis, J.-M. Chesneaux, Numerical 'health check' for scientific codes: the CADNA approach, *Computer Physics Communications*, 176(8), pages 507-521, 2007.
- [6] F. Jézéquel, A dynamical strategy for approximation methods, *C. R. Acad. Sci. Paris - Mécanique*, 334, pages 362-367, 2006.
- [7] F. Jézéquel, F. Rico, J.-M. Chesneaux, M. Charikhi, Reliable computation of a multiple integral involved in the neutron star theory, *Mathematics and Computers in Simulation*, 71(1), pages 44-61, 2006.
- [8] F. Jézéquel, Dynamical control of approximation methods, Habilitation à diriger des recherches, Université Pierre et Marie Curie, Paris, 2005.

- [9] J. Vignes, Discrete Stochastic Arithmetic for Validating Results of Numerical Software, *Special Issue of Numerical Algorithms*, 2004, 37, pp. 377-390
- [10] F. Jézéquel, J.-M. Chesneaux, Computation of an infinite integral using Romberg's method, *Numerical Algorithms*, 36 (3): 265-283, July 2004.
- [11] F. Jézéquel, Dynamical control of converging sequences computation, *Applied Numerical Mathematics*, 50(2): 147-164, 2004.
- [12] F. Jézéquel, J.-M. Chesneaux, For reliable and powerful scientific computations, *Sc. Comp. Val. Num.*, Krämer and Wolff von Gudenberg ed., Kluwer Academic/Plenum publishers (2001) 367-378,
- [13] J.-M. Chesneaux, F. Jézéquel, Dynamical control of computations using the Trapezoidal and Simpson's rules *Journal of Universal Computer Science*, Vol. 4 (1), 2-10, 1998.
- [14] M. Pichat, J. Vignes, Validité des résultats numériques dans les processus à comportement chaotique. Un outil d'évaluation : le logiciel CADNA. *CRAS*, Paris, Tome 322, Série 2b, 1996, pp. 681-688.
- [15] J.-M. Chesneaux, L'arithmétique stochastique et le logiciel CADNA, Habilitation à diriger des recherches, Université Pierre et Marie Curie, Paris, 1995.
- [16] J.-M. Chesneaux, The equality relations in scientific computing, *Num. Algo* 7, 1994, pp. 129-143.
- [17] M. Pichat, J. Vignes, Ingénierie du contrôle de la précision des calculs sur ordinateur. *Ed. Technip*, Paris 1993.
- [18] J. Vignes, A stochastic arithmetic for reliable scientific computation, *Math. and Comp. in Sim.* 35, 1993, pp. 233-261.
- [19] J.-M. Chesneaux, J. Vignes, Les fondements de l'arithmétique stochastique, *C.R Acad. Sci.*, Paris, Sér.I, 315, 1992, pp. 1435-1440.
- [20] J.-M. Chesneaux, Study of the computing accuracy by using probabilistic approach, *Contribution to comp. arithmetic and Self-Validating Numerical Methods*, C. Ullrich ed., IMACS, New Brunswick, NJ, 1990, pp. 19-30.

- [21] J. Vignes, Estimation de la précision des résultats de logiciels numériques. *La Vie des Sciences, Comptes Rendus, série générale*, 7, 1990, pp. 93-143.
- [22] J.-M. Muller, Arithmétique des ordinateurs, Masson, 1989.
- [23] J.-M. Chesneaux, Étude théorique et implémentation en ADA de la méthode CESTAC, *Thèse de l'université P. et M. Curie*, Paris, 1988.
- [24] J.-M. Chesneaux, J. Vignes, Sur la robustesse de la méthode CESTAC, *C.R. Acad. Sc. Paris, Sér. I Math.* 307, 1988, pp. 855-860.
- [25] J.-M. Chesneaux, Modélisation théorique et conditions de validité de la méthode CESTAC, *C.R.A.S., Paris, série 1, tome 307*, 1988 pp. 417-422.
- [26] S.M. Rump, Algorithms for Verified Inclusions - Theory and Practice. In R.E. Moore, editor, Reliability in Computing, volume 19 of Perspectives in Computing, pages 109-126. Academic Press, 1988.
- [27] J. Vignes, Zéro mathématique et zéro informatique. *La Vie des Sciences, C.R. Acad. Sci.*, Paris, 4, 1, janvier 1987, pp. 1-13.
- [28] J. Vignes, Implémentation des méthodes d'optimisation : test d'arrêt optimal, contrôle et précision de la solution (I) *R.A.I.R.O.*, 18, 1, février 1984, pp. 1-18; (II), *R.A.I.R.O.* 18, 2, mai 1984, pp. 103-129.
- [29] J. Vignes, M. La Porte, Error analysis in computing, *Information Processing 74*, North-Holland, 1974.
- [30] L.S. Yao, Computed chaos or numerical errors. *Nonlinear Analysis: Modelling and control* 15(1), 2010, pp. 109-126
- [31] Devaney, R.L.: An introduction to chaotic dynamical systems. Second edn. Addison-Wesley Studies in Nonlinearity. Addison-Wesley Publishing Company Advanced Book Program, Redwood City, CA, 1989
- [32] D. Stott Parker, Monte Carlo Arithmetic: Exploiting Randomness in floating-point arithmetic. Report of Computer Science Department, UCLA Los Angeles, March 30, 1997.
- [33] D. Stott Parker, Brad Pierce, P.R. Eggert, Monte Carlo Arithmetic: How to Gamble with Floating Point and Win. *Computing in Science and Engineering* 2000, pp. 58-68.

- [34] J. Vignes, Error analysis in computing. International Federation for Information Processing Congress, Stockholm, aug. 1974, pp. 610-614.
- [35] J. Vignes, New Methods for evaluating the validity of the results of mathematical computations. Math. and Comp. in Sim., 20, 1978, pp. 227-249.
- [36] J. Vignes, A Stochastic Approach to the Analysis of Round-off Error Propagation. A Survey of the CESTAC Method. Proceedings of Real Numbers and Computer Conference. Marseille, 1996 pp. 233-251.
- [37] W. Tucker, A rigorous ODE solver and Smale's 14th problem, Foundations of computational Mathematics, 2(1), pp.53-117, 2002
- [38] S. Smale, Mathematical problems for the next century, Math. Intell., 20, pp.7-15, 1998
- [39] IEEE Standard for Binary Floating-Point Arithmetic, ANSI/IEEE Standard 754-1985, Institute of Electrical and Electronics Engineers, August, 1985, reprinted in SIGPLAN 22, 2, pp. 9-25.