

# Static Analysis of the Accuracy in Control Systems : Principles and Experiments

Eric Goubault<sup>1</sup>, Sylvie Putot<sup>1</sup>, Philippe Baufreton<sup>2</sup>, and Jean Gassino<sup>3</sup>

<sup>1</sup> CEA/LIST 91191 Gif-sur-Yvette, {eric.goubault,sylvie.putot}@cea.fr

<sup>2</sup> Hispano-Suiza, philippe.baufreton@hispano-suiza-sa.com

<sup>3</sup> IRSN, jean.gassino@irsn.fr

**Abstract.** Finite precision computations can severely affect the accuracy of computed solutions. We present a complete survey of a static analysis based on abstract interpretation, and a prototype implementing this analysis for C codes, for studying the propagation of rounding errors occurring at every intermediary step in floating-point computations. In the first part of this paper we briefly present all the domains and techniques used in the implemented analyzer, called FLUCTUAT. We describe in the second part, the experiments made on real industrial codes, at Institut de Radioprotection et de Sûreté Nucléaire and at Hispano-Suiza, respectively coming from the nuclear industry and from aeronautics industry. This paper aims at filling in the gaps between some theoretical aspects of the static analysis of floating-point computations that have been described in [13, 14, 21], and the necessary choices of algorithms and implementation, in accordance with practical motivations drawn from real industrial cases.

**Keywords.** Static analysis, floating-point computations, control systems

## 1 Introduction

The use of floating-point arithmetic as a computable approximation of real arithmetic may introduce rounding errors at each arithmetic operation in a computation. Even though each of these errors is very small, their propagation in further computations, for example in loops, can produce a dramatic imprecision on a critical output. We propose an analysis that computes for each variable an over-approximation of the difference between a same computation in real and in floating-point arithmetic. Moreover, we decompose this difference over the operations that introduced errors, thus pointing out the operations responsible for the main losses.

*Principal contributions* This paper essentially surveys the different abstract interpretation domains and techniques used in the static analyzer FLUCTUAT. The emphasis is not put on the theoretical details of the main domain used for abstracting floating-point values, which can be found in [14]. But it describes the main techniques and answers to the difficulties we had to address in the practical use of these domains in an analyzer. Also, an important part is dedicated to some examples of analyzes which were led on industrial examples.

*Related work* Few complete static analyzers of C programs are fully described in the literature. This paper tries to match some of the available descriptions of the ASTREE analyzer [1, 2] which is probably the most complete analyzer today, with respect to the number of techniques and domains implemented, in particular. Also, some commercial abstract interpreters like CodeSonar [24] and PolySpace [25] are now available. All these tools analyze run-time errors mostly, whereas we analyze a subtle numerical property, namely the discrepancy introduced by the use of floating-point numbers instead of real numbers, in C programs. This requires very fine and specific abstract domains, and the difficulty lies mostly in the numerical subtleties of small parts of a code, and not in the size of the program.

*Overview* This paper is divided in two parts. In Section 2, we present the abstract domains and main techniques used in the FLUCTUAT analyzer. In subsection 2.1, we briefly introduce the abstract domain for representing floating point variables, which is described in more details in [14]. Then in subsection 2.2, we detail how this domain is extended to integer variables, and the specificities and difficulties of handling integers. We then describe in subsection 2.3, the simple aliasing model we are using, when it comes to abstract pointers, structures and arrays. The iteration strategy which is used to solve the semantic equations (and in particular specific widening operators), is described in subsection 2.4. Finally, some assertions in a language specific to the analyzer, allow to specify properties of variables, such as set of possible input values, but also more subtle properties such as bounds on the gradient of an input over iterations in a loop. These are presented in subsection 2.5.

We then discuss in Section 3 some experiments conducted with FLUCTUAT on industrial codes. We first describe the analysis of some programs developed at Hispano-Suiza in the aeronautics industry. We first concentrate on some interesting specific sub-functions, and then come to a full control application. In a second part, we describe the analysis of a code from the nuclear industry that IRSN has to expertise.

## 2 Abstract domains and techniques used in FLUCTUAT

### 2.1 Floating-point variables

**General principles** The analysis bounds at each operation the error committed between the floating-point and the real result. It relies for this on a model of the difference between the result  $x$  of a computation in real numbers, and the result  $f^x$  of the same computation using floating-point numbers, expressed as a sum of error terms

$$x = f^x + \bigoplus_{\ell \in L \cup \{hi\}} \omega_{\ell}^x. \quad (1)$$

We assume that a control point of a program is annotated by a unique label  $\ell \in L$ . In this relation, a term  $\omega_{\ell}^x$ ,  $\ell \in L$ , denotes the contribution to the global

error of the first-order error introduced by the operation labeled  $\ell$ . It expresses both the rounding error committed at label  $\ell$ , and its propagation due to further computations on variable  $x$ . Errors of order higher than one, coming from non-affine operations, are grouped in one term associated to special label  $hi$ . We let  $\mathcal{L}$  be the union of  $L$  and  $hi$ .

Let  $\mathbb{F}$  be either the set of simple or double precision floating-point numbers. Let  $\uparrow_{\circ}: \mathbb{R} \rightarrow \mathbb{F}$  be the function that returns the rounded value of a real number  $r$ , with respect to the rounding mode  $\circ$ . The function  $\downarrow_{\circ}: \mathbb{R} \rightarrow \mathbb{F}$  that returns the roundoff error is then defined for all  $f \in \mathbb{R}$ , by  $\downarrow_{\circ}(f) = f - \uparrow_{\circ}(f)$ . The result of an arithmetic operation  $\diamond^{\ell_i}$  contains the combination of existing errors on the operands, plus a new roundoff error term  $\downarrow_{\circ}(f^x \diamond f^y) \varepsilon_{\ell_i}$ . For addition and subtraction, the existing errors are added or subtracted componentwise, and a new error  $\downarrow_{\circ}(f^x + f^y)$  is associated to control point  $\ell_i$ :

$$x +^{\ell_i} y = \uparrow_{\circ}(f^x + f^y) + \bigoplus_{\ell \in \mathcal{L}} (\omega_{\ell}^x + \omega_{\ell}^y) + \downarrow_{\circ}(f^x + f^y).$$

The multiplication introduces higher order errors, we write :

$$x \times^{\ell_i} y = \uparrow_{\circ}(f^x f^y) + \bigoplus_{\ell \in \mathcal{L}} (f^x \omega_{\ell}^y + f^y \omega_{\ell}^x) + \sum_{(\ell_1, \ell_2) \in \mathcal{L}^2} \omega_{\ell_1}^x \omega_{\ell_2}^y + \downarrow_{\circ}(f^x f^y).$$

We refer the reader to [10, 19] for more details on this domain.

**Unstable tests** Our approach is that of abstract interpretation [7], and all control flows due to sets of possible inputs are considered. But there is one specific difficulty due to floating-point computations. Indeed, in tests, the branch followed by the floating-point and the corresponding real value of a variable can be different, we then call them unstable tests (as in [26]). Consider for example the following portion of code, supposing input  $x$  is in interval  $[1,3]$  with an error equal to  $1.0e-5$ :

```
if (x <= 2) x = x+2;
```

Then, for  $x$  equal to 2 for example, the floating-point result after this test is 4, whereas the result of this program if it were executed on the real semantics would be 2.00001. But handling this divergence in control flow in the general case would be complicated and costly, and quickly very imprecise. For example here, if we consider the different control flows, we find the floating-point value of  $x$  in  $[2,4]$ , with an error in interval  $[1.0e-5, 2]$ . Without any additional relation between values and errors, this result is highly imprecise.

We thus made the choice in the Fluctuat analyzer to make the assumption that the real and floating-point flows take the same branches. The result given here at the end of the program would thus be  $x = [2,4]$  with an error equal to  $1.0e-5$  (if we neglect the additional rounding error due to the addition).

However, when the analyzer detects, as is the case here, that the control flows may be different, it issues a warning.

**Relational domain** A natural abstraction of the coefficients in expression (1), is obtained using intervals. The machine number  $f^x$  is abstracted by an interval of floating-point numbers, each bound rounded to the nearest value in the type of variable  $x$ . The error terms  $\omega_i^x \in \mathbb{R}$  are abstracted by intervals of higher-precision numbers, with outward rounding. However, results with this abstraction suffer from the over-estimation problem of interval methods. If the arguments of an operation are correlated, the interval computed with interval arithmetic may be significantly wider than the actual range of the result.

We thus proposed and implemented a relational domain, relying on affine arithmetic [5, 22] for the computation of the floating-point value  $f^x$ . Affine arithmetic uses affine correlation between real variables, and allows to get much tighter results than classical interval arithmetic (the concretisation forms zonotopes : center-symmetric bounded convex polytopes). It relies on a representation of a quantity  $x$  by an affine form, which is a polynomial of degree one in a set of noise terms  $\varepsilon_i$  :

$$\hat{x} = \alpha_0^x + \alpha_1^x \varepsilon_1 + \dots + \alpha_n^x \varepsilon_n, \quad \text{with } \varepsilon_i \in [-1, 1] \text{ and } \alpha_i^x \in \mathbb{R}. \quad (2)$$

Each noise symbol  $\varepsilon_i$  stands for an independent component of the total uncertainty on the quantity  $x$ , its value is unknown but bounded in  $[-1, 1]$ ; the corresponding coefficient  $\alpha_i^x$  is a known real value, which gives the magnitude of that component. The sharing of noise symbols between variables expresses implicit dependencies. The full semantics is described in [14], and linearizes floating-point expressions dynamically (and not statically as in [20]). The semantics is memory-efficient: it needs only a small factor of the size that an (economic) interval analysis would take. No a priori decided packing of variables [2] is needed since the representation of relations is implicit [14]. Nevertheless, we use a sparse representation of the global environment, akin to the one described in Section 5 of [1].

The coefficients  $\alpha_i^x$  have no meaning relevant to our analysis, the decomposition is a mean for a more accurate computation. This is different from expression (1), where coefficient  $\omega_\ell^x$  represent the contribution of control point  $\ell$  to the total rounding error. However, they can be used for an analysis of the sensibility of a program to an input : when an input is taken in a small interval, a new noise symbol  $\varepsilon_i$  is created. The evolution of the corresponding  $\alpha_i^x$  in further computations indicates how this initial uncertainty is amplified or reduced [23].

These affine forms allow to represent results of real arithmetic. The analysis must be adapted to the case of floating-point arithmetic, where symbolic relations true in real arithmetic do no longer hold exactly. We thus decompose the floating-point value  $f^x$  of a variable  $x$  resulting from a trace of operations, in the real value of this trace of operations  $r^x$ , plus the sum of errors  $\delta^x$  accumulated along the computation,  $f^x = r^x + \delta^x$ . The real part is computed using affine arithmetic, and the error is computed using three intervals that respectively bound the error on the lower and upper bounds of the set of real values  $r^x$ , and the maximum error on all this set. Without going into too much detail, we can say that these errors on bounds allow to improve the estimates for the floating-

point bounds, compared to using the maximum error. But the maximum error can still be needed at each step to estimate the results of further computations.

This domain for the values of variables, is of course more expensive than interval arithmetic, but comparable to the domain used for the errors. And it allows us to accurately analyze non trivial numerical computations, as we will show in Section 3. We also plan to introduce a relational computation for errors, along the lines of [15]. First ideas on these relational semantics were proposed in [13], [21]. The relational semantics for the value  $f^x$  is described in detail in [14], with in particular the lattice operations such as join and meet.

## 2.2 Integer variables

We implemented modular integer arithmetic semantics, and a domain consisting of value (coded by affine forms) plus sum of errors is used as for floating-point variables. For example, when adding one to the greatest integer that can be represented in the `int` type, say `INT_MAX`, the value of the result is the smallest integer represented by an `int`, say `INT_MIN`, and an error of `INT_MAX-INT_MIN+1` is associated to this variable. Conversions between integers and floating-points are supported, and the errors are propagated.

**Bitwise operations** Some attention must be paid to the propagation of errors on operands in order to avoid losing too much precision. Indeed, the behavior of the `and`, `or` and `xor` operators is non affine with respect to the operands. In the general case, the errors on the operands  $x$  and  $y$  are propagated as follows:

- we compute the result of  $\diamond$  on the sets of floating-point values,  $f^z = f^x \diamond f^y$ ,
- we compute the result of the same operation on the interval bounds for the real values,  $r^z = r^x \diamond r^y$ , with  $r^x = f^x + \sum_l \omega_l^x$  and  $r^y = f^y + \sum_l \omega_l^y$ ,
- then the propagated error on  $z$  is  $r^z - f^z$ , and it is associated to the label of the current operation.

There are two consequences. First, we lose the decomposition of errors on operations executed before bitwise operations. Second, the larger the intervals  $f^x$  and  $f^y$ , the more over-approximated the propagated errors are. We thus propose an option of the analyzer to locally subdivide one of these intervals in the propagation of errors : the cost of a bitwise operation is approximately multiplied by the number of subdivisions, but this cost is in general negligible compared to the full analysis, and the results can be greatly improved.

Error terms are agglomerated for the same reason for the division and modulo operators on integers. The error is also computed as the difference between the floating-point interval result and the real interval result : local subdivisions can again greatly improve the estimation of errors.

**Conversions** The semantics for the conversion between floating-point numbers, and with integers differ by the meaning we give to each :

- in the conversion from double precision to simple precision floating-point numbers, we consider the difference between the initial double precision value and the result of the conversion, as an error on the result.
- in the conversion from a floating-point number to an integer, we consider that the truncation is wanted by the user, and is thus not an error. A new error can be introduced by such a conversion only when the floating-point number exceeds the capacity of the integer type. However, all errors are grouped in one integer term corresponding to the label of the conversion.
- in the conversion from an integer to a floating-point number, most of the time no precision is lost, and the sum value plus errors is transmitted as is. However, this is not always the case, and an error still has to be added in some cases : for example a 32 bits integer with all bits equal to 1 cannot exactly be represented by a simple precision floating-point number, which mantissa is represented on 23 bits.

We encountered some other cast operations that we included in the set of instructions understood by the analyzer, such as the ones used to decode and encode IEEE 754 format, directly by bitwise operations. Take for instance the following piece of code (assuming 64 bits little endian encoding for `double`):

```
double Var = ...; signed int *PtrVar = (signed int *) (&Var);
int Exp = (signed int) ((PtrVar[0] & 0x7FF00000) >> 20) - 1023;
```

We cast variable `Var` into an array of 32 bits types. Then we extract the first 32 bits of the 64 bits word. The rest of the manipulation of the program above, masks the bits of the mantissa, and shifts the value, to get in `Exp` the binary exponent in IEEE754 format of the value stored in `Var`.

In the interpretation of this case by FLUCTUAT, all error terms are agglomerated in one corresponding to the label of the cast, and local subdivisions of the values can be applied to improve the bounds for the errors, as for bitwise operations.

### 2.3 Aliases and arrays

Our alias and array analysis is based on a simple points-to analysis, like the ones of [18], or location-based alias analyses. An abstract element is a graph of locations, where arcs represent the points-to relations, with a label (which we call a selector) indicating which dereferencing operation can be used. Arrays are interpreted in two different ways, as already suggested by some of the authors in [8]: all entries are considered to be separate variables (called “expanded” in [1]) or the same one (for which the union of all possible values is taken - called “smash” in [1]). These abstractions have proven sufficient typically for SCADE generated C programs.

### 2.4 Iteration strategy

**Loops** The difficulty in loops is to get a good over-approximation of the least fixpoint without too many iterations of the analyzer. For that, we had to design adapted iteration strategies :

- in the case of nested loops, a depth first strategy was chosen : at each iteration of the outer loop, the fixpoint of the inner loop is computed,
- a loop is unfolded a number of times (similar to the “semantic loop unrolling” of [1]), before starting Kleene iterations (unions over iterates),
- some particularities of our domain require special care in loops : for example, noise symbols are potentially introduced at most operations, and there are new noise symbols for each iteration of a loop. But we can choose to reduce the level of correlation we want to keep, and for example keep correlations only between the last  $n$  iterations of a loop, where  $n$  is a parameter of the analyzer. Also, we can choose to agglomerate or not some noise symbols introduced in a loop when getting out of it. This allows to reduce the cost of the analysis while keeping accurate results.
- acceleration techniques (widening) adapted to our domain had to be designed. In particular, widenings are not always performed at the same time on integer or floating-point variables, and on values or error terms. Also, we have designed a widening specially adapted to floating-point numbers, by gradually reducing the precision of the numbers used to bound the terms : this accelerates the convergence of Kleene iteration compared to iteration with fixed precision, and allows to get very accurate results. This should be thought of as an improved method than the “staged widening with thresholds” of [1], in the sense that thresholds are dynamically chosen along the iterations, depending on the current values of the iterates. After a number of these iterations, a standard widening is used.

To illustrate this last point (progressive widening by reduction of the precision), let us consider the fixpoint computation of

```
while () x = 0.1*x;
```

with no unrolling of the loop, starting from  $x_0 \in [0, 1]$ . With our simple (non relational) semantics, we have, with  $ulp(1)$  denoting the machine rounding error around 1,

$$\begin{aligned} x_1 &= [0, 1] + \delta\varepsilon_2, \quad \delta = 0.1[-ulp(1), ulp(1)] \\ x_2 &= [0, 1] + (0.1\delta + \delta)\varepsilon_2 \\ x_n &= [0, 1] + \left(\sum_{k=0}^{n-1} 0.1^k\right)\delta\varepsilon_2. \end{aligned}$$

If real numbers are used to compute the error term, without any widening the computation does not terminate even though the error term remains finite. Now if floating point numbers are used to bound the error term, the convergence depends on the number of bits used to represent the mantissa. For simplicity's sake, let us consider  $\delta = [-1, 1]$ , and radix 10 numbers. With 3 significant digits, a fixpoint is got in 4 iterates :

$$\begin{aligned} \omega_1 &= \delta = [-1, 1] \\ \omega_2 &= \uparrow_{\infty} ([-0.1, 0.1] + [-1, 1]) = [-1.1, 1.1] \end{aligned}$$

$$\begin{aligned}\omega_3 &= \uparrow_\infty ([-0.11, 0.11] + [-1, 1]) = [-1.11, 1.11] \\ \omega_4 &= \uparrow_\infty ([-0.111, 0.111] + [-1, 1]) = \uparrow_\infty [-1.111, 1.111] = [-1.12, 1.12] \\ \omega_5 &= \uparrow_\infty ([-0.112, 0.112] + [-1, 1]) = \uparrow_\infty [-1.112, 1.112] = [-1.12, 1.12]\end{aligned}$$

More generally, we can show that with  $N$  significant digits, a fixpoint is got in  $N+1$  iterates. Thus reducing the precision of numbers accelerates the convergence towards a (larger) fixpoint.

Of course, this is a toy example, in practice the fixpoint is computed by unrolling the loop a certain number of times before beginning the unions, which here solves the problem. But we are confronted with this kind of computations in the general case. And in more complicated examples, when the optimal unrolling was not chosen, this widening allows to still compute an interesting fixpoint.

**Interprocedural analysis** In critical embedded systems, recursion is in general prohibited. Hence we chose to use a very simple interprocedural domain, with static partitioning, based on [17].

## 2.5 Assertions

A number of assertions can be added to the analyzed C source code, to specify the behavior of some variables of the program. For a simple precision floating-point variable  $x$ , the assertion :

```
x = __BUILTIN_DAED_FBETWEEN(a,b);
```

indicates that  $x$  can take any floating-point value in the interval  $[a, b]$ . The same assertions exist to define the range of double precision or integer variables.

One can also specify an initial error together with the range of values for a variable. For example,

```
x = __BUILTIN_DAED_FLOAT_WITH_ERROR(a,b,c,d);
```

specifies that variable  $x$  of type `float` takes its value in the interval  $[a, b]$ , and that it has an initial error in the interval  $[c, d]$ .

In some cases, bounds on the values are not sufficient to describe accurately the behavior of a system : we thus added an assertion that allows to bound, in a loop indexed by an integer variable  $i$ , the variation between two successive values of an input variable  $x$  :

```
for (i=i0 ; i<N ; i++)
  x = __BUILTIN_DAED_FGRADIENT(x0min,x0max,gmin(i),gmax(i),xmin,xmax,i,i0);
```

In this assertion,  $i_0$  is the value of variable  $i$  at first iteration. The value of  $x$  at first iteration is in interval  $[x0min,x0max]$ , the difference between two successive iterates is in the interval  $[gmin(i),gmax(i)]$ , which bounds may depend on the iterate, and the value of  $x$  is always bounded by  $[xmin,xmax]$ . Thus  $x(i_0) \in [x0min,x0max]$ , and for all  $i \in \{i_0, \dots, N\}$ , we have  $x(i) \in (x(i-1) + [gmin(i),gmax(i)]) \cap [xmin,xmax]$ . Our relational domain (subsection 2.1) is specially well adapted to dealing with these relations between iterates in a loop. An example of the use of this assertion is given in the worst-case scenario part of example presented in subsection 3.2.



**Subdivisions of inputs** In the example `SqrtR` of subsection 3.1, even with the relational domain, non-linearities of the studied iterative scheme produce too much imprecision, and the solver of the abstract equations does not prove the termination of the analyzed algorithm. A solution to this is to restrict the range of values of the inputs, for which we want to analyze the program, so that we are close enough to linear behaviors. This is done in `FLUCTUAT` by subdividing the domain of some inputs whose ranges are already bounded by assertions of the type `_BUILTIN_DAED_FBETWEEN`. The user can select one or two such variables to be subdivided by pointing in the program the corresponding assertions.

`FLUCTUAT` analyzes independently the program as many times as we subdivide some of the inputs. Suppose we subdivide  $n$  times an input variable  $x$  which has range, defined by an assertion, in  $[a, b]$ : the analyzer will analyze the program with  $x$  in  $[a, a + \frac{b-a}{n}]$ , then in  $[a + \frac{b-a}{n}, a + 2\frac{b-a}{n}]$ ,  $\dots$ ,  $[a + (n-1)\frac{b-a}{n}, b]$ . Hence it does not need more memory than needed for one analysis, but takes about  $n$  times the duration of one analysis, where  $n$  is the number of subdivisions. In the case when we subdivide two such assertions, the subdivisions are completely independent, hence leading to a quadratic factor time increase of the analysis. We chose not to offer the user the possibility to subdivide the values of more than two input variables because it would lead to too slow analyses. This would become reasonable only for parallel versions of `FLUCTUAT`.

This kind of subdivision cannot be used for an assertion in a loop, because it would be equivalent to choosing at all iterates of the loop the values of  $x$  to be in the same sub-interval. Indeed, subdividing independently all iterates would be far too costly, and maybe not either what is really intended by the user. We thus proposed the special assertion `_BUILTIN_DAED_FGRADIENT` for these cases of reactive programs, where inputs are acquired cyclically over time.

### 3 Experiments on control systems

#### 3.1 Hispano-Suiza

The Full Authority Digital Engine Control, better known as a FADEC, is one of the largest electronic control units on an aircraft. The FADEC continuously processes and analyzes key engine parameters (up to 70 times a second), to make sure the engine operates at maximum potential. It manages the startup phase (which takes only about 40 seconds on the latest models), and then the entire operating envelope, from idle to full throttle. The following test cases for `FLUCTUAT` are extracted from pieces of code which have been written during the development of reusable libraries, designed for the FADEC. They are representative of the code of the FADEC, and some of them present some hard numerical difficulties for static analyzers.

**Experiments on elementary symbols** We examined several elementary symbols used in applications at Hispano-Suiza. Elementary symbols are manually developed and coded independently from `SCADE` which is used as a design tool.

Among these symbols was the following code (slightly changed for convenience), which is intended to return `Output` equal to the square root of `Input` by a Householder method. This involves the iteration, until some residue is less than a small value `_EPS=e-6`, of a fifth-order polynomial. *The number of iterates for the algorithm to converge is thus not given by the syntax of the program and must be the result of an accurate analysis.*

```
void SqrtR (double Input, double *Output)
{
  double xn, xnp1, residu, lsup, linf;
  int i, cond;
  Input = __BUILTIN_DAED_DOUBLE_WITH_ERROR(0.1,20.0,0,0);
  if (Input <= 1.0) xn = 1.0; else xn = 1.0/Input;
  xnp1 = xn;   residu = 2.0*_EPS*(xn+xnp1);
  lsup = _EPS * (xn+xnp1);   linf = -lsup;
  cond = ((residu > lsup) || (residu < linf)) ;
  i = 0;
  while (cond)
  { xnp1 = xn * (15S8 + Input*xn*xn*(-5S4+3S8*Input*xn*xn));
    residu = 2.0*(xnp1-xn);   xn = xnp1;
    lsup = _EPS * (xn+xnp1);   linf = -lsup;
    cond = ((residu > lsup) || (residu < linf)) ; i++; }
  *Output = 1.0 / xnp1; }
```

When `Input` is in  $[0.1,20]$  as above, FLUCTUAT with 100000 subdivisions converges to a finite and precise estimate of the floating-point value of `Output` and of the number of iterates of the studied algorithm. Indeed, it finds `Output` to be in  $[3.16e-1, 4.48]$  with global error in  $[-2.56e-13, 2.56e-13]$ . And the number of iterates in the main loop `i` is found to be within 1 and 6. This number of iterates is an exact result, as can be confirmed by using FLUCTUAT in the symbolic execution mode for `Input` equal to 1 and to 20 respectively (or alternatively, by checking the execution of the real binary file). Note that no other static analyzer we know of is able to find even a good approximation of the floating-point enclosure of `Output`.

When perturbing the input by a small error, by the assertion

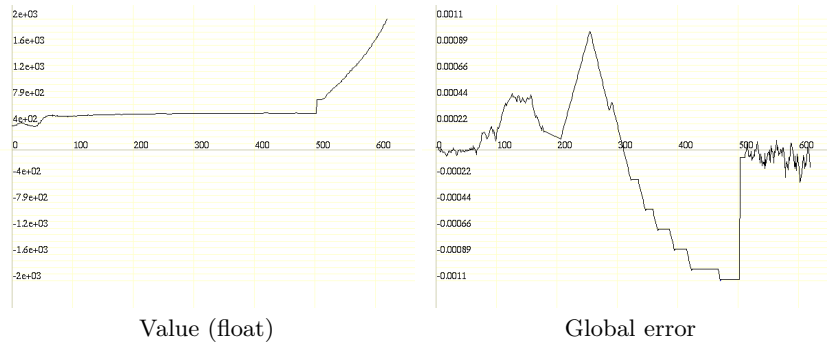
```
__BUILTIN_DAED_FLOAT_WITH_ERROR(0.1,20.0,-0.00000001,0.00000001)
```

and still subdividing 100000 times, we find in 603 seconds and 4Mb of memory on a 1 Ghz laptop pc, the same floating-point enclosure, and a global error in  $[-3.06e-06,3.06e-06]$  which is mainly due to the initial error of the order of  $e-8$ . This shows the good behavior of the algorithm. Even though the results on the global error seem satisfying, we hope to be able to improve them a lot with a new relational domain on the error terms, as sketched in [13].

In fact, FLUCTUAT does not need to subdivide equally for all ranges of the input. For instance, with `Input` restricted to  $[16,20]$ , it needs only 133 subdivisions to converge. Whereas, with `Input` restricted to  $[0.1,1]$ , it needs about 4500 subdivisions. Hence a dynamic subdivision mode is planned for a future version.

These results indicate a good behavior of the algorithm for `Input` in the range  $[0.1,20]$ . However, this function was designed to be used for `Input` in  $[1e-50,1e50]$ . Symbolic execution shows that the algorithm is much less satisfying for this extended range, and may need up to 95 iterations, which is too important for practical use, because of timing constraints. Since then, the algorithm for the square root has been changed.

**Representative code** The following test case for FLUCTUAT is extracted from pieces of SCADE code which have been generated during the development of a military engine controller. The control law named `asservxn2` is aimed to control the speed of the Low-Pressure Compressor called a fan. Therefore, the control loop should be stable inside the whole flight domain, including the fuel flow `wf32cb` and the motor regime `xn2`.



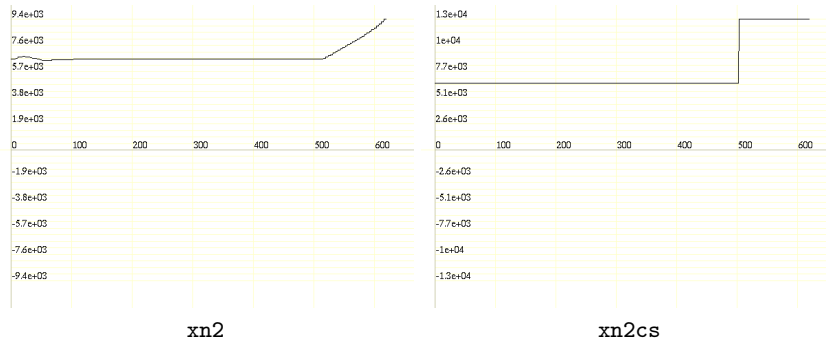
**Fig. 1.** Evolution of the fuel flow `wf32cb` for a given target motor regime

The code is 2358 lines long in C (44 functions, among which filters, interpolators, integrators etc.), uses complex nested compound structures containing arrays and pointers, and has 167 integer variables and 269 floating-point variables known at the end of the main function (and many more local variables).

We have first used the static analyzer as an abstract interpreter (hence delivering information about precision loss) on the test scenarios that have been used to certify the program. As an example, the first test scenario consists in showing that the rotation regime `xn2` is well controlled by the command on the flow `wf32cb`. FLUCTUAT has been run in symbolic execution mode (i.e. with the semantics described in Section 2.1, but on one control flow only) on the sequence of 2500 consecutive inputs, on a 50 seconds duration. The scenario corresponds to a target low pressure regime shown in Figure 2 for the first 650 inputs. The control program computes the fuel flow necessary to reach this regime, see Figure 1. As shown in the excerpt of the test scenario, the motor regime is well controlled by `wf32cb`: when `wf32cb` increases, `xn2` increases as

well until it reaches its target value `xn2cs`, in which case `wf32cb` stabilizes. At iteration 500, the target regime is increased and the control begins. The error in the command `wf32cb` is shown to be always bounded by  $10^{-3}$  in absolute value, which indicates a good (relative) precision of the control algorithm.

Other similar tests have been carried out. We are in the process of studying the code for more general inputs (i.e. infinite sets of inputs, corresponding to ranges of target motor regimes), in a similar manner as done in next section (using gradient constraints on the inputs). It is to be noted though that the control mechanism uses an integrator, which is known to be hard to analyze, see for instance [12] for more explanations.



**Fig. 2.** Evolution of the motor regime `xn2`, and of its target `xn2cs`

### 3.2 Institut for Radiological protection and Nuclear Safety (IRSN)

Computer systems are increasingly used for safety purposes in research and industrial nuclear reactors. For example, on the latest French power reactor series, representing 24 units, software is used to perform safety functions including critical ones like protection. The Protection System does not control the process but monitors it, by continuously acquiring parameters like water pressure, temperatures in different pipes, neutron flux in different locations, positions of control rods, and so on. From these inputs, the system computes tens of values using classical data processing techniques: filters, arithmetic and logic operations, non-linear functions, thresholds and so on. The system then checks that these computed values remain within the authorized domain. If not, it has to automatically shutdown the reactor within half a second, and to trigger safety systems like water injection or spraying, depending on the nature of the incident or accident.

The following test case (mean-square filter), was taken from a representative piece of code that IRSN has to give expertise on, for the French nuclear certification body.

Several key process parameters are sampled every 50 milliseconds by the protection system, which stops the reactor if a given threshold is exceeded. Unfortunately, the readings are affected by noise, which could delay a necessary stop or, on the contrary, cause a spurious one. The actual power of the plant below the threshold to provide a noise margin is not adequate for obvious economic reasons. On the other hand a spurious stop is also undesirable because it induces strong promptings to the mechanical structures and prevents the plant to produce electricity for several hours. Thus, a least-square linear regression filter is applied to improve the estimate of the most sensitive parameter. It is important however to make sure that this filtering step does not add too much numerical error due to rounding, we thus used FLUCTUAT to bound the error committed in it, and study the propagation of existing errors.

The filter is adaptative, that means its depth  $D_k$  can vary at each cycle  $k$ , according to a formula that depends on parameters which are known exactly, but can take different values, depending on the signal value. The input sample  $(e_k)_k$  is in an interval  $[1e^2, 1.5e^8]$ , and is transformed to give the input of the filter, by  $Y_k = \log(ae_k + b)$ . In this transformation, only ranges are known for parameters  $a$  and  $b$ , with nominal values.

**Worst-case scenario.** We first consider a reduced version of the filter, using the fact that the filter can be written in such a way that outputs are independent, except that the depth of the filter depends on the previous values. It can be shown that, with the parameters used, the depth of the filter is always bounded. We thus study a worst-case scenario, that allows to get bounds for the outputs and errors on the outputs that hold true for any step of the filter. For that, we take a bounded filter depth, and the inputs in the maximum possible range, that is  $Y_k \in [8.42, 22.4]$ .

- We first suppose all inputs are independent and can take any value in this range at any step, by

```
for (int k=1 ; k<=N ; k++) {
    Yk = __BUILTIN_DAED_FBETWEEN(8.42,22.4); ... }
```

Then we get with the relational domain, the following enclosures of the filtered value  $O$ , and of a value  $S$  related to the variation speed :

$$O = [3.912, 26.907] + [-1.91e - 4, 1.91e - 4]\varepsilon$$

$$S = [-0.350, 0.350] + [-5.37e - 6, 5.37e - 6]\varepsilon$$

In these two expressions, the first interval bounds the floating-point value, the second one bounds the rounding error, the filter being implemented using simple precision floating-point numbers.

- In order to have a more representative model of the inputs, we then used the assertion on the gradient to limit the variation between two successive inputs : we still take the range of inputs equal to  $[8.42, 22.4]$ , but also bound the difference between two successive inputs by  $[0, 0.01]$ , by

```
for (int k=1 ; k<=N ; k++) {
  Yk = __BUILTIN_DAED_FGRADIENT(8.42,22.4,0,0.01,8.42,22.4,k,1);...}
```

Then we get, with the relational domain, much tighter enclosures :

$$O = [8.33, 22.5] + [-1.91e - 4, 1.91e - 4]\varepsilon$$

$$S = [-1.18e - 5, 0.01] + [-5.34e - 6, 5.34e - 6]\varepsilon$$

The error is of the same order as previously, but now the bounds for the value of the output are very close to the input bounds, and we get back the information on the variation speed.

**Complete filter.** We now want to study more closely the behavior of the output. We choose here a plausible scenario for the inputs, that is a sampling of function  $e(x)$  defined by

$$x \leq 0 : e(x) = 1.e^2,$$

$$x > 0 \text{ et } e(x) \leq 1.5e^8 : e(x) = e(0) * 2^{50*x/60}.$$

We take intervals for the coefficients of the transformation, and add a perturbation to the input of the filter thus obtained (note that we could also have taken small intervals around these inputs). We present in Figures 3 and 4, the results got with FLUCTUAT, for the evolution over time of the bounds on the values and errors on the input  $Y_k$  of the filter, and of its output.

The error on the input is due partly to the logarithm computation, partly to the addition of a perturbation depending on previous inputs. For the time being, we have parameterized the error due to the logarithm computation, which is not yet specified in the IEEE 754 norm.

The output is approximately in the same range as the input and looks indeed smoothed. We represent in Figure 4 right, the evolution of the depth of the filter, it must be noted that the range at a given time depends on the values of the parameters, it does not have one fixed value. The error on the output is overestimated (relational analysis for errors not implemented), but we can still observe that the error is not too much amplified. We can also note that the variation of the error on the output is related to the variation of the filter depth.

Finally, we can note that the magnitude of the maximum error on the output, is of the same order as the magnitude on the output by the worst-case analysis. This confirms the relevance of the worst-case analysis.

## 4 Conclusion and Future Work

We have shown in this paper how we designed a static analyzer for bounding the imprecision error in numerical programs. This design relied on a careful study of the semantics of the IEEE754 standard, of numerical convergence of the kind of iterative schemes we encountered in control systems, and of specificities of the

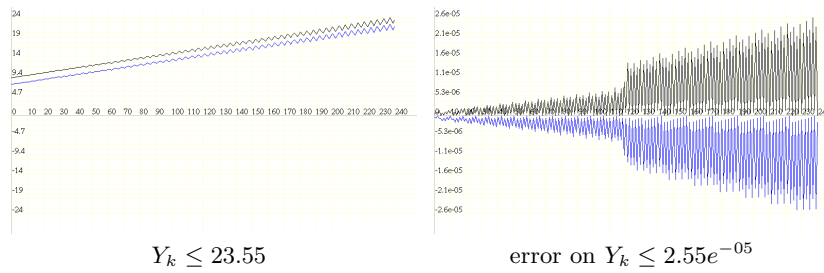


Fig. 3. Evolution of the filter input over time

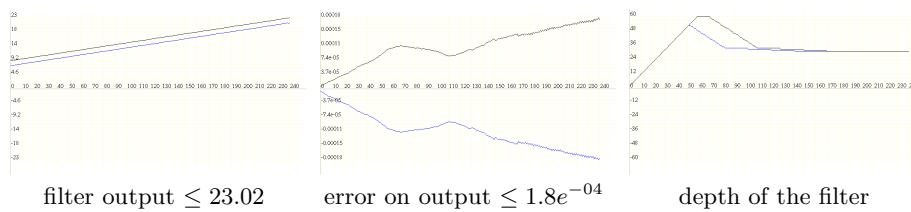


Fig. 4. Evolution of the filter output over time

programming of these systems (SCADE generated code in general). Some real industrial examples were given, which were sufficiently simple to explain in a few pages. The analyzer has already been used (on a low-end PC with 512Mb of memory) for code of the order of 10 thousand lines of C code, and seems to scale up well.

In the future, we plan to invest more on domains dealing precisely with integrators, which make precision analysis hard to carry out, see [12] for example. We also plan to experiment relational domains for error computations as briefly sketched in [13, 15]. We also would like to improve the precision of the least fixed point computation of the abstract equations in our tool, using policy iteration mechanisms, see [6, 9]. Last but not least, our abstract domain is well adapted to under-approximations as well [16]. Our intention is to derive also this kind of information, so as to assess the quality of the results we give, and give some indications in some cases that the control system under analysis is definitely not implemented in a sufficiently accurate way.

## References

1. B. Blanchet, P. Cousot, R. Cousot, J. Feret, L. Mauborgne, A. Miné, D. Monniaux and X. Rival. Design and Implementation of a Special-Purpose Static Program Analyzer for Safety-Critical Real-Time Embedded Software. *Essays Dedicated to Neil D. Jones*. LNCS 2566, Springer-Verlag, Berlin, 2002.
2. B. Blanchet, P. Cousot, R. Cousot, J. Feret, L. Mauborgne, A. Miné, D. Monniaux, and X. Rival. A Static Analyzer for Large Safety-Critical Software. *Proc. PLDI'03*.

3. F. Bourdoncle. Abstract Interpretation by dynamic partitioning. In *Journal of Functional Programming*, 2(4):407-435, 1992.
4. F. Bourdoncle. Efficient chaotic iteration strategies with widenings. LNCS 735.
5. J. L. D. Comba and J. Stolfi. Affine arithmetic and its applications to computer graphics. In *SIBGRAPI'93*.
6. A. Costan and S. Gaubert and E. Goubault and M. Martel and S. Putot. A policy iteration algorithm for computing fixed points in static analysis of programs. In *Proceeding of CAV'05*, LNCS 3576, 2005.
7. P. Cousot and R. Cousot. Abstract interpretation frameworks. In *Journal of Logic and Symbolic Computation*, 2(4):511-547, 1992.
8. E. Goubault, D. Guilbaud, A. Pacalet, B. Starynkévitch and F. Védryne. A Simple Abstract Interpreter for Threat Detection and Test Case Generation, WAPATV'01.
9. S. Gaubert, E. Goubault, A. Taly and S. Zennou. Static Analysis by Policy Interaction on Relational Domains In *Proceedings of ESOP'07*.
10. E. Goubault. Static analyses of the precision of floating-point operations. In *Proceedings of SAS'01*, LNCS 2126.
11. E. Goubault, M. Martel, and S. Putot. Asserting the precision of floating-point computations : a simple abstract interpreter. In *ESOP'02*, LNCS 2305.
12. E. Goubault and M. Martel and S. Putot. Some future challenges in the validation of control systems. In *Proceedings of ERTS'06*.
13. E. Goubault and S. Putot. Weakly Relational Domains for Floating-Point Computation Analysis. *Presented at NSAD'05*.
14. E. Goubault and S. Putot. Static Analysis of Numerical Algorithms. In *Proceedings of SAS'06*, LNCS 4134, 2006.
15. E. Goubault and S. Putot. Automatic analysis of imprecision errors in software, <http://www.di.ens.fr/~goubault/papers/abstract.pdf>, 2007.
16. E. Goubault and S. Putot. Under-approximations of computations in real numbers based on generalized affine arithmetic. Submitted <http://www.enseignement.polytechnique.fr/profs/informatique/Sylvie.-Putot/Publications/under.pdf>, 2007.
17. N. D. Jones and S. S. Muchnick. A flexible approach to interprocedural flow analysis and programs with recursive data structures. In *POPL'82*.
18. W. Landi and B. Ryder. A safe approximate algorithm for inter-procedural pointer aliasing In *Proceedings of PLDI*, ACM, 1992.
19. M. Martel. Propagation of roundoff errors in finite precision computations : a semantics approach. In *Proceedings of ESOP'02*, LNCS 2305.
20. A. Miné. Relational Abstract Domains for the Detection of Floating-Point Run-Time Errors. In *Proceedings of ESOP'04*, LNCS 2986.
21. S. Putot, E. Goubault and M. Martel. Static Analysis-Based Validation of Floating-Point Computations. In *Proceedings of Dagstuhl Seminar on Numerical Software with Result Verification*, LNCS 2991, 2004.
22. J. Stolfi and L. H. de Figueiredo. An introduction to affine arithmetic. In *TEMA Tend. Mat. Apl. Comput.*, 4, No. 3 (2003), 297-312.
23. E. Goubault and S. Putot. Fluctuat user manual, CEA report 2006, available by asking the authors.
24. Grammatech Inc. CodeSonar, overview. See <http://www.grammatech.com/products/codesonar/overview.html>.
25. PolySpace Technologies. PolySpace for hand-written code. See <http://www.polyspace.fr/products.htm>.
26. LIP6. The CADNA Library. See <http://www-anp.lip6.fr/cadna/Accueil.php>.