

Title	Combining Testing and Static Analysis to Overflow and Roundoff Error Detection
Author(s)	Do, Ngoc Thi Bich; Ogawa, Mizuhito
Citation	Research report (School of Information Science, Japan Advanced Institute of Science and Technology), IS-RR-2010-004: 1-27
Issue Date	2010-06-21
Type	Technical Report
Text version	publisher
URL	http://hdl.handle.net/10119/9018
Rights	
Description	リサーチレポート (北陸先端科学技術大学院大学情報科学研究科)

Combining Testing and Static Analysis to Overflow and Roundoff Error Detection

Do Thi Bich Ngoc · Mizuhito Ogawa.

Abstract This paper proposes a technique for automatic detection of overflow and roundoff errors, caused by the floating-point number to fixed-point number conversion. First, a new range representation, “extended affine interval”, is proposed to overapproximate overflow and roundoff errors. Second, the overflow and roundoff error analysis problem is encoded as a weighted model checking, which is implemented as a static analyzer CANA. Last, we propose a new testing refinement loop, called “counterexample-guided narrowing”, by combining the static analysis and testing. They are composed and implemented in a prototype tool, CANAT, in which analysis results are used not only for possible roundoff error detection, but also for finding dominant error factors in input parameters. To avoid widening, currently we focus on programs with bounded loops and arrays with fixed length, which typically appear in encoder/decoder reference algorithms. Experimental results on small programs show that the extended affine interval is much more precise than classical interval, and the counterexample-guided narrowing approach outperforms the random testing technique.

Keywords Software verification · Static analysis · Testing · Roundoff error · Overflow error · Affine interval · Classical interval

1 Introduction

In the computers, a real number has a finite representation, i.e., either a floating-point number, or a fixed-point number, which may introduce overflow and roundoff errors (ORE). OREs cause tricky behavior; for instance, consider

$$(333.75 - a^2)b^6 + a^2(11a^2b^2 - 121b^4 - 2) + 5.5b^8 + (a/(2b))$$

for $a = 77617$, $b = 33096$. This is known as Rump’s example [31] and IEEE 754 standard floating operations [22] return the results

Single precision	1.172604
Double precision	1.1726039400531786
Fourfold precision	1.17260394005317863185883490452011838

which seem that the single precision is enough. However, if we symbolically compute with rational number representations, it will result $-54767/66192$ (approx. -0.8273960599).

Overflow and roundoff errors problem

A famous example that caused a serious disaster by roundoff errors is “The Patriot Missile Failure”¹.

On February 25, 1991, during the Gulf War, an American Patriot Missile battery in Dharan, Saudi Arabia, failed to track and intercept an incoming Iraqi Scud missile. The Scud struck an American Army barracks, killing 28 soldiers and injuring around 100 other people.

Do Thi Bich Ngoc
School of Information Science - Japan Advanced Institute of Science and Technology
Tel.: +84-0761-51-1247
Fax: +84-0761-51-1149
E-mail: dongoc@jaist.ac.jp

Mizuhito Ogawa
School of Information Science - Japan Advanced Institute of Science and Technology
Tel.: +84-0761-51-1247
Fax: +84-0761-51-1149
E-mail: mizuhito@jaist.ac.jp

¹ <http://www.ima.umn.edu/arnold/disasters/patriot.html>


```

/* CANAT
CANAT ALL sign 11 4
maintest x range -1 3
maintest y range -10 10
_test global rst 0.26
*/
typedef float Real;
Real rst;
Real maintest(Real x, Real y){
1.   if (x>0)
2.     {rst=x*x;}
3.   else rst = 3*x;
4.   rst = rst - y;
5.   return rst;
}

```

Fig. 2 An example of a C program

1. Whether the largest RE of a result lies within given threshold?
2. Whether overflow error may occur?
3. If they occur, where?

We say that the program is “safe” if for all inputs, REs of the result lie in $[-\theta, \theta]$.

Example 1 Fig. 2 shows a C program with annotations that:

- initial ranges of x, y : $x \in [-1, 3]$, $y \in [-10, 10]$,
- fixed-point format (11 : 4), and
- RE threshold is $\theta = 0.26$

Note that base $b = 2$.

The questions are:

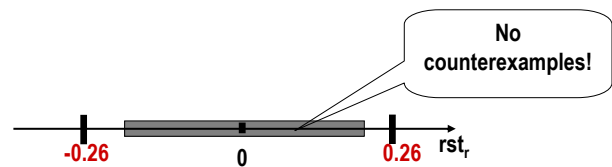
1. Does RE of rst lie within $[-0.26, 0.26]$?
2. May overflow error occur? Where?

The OREs will be propagated through computations of the program. Further, the computations themselves cause OREs because the arithmetic needs to round the result to fit the number format. Besides, OREs are also affected by types of statements (e.g., branch, loop, assignment).

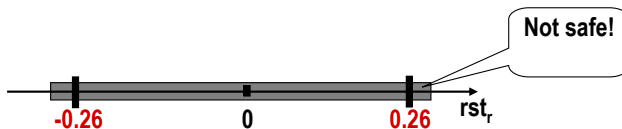
Testing vs static analysis

We may face a situation that an ORE analysis reports that the roundoff error of the result exceeds the roundoff error threshold, but a test cannot find any counterexamples.

Example 2 Assume that for the program in Fig. 2, the initial ranges of x, y are $[-1, 3]$, $[-10, 10]$, respectively. The conversion from the floating-point type to fixed-point type such that the width of integer part is 11 and the width of fraction part is 4. It is *safe* if no overflow errors occur and no roundoff errors of rst go beyond $[-0.26, 0.26]$.



a. roundoff error of rst found by testing $\subseteq [-0.20, 0.21]$



b. roundoff error of rst found by analysis $\subseteq [-0.28, 0.28]$

Fig. 3 Results of analyzing and testing C program in Fig. 2

By random 100 test cases, all roundoff errors lie in the range $[-0.20, 0.21] \subset [-0.26, 0.26]$, which means no counterexamples are found (Fig. 3 a).

The ORE analysis (in Section 4) reports that the roundoff error of rst lies in $[-0.28, 0.28]$, which exceeds the roundoff error threshold $[-0.26, 0.26]$ (Fig. 3 b).

Then, both testing and analysis cannot clarify whether the program in Fig. 2 is safe.

The challenge is how to bridge the gap between testing and static analysis?

Main results and paper structure

This paper proposes a technique for automatic detection of overflow and roundoff errors, caused by the floating-point number to fixed-point number conversion. First, inspired by the range representation to under approximate OREs [21], a new range representation, *extended affine interval* (EAI) is proposed to over approximate OREs. EAI does not increase the number of noise symbols.

Second, the overflow and roundoff error analysis problem is encoded as a weighted model checking, which is implemented as a static analyzer CANA. A weight domain is designed based on the EAI arithmetic.

Last, we propose a new testing refinement loop, called “counterexample-guided narrowing”, by combining the static analysis and testing. This refinement of testing is performed by detected OREs by CANA, in which EAI shows dominant error factors among input parameters. They are composed as a prototype tool, CANAT.

To avoid widening, currently we focus on programs with bounded loops and arrays with fixed length, which typically appear in encoder/decoder reference algorithms.

Experimental results on small programs show that the extended affine interval is much more precise than

classical interval, and the counterexample-guided narrowing approach clearly outperforms the random testing technique.

The rest of this paper is organized as follows. Section 2 formally presents ORE problems. In Section 3, we introduce various interval arithmetic. We also propose a new range representation, *EAI*, adding to traditional CI and AI. The ORE analysis based on weighted model checking is introduced in Section 4. We also describe the implementation of the proposed framework in a tool, *CANA*. Section 5 proposes the counterexample-guided narrowing approach to detect REs and its implementation, *CANAT*. Section 6 overviews related works. Finally, Section 7 concludes the paper and indicates future works. The content of Section 3 to 4 and that of Section 5 were preliminary reported in [44] and [45], respectively.

2 Representation of Real Numbers and the ORE Problem

We first present overflow and roundoff error (ORE) problem when represent real numbers in computers such as floating-point numbers (Subsection 2.1) and fixed-point numbers (Subsection 2.2). In Subsection 2.3, we introduce ORE arithmetics, which decomposes a number into a pair of floating-point (or fixed-point) and a roundoff error evaluation.

2.1 Floating-point Numbers and ORE problem

2.1.1 Floating-point Numbers

Floating-point numbers are often used to represent real numbers in numerical computation. In a floating-point number, the position of the radix point is dynamic.

Definition 1 A **floating-point number** x has a representation in **base** b , with **sign** s , **significand** m , and **exponent** e , such that

$$x = (-1)^s \times m \times b^e \quad (1)$$

where s is 0 or 1, $m = d_0.d_1\dots d_{p-1}$ with $0 \leq d_i < b$, and e is an integer. The **set of floating-point numbers** is denoted by R_{fl}

Remark 1 In order to optimize the quantity of representable numbers, floating-point numbers are typically in *normalized form*, which puts the radix point after the first non-zero digit (e.i., $d_0 \neq 0$).

Example 3 The **decimal number** $x = 8.75$, represented as $(-1)^0 \times 0.875 \times 10^1$, has $s = 0, m = 0.875, e = 1$. Its equivalent **binary** format is $x = (-1)^0 \times (0.100011) \times 2^{100}$ with $s = 1, m = 0.100011, e = 100$. The corresponding **normal floating-point number** is $x = (-1)^0 \times (1.00011) \times 2^{10}$ with $s = 1, m = 1.00011, e = 10$.

The *floating-point format* $(b, p, emax)$ determines a set of representable floating-point numbers, in which:

- b is base (e.g. 2 or 10)
- p is number of digits in the significand
- $emax$ is the maximum value of exponent e (the minimum value of e is $emin = 1 - emax$).

We basically follow the IEEE7542008 standard [22], shown in Table 1. Thus, a normal floating-point number closest to zero is $\pm b^{emin}$ and a number farthest from zero is $\pm(b - b^{1-p}) \times b^{emax}$. For instance, in the binary 64 floating-point format,

- The number closest to zero is

$$\pm 2^{-1022} \approx \pm 2.225073858507202010^{-308}$$

- The number farthest from zero is

$$\pm ((1 - (1/2)^{53}) 2^{1024}) \approx \pm 1.7976931348623157 \times 10^{308}$$

2.1.2 OREs of Floating-point Numbers

Since floating-point numbers have the finite precision, *roundoff error* may occur due to the finite fraction part, and *overflow error* may occur due to the finite integer part.

Roundoff error (RE) If the significand m of x is represented by more than p bits, x will be truncated (or chopped) in some way. The IEEE7542008 standard defines four rounding algorithms [22].

- *Round to Nearest*: This is the default mode. In this mode results are rounded to the nearest representable value. If the result is midway between two representable values, the even representable is chosen. Even here means the lowest-order bit is zero.
- *Round toward 0*: All results are rounded to the largest representable value whose magnitude is less than that of the result. In other words, if the result is negative it is rounded up; if it is positive, it is rounded down.
- *Round toward $+\infty$* : All results are rounded to the smallest representable value, which is greater than the result.
- *Round toward $-\infty$* : All results are rounded to the largest representable value, which is less than the result.

Definition 6 (Real-to-fixed ORE comparison operators) Let (x_f, x_r) , and (y_f, y_r) be pairs of fixed-point parts and REs of real numbers x, y .

$$((x_f, x_r) < (y_f, y_r)) =$$

$$\begin{cases} \text{true} & \text{if } (x_f + x_r < y_f + y_r) \wedge (x_f < y_f) \\ \text{false} & \text{if } (x_f + x_r \geq y_f + y_r) \wedge (x_f \geq y_f) \\ \text{unknown} & \text{otherwise} \end{cases}$$

$$((x_f, x_r) == (y_f, y_r)) =$$

$$\begin{cases} \text{true} & \text{if } (x_f = y_f \wedge x_r = y_r) \\ \text{false} & \text{if } (x_f, x_r) < (y_f, y_r) \vee (y_f, y_r) < (x_f, x_r) \\ \text{unknown} & \text{otherwise} \end{cases}$$

Other comparison operators, such as $>$, $!$, $=$, are defined using the operators above.

Example 9 Let $x = 34.5678$, $y = 98.76543$. We assume the fixed-point format $(b = 10, ip = 3, fp = 2)$, “round toward $-\infty$ ”, and the RE threshold $\theta = 0.01$. We have:

- The fixed-point value of x is $x_{fx} = 34.56$ and the corresponding RE is $x_r = 0.0078$
- The fixed-point value of y is $y_{fy} = 98.76$ and the corresponding RE is $y_r = 0.00543$

We next show how to evaluate ORE arithmetic:

- **Addition:**

$$\begin{aligned} (x_f, x_r) \boxplus (y_f, y_r) &= (rd_{fp}(34.56 + 98.76), 0.0078 + \\ &\quad 0.00543 + re_{fp}(34.56 + 98.76)) \\ &= (133.32, 0.01323) \end{aligned}$$

That means the result of addition is 133.32 and its RE is $0.01323 > \theta$. Thus, an RE is detected.

- **Multiplication:**

$$\begin{aligned} (x_f, x_r) \boxtimes (y_f, y_r) &= \\ (rd_{fp}(34.56 \times 98.76), 0.0078 \times 98.76 + 34.56 \times 0.00543 \\ + 0.0078 \times 0.00543 + re_{fp}(34.56 \times 98.76)) \\ &= (3413.14, 0.963631154) \end{aligned}$$

That means the result of multiplication is 3413.14 ($> 10^3$) and its RE is $0.963631154 > \theta$. Thus, both an OE and an RE are detected.

2.3.2 Real-to-Float ORE Arithmetic

For a floating-point format $(b, p, emax)$ and a real number x , we denote the **floating-point part** by $rd_{fl}(x)$ and the **RE** by $re_{fl}(x) (= x - rd_{fl}(x))$. If $rd_{fl}(x) > (b - b^{1-p} \times b^{emax})$, we conclude that OE occurs, and if $re_{fl}(x) > \theta$ (where θ is predefined threshold) we conclude that RE occurs. The following definition describes the rules of propagating ORE when converting real numbers to floating-point numbers.

Definition 7 (Real-to-Float ORE arithmetic) Let (x_f, x_r) and (y_f, y_r) be pairs of floating-point parts and REs of real numbers x, y . Real-to-Float ORE arithmetic $\boxplus = \{\boxplus, \boxminus, \boxtimes, \boxdiv\}$ is defined below.

$$\begin{aligned} (x_f, x_r) \boxplus (y_f, y_r) &= \\ (rd_{fl}(x_f + y_f), x_r + y_r + re_{fl}(x_f + y_f)) \\ (x_f, x_r) \boxminus (y_f, y_r) &= \\ rd_{fl}(x_f - y_f), x_r - y_r + re_{fl}(x_f - y_f) \\ (x_f, x_r) \boxtimes (y_f, y_r) &= \\ (rd_{fl}(x_f \times y_f), x_r \times y_f + x_f \times y_r + x_r \times y_r + re_{fl}(x_f \times y_f)) \\ (x_f, x_r) \boxdiv (y_f, y_r) &= \\ (rd_{fl}(x_f \div y_f), (x_f + x_r) \div (y_f + y_r) - x_f \div y_f + re_{fl}(x_f \div y_f)) \end{aligned}$$

Real-to-float ORE comparison operators are defined similar Definition 6.

2.3.3 Float-to-Fixed ORE Arithmetic

For a floating-point number x , the floating-point format $(b, p, emax)$, and the fixed-point format (b, ip, fp) , we denote the **fixed-point part** by $rd_{fx}(x)$ and the **RE** by $re_{fx}(x) (= re_{fx}(x) - re_{fl}(x))$. If $rd_{fx}(x) > b^{ip}$ we conclude that an OE occurs, and if $re_{fx}(x) > \theta$ (where θ is predefined threshold) we conclude that an RE occurs. The following definition describes the rules of propagating ORE between floating-point numbers and fixed-point numbers.

Definition 8 (Float-to-Fixed ORE arithmetic) Let (x_f, x_r) and (y_f, y_r) be pairs of fixed-point parts and REs of floating-point numbers x, y . Float-to-Fixed ORE arithmetic $\boxplus = \{\boxplus, \boxminus, \boxtimes, \boxdiv\}$ is defined below.

$$\begin{aligned} (x_f, x_r) \boxplus (y_f, y_r) &= \\ (rd_{fx}(x_f + y_f), x_r + y_r + re_{ff}(x_f + y_f)) \\ (x_f, x_r) \boxminus (y_f, y_r) &= \\ (rd_{fx}(x_f - y_f), x_r - y_r + re_{ff}(x_f - y_f)) \\ (x_f, x_r) \boxtimes (y_f, y_r) &= \\ (rd_{fx}(x_f \times y_f), x_r \times y_f + x_f \times y_r + x_r \times y_r + re_{ff}(x_f \times y_f)) \\ (x_f, x_r) \boxdiv (y_f, y_r) &= \\ (rd_{fx}(x_f \div y_f), \\ (x_f + x_r) \div (y_f + y_r) - x_f \div y_f + re_{ff}(x_f \div y_f)) \end{aligned}$$

Float-to-fixed ORE comparison operators are defined similar Definition 6.

More precise ORE estimation

When we fix the conversion, such as from the floating-point IEEE 754 binary64 (2, 53, 1024) to the fixed-point (2, ip, fp) with size 2 bytes (e.i., $ip + fp = 16$) (which frequently appears in practice), we can obtain better

estimation of OREs. Assume “round to nearest” in Definition 8.

Let $\delta_\circ = re_{fx}(x \circ y) - re_{fl}(x_f \circ y_f)$ where $\circ \in \{+, -, \times, \div\}$. We now find the bound of δ_\circ by considering the bound of $re_{fl}(x_f \circ y_f)$ and $re_{fx}(x \circ y)$:

- *Floating-point roundoff error* $rd_{fl}(x \circ y)$:
Assume $rd_{fl}(x \circ y) = -(1)^s \times m \times b^e$, we have $|re_{fl}(x \circ y)| < 2^{-53+e}/2$. Without loss of generality, we can assume $e \leq ip$ (otherwise an OE occurs in the fixed-point operator $(x_f \circ y_f)$). Thus, we have:

$$\begin{aligned} |re_{fl}(x \circ y)| &< 2^{-53+e}/2 \\ &< 2^{-53+ip}/2 \\ &< 2^{-53+16-fp}/2 \text{ (because } ip + fp = 16) \\ &< 2^{-38-fp} \end{aligned}$$

- *Fixed-point roundoff error* $re_{fx}(x_f \circ y_f)$:
Because the fixed-point format is unique, the results of the addition and the subtraction have the same format. Thus, $re_{fx}(x_f \circ y_f) = 0$ for $\circ \in \{+, -\}$.
For the multiplication, the fraction part of the result has $2 \times fp$ digits. The fraction part is round to fp digits, and $|re_{fx}(x_f \times y_f)| < 2^{fp}/2 - 2^{2 \times fp}/2$.
For the division, similarly $|re_{fx}(x_f \div y_f)| < 2^{fp}/2$.

Hence, we obtain the **Float64-to-Fixed16 ORE arithmetic** by replacing $re_{ff}(x \circ y)$ with δ_\circ , where $\circ \in \{+, -, \times, \div\}$ and

$$\begin{cases} |\delta_+| < 2^{-38-fp} \\ |\delta_-| < 2^{-38-fp} \\ |\delta_\times| < 2^{fp-1} - 2^{2 \times fp-1} + 2^{-38-fp} \\ |\delta_\div| < 2^{fp-1} + 2^{-38-fp} \end{cases} \quad (2)$$

3 Interval Arithmetics

In order to estimate OREs of arithmetic operations, there are two known range representations: classical interval [40] and affine interval [54, 55]. In this section, we firstly describe these two methods in detail in Subsections 3.1 and 3.2. Then, inspired by the idea in [21] for under approximation, we propose an “extended affine interval” for overapproximation in Subsection 3.3. Lastly, we represent how to implement these intervals on computers using floating-point type in Subsection 3.4.

3.1 Classical Interval

Classical interval (CI) was introduced in the 1960s by Moore [40] as an approach to putting bounds on rounding errors in mathematical computations. In CI, the upper and the lower bounds describe possible values.

Definition 9 A **classical interval** of x is an interval $\bar{x} = [x_l, x_h]$ with $x_l \leq x \leq x_h$. The set of classical intervals is denoted by \bar{R} .

Definition 10 **CI arithmetic** consists of operations $\{\bar{+}, \bar{-}, \bar{\times}, \bar{\div}\}$ on pairs of CIs defined below:

$$\begin{aligned} [x_l, x_h] \bar{+} [y_l, y_h] &= [x_l + y_l, x_h + y_h] \\ [x_l, x_h] \bar{-} [y_l, y_h] &= [x_l - y_h, x_h - y_l] \\ [x_l, x_h] \bar{\times} [y_l, y_h] &= [\min(x_l y_l, x_l y_h, x_h y_l, x_h y_h), \\ &\quad \max(x_l y_l, x_l y_h, x_h y_l, x_h y_h)] \\ [x_l, x_h] \bar{\div} [y_l, y_h] &= [x_l, x_h] \bar{\times} [\frac{1}{y_h}, \frac{1}{y_l}] \text{ if } 0 \notin [y_l, y_h] \end{aligned}$$

The following example demonstrates how to compute CI operations:

Example 10 For $x \in \bar{x} = [-1, 3]$, $y \in \bar{y} = [-6, 10]$. Let us compute the bound of $z = x \circ y$ ($\circ \in \{+, -, \times, \div\}$) using CI:

- Addition $z = x + y$:

$$\begin{aligned} \bar{z} &= \bar{x} \bar{+} \bar{y} \\ &= [-1, 3] \bar{+} [-6, 10] \\ &= [-1 - 6, 3 + 10] \\ &= [-7, 13] \end{aligned}$$

- Subtraction $z = x - y$:

$$\begin{aligned} \bar{z} &= \bar{x} \bar{-} \bar{y} \\ &= [-1, 3] \bar{-} [-6, 10] \\ &= [-1 - 10, 3 - (-6)] \\ &= [-11, 9] \end{aligned}$$

- Multiplication $z = x \times y$:

$$\begin{aligned} \bar{z} &= \bar{x} \bar{\times} \bar{y} \\ &= [-1, 3] \bar{\times} [-6, 10] \\ &= [\min\{6, -10, -18, 30\}, \max\{6, -10, -18, 30\}] \\ &= [-18, 30] \end{aligned}$$

- Division $z = x \div y$:

$$\begin{aligned} \bar{z} &= \bar{x} \bar{\div} \bar{y} \\ &= [-1, 3] \bar{\div} [-6, 10] \end{aligned}$$

Since $0 \in [-6, 10]$, we cannot compute the bound of z ; instead a “division by zero” warning occurs.

For $\bar{x}, \bar{x}_1, \dots, \bar{x}_n \in \bar{R}$, $\circ \in \{\bar{+}, \bar{-}, \bar{\times}, \bar{\div}\}$, and a constant c , we denote:

- $\bar{x}_1 \bar{x}_2 = \bar{x}_1 \bar{\times} \bar{x}_2$, $c \bar{x} = \bar{x} c = \bar{x} \bar{\times} [c, c]$,
- $c \circ \bar{x} = [c, c] \circ \bar{x}$, $\bar{x} \circ c = \bar{x} \circ [c, c]$, and
- $\sum_{i=1}^n \bar{x}_i = \bar{x}_1 \bar{+} \bar{x}_2 \bar{+} \dots \bar{+} \bar{x}_n$.

CI arithmetic assumes that all intervals are independent, even if exact values are dependent. The next example illustrates such a problem.

Example 11 Let $x \in \bar{x} = [-1, 3]$. It is easy to see that:

$$\begin{aligned}\bar{x} - \bar{x} &= [-1, 3] - [-1, 3] \\ &= [-4, 4]\end{aligned}$$

CI arithmetic assumes the first operand and the second operand to be independent, while in fact, they represent the same value x and the exact result is $[0, 0]$.

This leads to a great loss of precision in a long computation chain, which is called “error explosion”.

3.2 Affine Interval

Affine interval (AI) was introduced by Stolfi [54,55] as a model for self-validated numerical analysis. It was proposed to address the “error explosion” problem in conventional CI. In AI, the value are represented as affine combinations (affine forms) of certain primitive noise symbols, which stand for sources of uncertainty in the data or approximations made during the computation.

Definition 11 An **Affine interval** of x is a formula

$$\ddot{x} = x_0 + x_1\varepsilon_1 + x_2\varepsilon_2 + \cdots + x_n\varepsilon_n$$

with $x \in [x_0 - \sum_i^n |x_i|, x_0 + \sum_i^n |x_i|]$. x_0 is called the **central value**. For each $i \in [1, n]$, $\varepsilon_i \in [-1, 1]$ is a **noise symbol**, which stands for an independent component of the total uncertainty. The set of affine interval forms is denoted by \ddot{R} .

In AI arithmetic, the results of linear operations (e.i., addition, subtraction) are straightforward operations on AIs. However, the results of nonlinear operations (e.i., multiplication, division) are not AI forms. Hence, we need to approximate the nonlinear parts of the results by introducing new noise symbols.

Definition 12 **AI arithmetic** consists of operations $\{\ddot{+}, \ddot{-}, \ddot{\times}, \ddot{\div}\}$ on pairs of AIs as defined below. Let $\ddot{x} = x_0 + \sum_{i=1}^n x_i\varepsilon_i$ and $\ddot{y} = y_0 + \sum_{i=1}^n y_i\varepsilon_i$. AI operations are as defined below:

$$\begin{aligned}\ddot{x} \ddot{+} \ddot{y} &= (x_0 + y_0) + \sum_{i=1}^n (x_i + y_i)\varepsilon_i \\ \ddot{x} \ddot{-} \ddot{y} &= (x_0 - y_0) + \sum_{i=1}^n (x_i - y_i)\varepsilon_i \\ \ddot{x} \ddot{\times} \ddot{y} &= x_0y_0 + \sum_{i=1}^n (x_0y_i + x_iy_0)\varepsilon_i + B\varepsilon_{n+1} \\ \ddot{x} \ddot{\div} \ddot{y} &= \ddot{x} \ddot{\times} \left(\frac{1}{\ddot{y}}\right), \text{ if } 0 \notin [x_0 - \sum_i^n |x_i|, x_0 + \sum_i^n |x_i|]\end{aligned}$$

where $\varepsilon_{n+1} \in [-1, 1]$ is a new noise symbol, B is the maximum value of $(\sum_{i=1}^n x_i\varepsilon_i)(\sum_{i=1}^n y_i\varepsilon_i)$, and $\frac{1}{\ddot{y}}$ is computed by Chebyshev approximation [54].

The range of values described an AI is evaluated by replacing each noise symbol ε_i with $[-1, 1]$.

The advantage of AI is precision on linear operations, compared to CI. For instance, in Example 11, values in a CI $\bar{x} = [-1, 3]$ are equivalently described by an AI $\ddot{x} = 1 + 2\varepsilon_x$. Then,

$$\ddot{z} = \ddot{x} \ddot{-} \ddot{x} = 0$$

which shows the exact result of the subtraction $(x - x)$.

For multiplication, there are choices to approximate B , and a direct approximation of B is $(\sum_{i=1}^n |x_i|)(\sum_{i=1}^n |y_i|)$. For division, we apply Chebyshev approximation below.

Chebyshev approximation in division

Chebyshev approximation aims to minimize the maximum absolute error. Let \mathcal{F} be some space of functions, e.g., polynomials, affine forms. An element of \mathcal{F} that minimizes the maximum absolute difference from a given function f over a specified domain Ω is known as a Chebyshev (or minimax) \mathcal{F} -approximation to f over Ω . We briefly overview the results in [54].

For univariate functions, the minimax affine approximation is characterized by the following property.

Theorem 1 [54] *Let f be a bounded and continuous function from some closed and bounded interval $I = [a, b]$ to R . Let h be the affine function that best approximates f in I under the minimax error criterion. Then, there exist three distinct points $u, v, w \in I$ where the error $f(x) - h(x)$ has maximum magnitude; and the sign of the error alternates when the three points are considered in ascending order.*

This theorem provides an algorithm for finding the optimum approximation in many cases, via the following corollary:

Corollary 1 [54] *Let f be a bounded and twice differentiable function defined on some interval $I = [a, b]$, whose seconde derivative f'' does not change sign inside I . Let $f^a(x) = \alpha x + \zeta$ be its minimax affine approximation in I . Then:*

- *The coefficient α is simply $(f(b) - f(a))/(b - a)$, the slope of the line $r(x)$ that interpolates the points $(a, f(a))$ and $(b, f(b))$.*
- *The maximum absolute error will occur twice (with the same sign) at the endpoints a and b of the range, and once (with the opposite sign) at every interior point u of I where $f'(u) = \alpha$.*
- *The independent term ζ is such that $\alpha u + \zeta = (f(u) + r(u))/2$, and the maximum absolute error is $\delta = |f(u) - r(u)|/2$.*

This result gives us a method for finding the optimum coefficients α and ζ , as long as we can solve the equation $f'(u) = \alpha$.

If an AI \ddot{y} includes zero, a “division by zero” warning occurs. Thus, we only consider the cases $\bar{y} = [l, h]$ are entirely either positive or negative (i.e., $l > 0$ or $h < 0$). The Chebyshev approximation of $\frac{1}{\ddot{y}}$ (Fig. 4) is computed as follows:

- $a = \min\{|l|, |h|\}$, $b = \max\{|l|, |h|\}$.
- $\alpha = -1/b^2$.
- $d_{max} = 1/a - \alpha a$, $d_{min} = 1/b - \alpha b$.
- $\zeta = (d_{min} + d_{max})/2$, if $l < 0$ then $\zeta = -\zeta$.
- $\delta = (d_{max} - d_{min})/2$.
- $\frac{1}{\ddot{y}} = \alpha\ddot{y} + \zeta + \delta\varepsilon_k$, where ε_k is a new noise symbol.

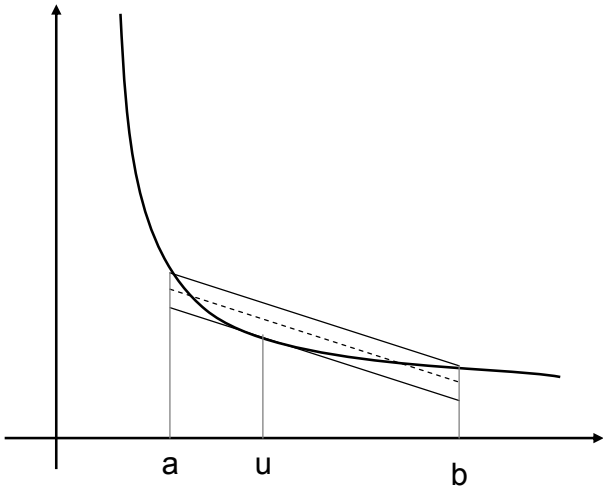


Fig. 4 Chebyshev approximation for $\frac{1}{\ddot{y}}$

Conversion between CI and AI

Standard range representation is CI. To apply AI, conversion between them are needed.

- *CI to AI*: Given a CI $\bar{x} = [l, h]$, a corresponding AI is $\ddot{x} = \frac{l+h}{2} + \frac{h-l}{2}\varepsilon_k$. Under valuations of noise symbol ε_k to $[-1, 1]$, they represent the same range. This is called *AI coercion*.
- *AI to CI*: An AI $\ddot{x} = x_0 + \sum_{i=1}^n x_i \varepsilon_i$ is projected to a CI $\bar{x} = [x_0 - \sum_{i=1}^n |x_i|, x_0 + \sum_{i=1}^n |x_i|]$. This projection loses information on dependency of uncertainty. This is called *AI projection*.

The following example demonstrates how to propagate the ranges by using AI:

Example 12 For $x \in \bar{x} = [-1, 3]$, $y \in \bar{y} = [-6, 10]$. The corresponding AI coercions of \bar{x} , \bar{y} are:

- $\ddot{x} = 1 + 2\varepsilon_x$

$$- \ddot{y} = 2 + 8\varepsilon_y$$

Let us compute the bound of $z = x \circ y$ ($\circ \in \{+, \times\}$) using AI:

- Addition $z = x + y$:

$$\begin{aligned} \ddot{z} &= \ddot{x} \dot{+} \ddot{y} \\ &= (1 + 2\varepsilon_x) \dot{+} (2 + 8\varepsilon_y) \\ &= 3 + 2\varepsilon_x + 8\varepsilon_y \end{aligned}$$

The AI projection of \ddot{z} is $[3 - 2 - 8, 3 + 2 + 8] = [-7, 13]$.

- Multiplication $z = x \times y$:

$$\begin{aligned} \ddot{z} &= \ddot{x} \dot{\times} \ddot{y} \\ &= (1 + 2\varepsilon_x) \dot{\times} (2 + 8\varepsilon_y) \\ &= 2 + 4\varepsilon_x + 8\varepsilon_y + 16\varepsilon_1 \end{aligned}$$

where ε_1 is new noise symbol standing for $\varepsilon_x \varepsilon_y$.

The AI projection of \ddot{z} is $[2 - 4 - 8 - 16, 2 + 4 + 8 + 16] = [-26, 30]$.

3.3 Extended Affine Interval

AI is more precise than CI for linear operations, but each time we perform a nonlinear operation, it introduces a new noise symbol. This would be problematic for a program with a large number of nonlinear operations.

In [21], instead of introducing new noise symbols, coefficients of noise symbols are replaced with CIs. Arithmetic operations are designed for under approximation, and we apply similar ideas for overapproximation. This is called an *extended affine interval* (EAI), which also avoids introduction of new noise symbols for nonlinear operations.

Definition 13 An **extended affine interval** of x is a formula

$$\hat{x} = \bar{x}_0 \dot{+} \sum_{k=1}^n \bar{x}_k \varepsilon_k$$

with $x \in \bar{x}_0 \dot{+} \sum_{k=1}^n \bar{x}_k [-1, 1]$, where $\varepsilon_i \in [-1, 1]$ is a noise symbol for each $i \in [1, n]$ and $\bar{x}_j \in \bar{R}$ for each $j \in [0, n]$. The set of extended affine intervals is denoted by \hat{R} .

The linear operations of EAI arithmetic are designed similarly to those of AI arithmetic. For nonlinear operations, unlike AI, EAI arithmetic does not need to introduce new noise symbols. The results of nonlinear operations approximate nonlinear parts, with CI coefficients. For example, let us consider the multiplication of two EAIs. Let $\hat{x} = \bar{x}_0 \dot{+} \sum_{i=1}^n \bar{x}_i \varepsilon_i$, $\hat{y} =$

$\bar{y}_0 \bar{\mp} \sum_{i=1}^n \bar{y}_i \varepsilon_i$. Without loss of generality, assume that $\sum_{k=1}^n \bar{y}_k [-1, 1] \subseteq \sum_{k=1}^n \bar{x}_k [-1, 1]$. We have:
 $\hat{x} \hat{\times} \hat{y} = (\bar{x}_0 \bar{\mp} \sum_{i=1}^n \bar{x}_i \varepsilon_i) \hat{\times} (\bar{y}_0 \bar{\mp} \sum_{i=1}^n \bar{y}_i \varepsilon_i)$
 $= \bar{x}_0 \bar{y}_0 \bar{\mp} \sum_{i=1}^n (\bar{x}_0 \bar{y}_i \bar{\mp} \bar{x}_i \bar{y}_0 \bar{\mp} \bar{x}_i B) \varepsilon_i$, where $B = \sum_{i=1}^n \bar{y}_i \varepsilon_i$. A direct approximation of B is $\sum_{k=1}^n \bar{y}_k [-1, 1]$. Formally, EAI arithmetic is defined as follows:

Definition 14 EAI arithmetic consists of operations $\{\hat{+}, \hat{-}, \hat{\times}, \hat{\div}\}$ on pairs of EAI.

Let $\hat{x} = \bar{x}_0 \bar{\mp} \sum_{i=1}^n \bar{x}_i \varepsilon_i$, $\hat{y} = \bar{y}_0 \bar{\mp} \sum_{i=1}^n \bar{y}_i \varepsilon_i$, $\bar{X} = \sum_{k=1}^n (\bar{x}_k [-1, 1])$, and $\bar{Y} = \sum_{k=1}^n (\bar{y}_k [-1, 1])$. Then,

$$\begin{aligned} \hat{x} \hat{+} \hat{y} &= (\bar{x}_0 \bar{\mp} \bar{y}_0) \bar{\mp} \sum_{i=1}^n (\bar{x}_i \bar{\mp} \bar{y}_i) \varepsilon_i \\ \hat{x} \hat{-} \hat{y} &= (\bar{x}_0 \bar{-} \bar{y}_0) \bar{\mp} \sum_{i=1}^n (\bar{x}_i \bar{-} \bar{y}_i) \varepsilon_i \\ \hat{x} \hat{\times} \hat{y} &= \bar{x}_0 \bar{y}_0 \bar{\mp} \sum_{i=1}^n (\bar{x}_0 \bar{y}_i \bar{\mp} \bar{x}_i \bar{y}_0) \varepsilon_i \bar{\mp} B \\ \hat{x} \hat{\div} \hat{y} &= \hat{x} \hat{\times} (\frac{1}{\hat{y}}) \quad \text{if } 0 \notin \bar{x}_0 \bar{\mp} \sum_{k=1}^n \bar{x}_k [-1, 1] \end{aligned}$$

where:

$$B = \begin{cases} (\sum_{i=1}^n \bar{x}_i \varepsilon_i) \bar{Y} & \text{if } \bar{Y} \subseteq \bar{X} \\ \bar{X} (\sum_{i=1}^n \bar{y}_i \varepsilon_i) & \text{otherwise} \end{cases}$$

and $\frac{1}{\hat{y}}$ is computed by Chebyshev approximation [54].

Similar to AI arithmetic, the commutative property holds for both addition and multiplication; the associative property only holds for addition; and the distributive property does not hold.

Remark 2 The overapproximation B may conceal some noise symbols. If we are sensitive to this matter, B can be modified as:

$$B = \alpha (\sum_{i=1}^n \bar{x}_i \varepsilon_i) \bar{Y} \bar{\mp} \beta \bar{X} (\sum_{i=1}^n \bar{y}_i \varepsilon_i)$$

$$\text{with } \alpha = \frac{|\bar{X}|}{(|\bar{X}| + |\bar{Y}|)} \text{ and } \beta = (1 - \alpha).$$

Conversion between CI and EAI

To apply EAI, the conversion between them are needed.

- *CI to EAI*: Given a CI $\bar{x} = [l, h]$, a corresponding EAI is $\hat{x} = \frac{l+h}{2} \bar{\mp} \frac{h-l}{2} \varepsilon_k$. Under valuations of a noise symbol ε_k to $[-1, 1]$, they represent the same range. This is called *EAI coercion*.
- *EAI to CI*: An EAI $\hat{x} = \bar{x}_0 \bar{\mp} \sum_{i=1}^n \bar{x}_i \varepsilon_i$ is projected to a CI $\bar{x} = \bar{x}_0 \bar{\mp} \sum_{i=1}^n \bar{x}_i [-1, 1]$. Replacement of a noise symbol ε_k to $[-1, 1]$ loses information on dependency of uncertainty. This is called *EAI projection*.

Example 13 For $x \in \bar{x} = [-1, 3]$, $y \in \bar{y} = [-6, 10]$. The corresponding EAI coercions of \bar{x} , \bar{y} are:

$$\begin{aligned} - \hat{x} &= 1 \bar{\mp} 2\varepsilon_x \\ - \hat{y} &= 2 \bar{\mp} 8\varepsilon_y \end{aligned}$$

- Addition $z = x + y$:

$$\begin{aligned} \hat{z} &= \hat{x} \hat{+} \hat{y} \\ &= (1 \bar{\mp} 2\varepsilon_x) \hat{+} (2 \bar{\mp} 8\varepsilon_y) \\ &= 3 \bar{\mp} 2\varepsilon_x \bar{\mp} 8\varepsilon_y \end{aligned}$$

The EAI projection of \hat{z} is

$$3 \bar{\mp} 2[-1, 1] \bar{\mp} 8[-1, 1] = [-7, 13]$$

- Multiplication $z = x \times y$:

$$\begin{aligned} \hat{z} &= \hat{x} \hat{\times} \hat{y} \\ &= (1 \bar{\mp} 2\varepsilon_x) \hat{\times} (2 \bar{\mp} 8\varepsilon_y) \\ &= 2 \bar{\mp} 4\varepsilon_x \bar{\mp} 8\varepsilon_y \bar{\mp} B \end{aligned}$$

where $\bar{X} = 2[-1, 1] = [-2, 2]$ and $\bar{Y} = 8[-1, 1] = [-8, 8]$.

Because $\bar{X} \subset \bar{Y}$. Hence, $B = 8\varepsilon_y \bar{X} = [-16, 16]\varepsilon_y$. Then, $\hat{z} = 2 \bar{\mp} 4\varepsilon_x \bar{\mp} [-8, 24]\varepsilon_y$

The EAI projection of \hat{z} is

$$2 \bar{\mp} 4[-1, 1] \bar{\mp} [-8, 24]\bar{\times}[-1, 1] = [-26, 30]$$

Although EAI does not introduce new noise symbols, this does not mean EAI arithmetic is always less precise than AI arithmetic. AI arithmetic only advances in cases when we reuse the results of nonlinear parts.

Example 14 Let $z = x \times x$; $t = z - z$ and $x \in [-1, 31]$.

- AI arithmetic: $\ddot{x} = 1 + 2\varepsilon_x$, $\ddot{z} = 1 + 4\varepsilon_x + 4\varepsilon_1$ where ε_1 is introduced for multiplication $\varepsilon_x \varepsilon_x$.
 $\hat{t} = \ddot{z} - \ddot{z} = 0$
- EAI arithmetic: $\hat{x} = 1\bar{2}\varepsilon_x$, $\hat{z} = 1\bar{\mp}[0, 8]\varepsilon_x$.
 $\hat{t} = \hat{z} - \hat{z} = [-8, 8]\varepsilon_x$

In this case, AI is more precise than EAI. However, if we compute the bound of $t = x \times x - x \times x$ (without reusing the multiplication $x \times x$), both AI and EAI return the same bound.

3.4 Interval Representations by Floating-point Numbers

As pointed in [54, 55], the interval representations themselves are affected by OREs, since boundaries and coefficients are represented by floating-point numbers. We will briefly overview how to overapproximate CI, AI, and EAI.

For $x \in R$, we define:

- $\downarrow x \in R_{fl}$ is the round toward $+\infty$ of x ,
- $\uparrow x \in R_{fl}$ is the round toward $-\infty$ of x , and

Floating-point classical interval

For CI, a safe approximation is to truncate down the lower bound and truncate up the upper bound of an interval.

Definition 15 A *floating-point classical interval* of CI $\bar{x} = [l, h]$ for $l, h \in R$ is

$$\Downarrow \bar{x} = \Downarrow [l, h] = [\Downarrow l, \Uparrow h].$$

The set of floating-point classical intervals is denoted by $\Downarrow \bar{R}$.

The floating-point CI arithmetic is obtained by applying the \Downarrow operator for each operation $\circ \in \{\bar{+}, \bar{-}, \bar{\times}, \bar{\div}\}$. For $\Downarrow \bar{x}, \Downarrow \bar{y} \in \Downarrow \bar{R}$, we define $\Downarrow \bar{x} \circ \Downarrow \bar{y} = \Downarrow (\bar{x} \circ \bar{y})$.

Note that, since $\Downarrow x \leq x \leq \Uparrow x$, $\bar{x} \subseteq \Downarrow \bar{x}$, and the extended \circ gives an overapproximation, i.e., $\bar{x} \circ \bar{y} \subseteq \Downarrow (\bar{x} \circ \bar{y})$.

This is confirmed by two steps:

- For $\bar{x}, \bar{y} \in \bar{R}$, $\bar{x} \subseteq \Downarrow \bar{x}$ and $\bar{y} \subseteq \Downarrow \bar{y}$. Hence, for $\circ \in \{+, -, \times, \div\}$, $(\bar{x} \circ \bar{y}) \subseteq (\Downarrow \bar{x} \circ \Downarrow \bar{y})$.
- By definition, $(\Downarrow \bar{x} \circ \Downarrow \bar{y}) \subseteq \Downarrow (\bar{x} \circ \bar{y})$.

Floating-point affine interval

For AI, instead of truncating coefficients in an appropriate way, we simply introduce a new noise symbol.

Definition 16 An *Floating-point affine interval* of AI $\ddot{x} = x_0 + \sum_{k=1}^n x_k \varepsilon_k \in \hat{R}$ is a formula

$$\Downarrow \ddot{x} = \Downarrow x_0 + \sum_{k=1}^n \Downarrow x_k \varepsilon_k + B \varepsilon_{n+1}$$

where new noise symbol ε_{n+1} is introduced for REs and $B = \sum_{k=0}^n (\Uparrow x_k - \Downarrow x_k)$. The set of floating-point extended affine intervals is denoted by $\Downarrow \hat{R}$.

The floating-point AI arithmetic is obtained by introducing a new noise symbol for each operation of AI arithmetic. For example, let $\Downarrow \ddot{x} = \Downarrow x_0 + \sum_{k=1}^n \Downarrow x_k \varepsilon_k$, $\Downarrow \ddot{y} = \Downarrow y_0 + \sum_{k=1}^n \Downarrow y_k \varepsilon_k$ be two floating-point AI. The addition is

$$\Downarrow \ddot{x} + \Downarrow \ddot{y} = \Downarrow (\Downarrow x_0 + \Downarrow y_0) + \sum_{k=1}^n \Downarrow (\Downarrow x_k + \Downarrow y_k) \varepsilon_k + B \varepsilon_{n+1}$$

where $B = \sum_{k=1}^n (\Uparrow (\Downarrow x_k + \Downarrow y_k) - \Downarrow (\Downarrow x_k + \Downarrow y_k))$.

Floating-point extended affine interval

For EAI, we safely approximate CI coefficients by the floating-point CI.

Definition 17 A *floating-point extended affine interval* of EAI $\hat{x} = \bar{x}_0 \bar{+} \sum_{k=1}^n \bar{x}_k \varepsilon_k \in \hat{R}$ is

$$\Downarrow \hat{x} = \Downarrow \bar{x}_0 \bar{+} \sum_{k=1}^n \Downarrow \bar{x}_k \varepsilon_k.$$

The set of floating-point extended affine intervals is denoted by $\Downarrow \hat{R}$.

The floating-point EAI arithmetic is obtained by replacing each CI at a coefficient by the floating-point CI.

From now on, we will apply floating-point CI, floating-point AI, and floating-point EAI, instead of CI, AI, and EAI, respectively.

4 ORE Analysis as Weighted Model Checking

It has been suggested intimate connections between dataflow analysis and model checking [41, 50]. A program is firstly encoded into a model (*transition system*) by abstraction, and a program analysis is formulated as a model checking problem. This is nicely adopted for control flow analysis and/or classical dataflow analysis in Dragon book [1, 28]. However, as natural requests, we intend more richer dataflow, such as quantity properties with more precise treatments on conditional branches. For instance, linear constraint propagation [41], affine relation analysis [49], or ORE constraint analysis [44] are such examples. In these cases, the direct encoding will be a transition-as-an-environment-transformer, which requires all possible environments as states. This will lead the state explosion problem in model checking.

In 2003, Rep [48] proposed *weighted pushdown model checking*, in which each transition is associated with a weight. A weight directly represents dataflow, that is, how an abstract environment will be transformed, without generating explicit environments as states. This will not improve complexity in theory, but in practice we can combine with an on-the-fly generation of weights, which drastically reduces the search space during model checking.

We follow this weighted model checking approach (but without using a pushdown stack). Due to infinity (or unboundedness) of weight domains for the ORE analysis, we restrict ourselves to acyclic models only. This is a strong limitation, but our main application target is DSP algorithms, which typically consists of loops with bounded number of iterations and arrays with fixed lengths. An effective widening operation design to avoid this restriction is left for future work. We also put an input range for the ORE analysis, since a narrower input range results more precise analysis results. As side effect, this also enables us to generate weights in an on-the-fly manner.

4.1 Dataflow Analysis as Weighted Model Checking

4.1.1 Weighted Model Checking

Weighted model checking computes dataflow (or, an update of environments) by associating a *weight* to each transition in the model, and the goal is to determine the weight summary of the meet-over-all-path.

Weight domain and weighted transition system. In weighted model checking, the weight domain D is an idempotent semiring.

Definition 18 An **idempotent semiring** is a quintuple $(D, \oplus, \otimes, \mathbf{0}, \mathbf{1})$, where $\mathbf{0}, \mathbf{1} \in D$ and \oplus, \otimes are binary operators on D such that, for $a, b, c \in D$,

- (D, \oplus) is a commutative monoid with the unit $\mathbf{0}$,
- (D, \otimes) is a monoid with the unit $\mathbf{1}$,
- \otimes distributes over \oplus , i.e., $a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$ and $(a \oplus b) \otimes c = (a \otimes c) \oplus (b \otimes c)$,
- \oplus is idempotent, i.e., $a \oplus a = a$, and
- $\mathbf{0}$ is the zero element of \otimes , i.e., $a \otimes \mathbf{0} = \mathbf{0} \otimes a = \mathbf{0}$.

In the context of dataflow analysis, each element of an idempotent semiring is regarded as follows:

- $\mathbf{0}$ stands for interruption of dataflow,
- $\mathbf{1}$ stands for the identity function (i.e., no state update),
- \otimes is the composition of two successive dataflow, and
- \oplus merges two dataflow at the meet of two transition sequences.

The *weighted transition system* is then defined as a *transition system “plus” a weight domain*.

Definition 19 Let $\mathcal{P} = (P, \Delta, s_0)$ be a transition system with P to be a finite set of states, $\Delta (\subseteq P \times P)$ to be a set of transitions, and $s_0 (\in P)$ to be an initial state. A **weighted transition system** (WTS) is a triplet $\mathcal{W} = (\mathcal{P}, \mathcal{S}, f)$, where $\mathcal{S} = (D, \oplus, \otimes, \mathbf{0}, \mathbf{1})$ is an idempotent semiring and $f : \Delta \rightarrow D$ is a map that assigns a weight to each transition.

Let Δ^* be the set of all **sequences of transitions**. For $\sigma = [r_1, \dots, r_k] \in \Delta^*$, we define $v(\sigma) =_{\Delta} f(r_1) \otimes \dots \otimes f(r_k)$. If σ is a transition sequence from a state c to a state c' , we denote $c \Rightarrow^{\sigma} c'$. The set of all such sequences is denoted by $paths(c, c')$, i.e.,

$$paths(c, c') = \{\sigma \mid c \Rightarrow^{\sigma} c'\}$$

Weighted model checking. Weighted model checking finds the weight summary of $paths(c, c')$, which is the summation $\bigoplus_{\sigma \in paths(c, c')} v(\sigma)$.

There are two kinds of generalized reachability problems:

Definition 20 Let $\mathcal{W} = (\mathcal{P}, \mathcal{S}, f)$ be a weighted transition system with $\mathcal{P} = (P, \Delta, s_0)$. Let $C \subseteq P$ and $c \in P$.

- The **generalized predecessor problem** is to find $\delta(c) = \bigoplus \{v(\sigma) \mid \sigma \in path(c, c'), c' \in C\}$.
- The **generalized successor problem** is to find $\delta(c) = \bigoplus \{v(\sigma) \mid \sigma \in path(c', c), c' \in C\}$

If a cycle exists in a weighted model, $paths(c, c')$ becomes infinite. For the termination of a weighted model checking, an idempotent semiring needs to be *bounded*.

Definition 21 An idempotent semiring is **bounded** if there are no infinite descending chains wrt \sqsubseteq , where $a \sqsubseteq b$ if, and only if, $a \oplus b = a$.

4.2 Weight Domain for ORE Analysis

For an ORE problem, we abstract a concrete environment as an abstract environment by using intervals.

Definition 22 Let Var be the set of all variables of the program. An **abstract environment** at a program location is the set of functions $AbsEnv = \{Var \rightarrow \Phi_{\perp}^k\}$, where $k = |Var|$ and $\Phi_{\perp} \in \{\bar{\Phi}_{\perp}, \check{\Phi}_{\perp}, \hat{\Phi}_{\perp}\}$. We define the *zero environment* $e_0 \in AbsEnv$ by $e_0(x) = \perp$ for $x \in Var$. Let $e, e' \in AbsEnv$, and **environment meet operation** is defined below:

$$e \sqcup e' = \lambda x. e(x) \sqcup e'(x)$$

where $\sqcup \in \{\square, \check{\square}, \hat{\square}\}$.

Weight design

The standard definition of a weight domain has the base set of weights $D = AbsEnv \rightarrow AbsEnv$. We then theoretically define the weight domain for D as follows:

Definition 23 The weight domain (bounded idempotent semiring) $\mathcal{S} = (D, \oplus, \otimes, \mathbf{0}, \mathbf{1})$ with

$$D = AbsEnv \rightarrow AbsEnv,$$

$$\mathbf{1} = \lambda x. x,$$

$$\mathbf{0} = \lambda x. e_0,$$

$$w_1 \oplus w_2 = \begin{cases} \lambda x. w_1(x) \sqcup w_2(x) & \text{if } w_1, w_2 \neq \mathbf{0} \\ w_1 & \text{if } w_2 = \mathbf{0} \\ w_2 & \text{if } w_1 = \mathbf{0} \end{cases}$$

$$w_1 \otimes w_2 = \begin{cases} w_2 \cdot w_1 & \text{if } w_1, w_2 \neq \mathbf{0} \\ \mathbf{0} & \text{otherwise} \end{cases}$$

where $\sqcup \in \{\sqcup, \ddot{\sqcup}, \hat{\sqcup}\}$.

However, this does not satisfy the descending chain condition (boundedness), since intervals are infinitely many (thus the abstract domain is infinite). To cope with this problem, we:

- restrict the models to be acyclic,
- fix an initial abstract environment I , and
- generate weight on-the-fly.

In the context of our ORE analysis, the intuition behind the first two is,

- a target program has bounded loops only; thus after unfolding loops, abstraction produces an acyclic transition system, and
- the result of ORE analysis depends heavily on the input value; we will set a possible range of inputs at the program entry in advance.

On-the-fly weight generation

We first introduce the augmented weight domain to associate an input abstract environment to each weight. “ $_$ ” means any input.

Definition 24 The *augmented weight domain* $\mathcal{S}^+ = (D^+, \oplus, \otimes, \mathbf{0}^+, \mathbf{1}^+)$ consists of $D^+ = \{(W, w) \mid W \in \text{AbsEnv}, w \in D\}$, $\mathbf{0}^+ = (-, \mathbf{0})$, $\mathbf{1}^+ = (-, \mathbf{1})$, and

$$w_1^+ \oplus w_2^+ = \begin{cases} (W_1, w_1 \oplus w_2) & \text{if } W_1 = W_2 \\ \mathbf{0}^+ & \text{otherwise} \end{cases}$$

$$w_1^+ \otimes w_2^+ = \begin{cases} (W_2, w_1 \otimes w_2) & \text{if } W_1 = w_2(W_2) \\ \mathbf{0}^+ & \text{otherwise} \end{cases}$$

for $w_1^+ = (W_1, w_1), w_2^+ = (W_2, w_2) \in D^+$.

Now we are ready to define the on-the-fly weight domain $\mathcal{S}_{\mathcal{P}, I}^+$ for a transition system \mathcal{P} and $I \in \text{AbsEnv}$. The intuition is, starting from the initial abstract environment I , only reachable instances of weights are computed in on-the-fly manner.

Definition 25 For a transition system \mathcal{P} and $I \in \text{AbsEnv}$, the weight domain $\mathcal{S}_{\mathcal{P}, I}^+ = (D_{\mathcal{P}, I}^+, \oplus, \otimes, \mathbf{0}, \mathbf{1})$ is a sub semiring of \mathcal{S}^+ with $D_{\mathcal{P}, I}^+ \subseteq D^+$. $D_{\mathcal{P}, I}^+$ is given by

$$\left\{ (W, w) \mid \begin{array}{l} \exists \sigma, \sigma' \in \Delta^* \exists c, c' \in P. s_0 \Rightarrow^\sigma c \Rightarrow^{\sigma'} c' \\ \wedge W = v(\sigma)(I) \wedge w = v(\sigma') \end{array} \right\}$$

In implementation, we will identify $D^+ \subseteq \text{AbsEnv} \times D$ with $D^+ \subseteq \text{AbsEnv} \times \text{AbsEnv}$ by

$$(W, w) \equiv (W, w(W))$$

for $W \in \text{AbsEnv}, w \in D$.

4.3 ORE Analysis

The ORE analysis problem will be solved as weighted model checking on acyclic models by the following steps (Fig. 5).

1. As preprocessing, translate a C program into CIL (three address code language). Then, each loop are unfolded and each array is replaced with a set of variables (as many as its length). We obtain an acyclic program without arrays.
2. Generate weighted transition system, which is a control flow graph with an associated ORE arithmetics operation corresponding to a CIL instruction. ORE arithmetics is prepared for three types (CI, AI, EAI).
3. Apply weighted model checking. During model checking, weights are generated by an on-the-fly manner from given initial ranges of input parameters.

4.3.1 Abstract domain for ORE problem

Abstract domain

The abstract value of a variable aims to cover all of its possible values at one program location. For the ORE problem, the abstract value is a pair of fixed point and roundoff error ranges. We will show three kinds of abstractions based on CI, AI, and EAI range representations.

Definition 26 Let fxp and rdf be corresponding range representations of fixed point and roundoff error.

CI abstract domain $\bar{\Phi} = \{(fxp, rdf) \mid fxp, rdf \in \bar{R}\}$

AI abstract domain $\hat{\Phi} = \{(fxp, rdf) \mid fxp, rdf \in \hat{R}\}$

EAI abstract domain $\hat{\Phi} = \{(fxp, rdf) \mid fxp, rdf \in \hat{R}\}$

For a fresh symbol \perp (which stands for *undefined* or *uninitialized*), we define $\Phi_{\perp} = \Phi \cup \{\perp\}$.

Abstract arithmetic

Abstract arithmetic aims to propagate both fixed point ranges and roundoff error ranges of variables.

Definition 27 Replacing $(x_f, x_r), (y_f, y_r), \boxtimes$, and ϵ in the definition of ORE arithmetic (Definition 8) with

- $(\bar{x}_f, \bar{x}_r), (\bar{y}_f, \bar{y}_r), \bar{\boxtimes} = \{\bar{\boxplus}, \bar{\boxminus}, \bar{\boxtimes}, \bar{\boxdiv}\}$, and $\bar{\epsilon}$, we obtain **CI abstract arithmetic**,
- $(\hat{x}_f, \hat{x}_r), (\hat{y}_f, \hat{y}_r), \hat{\boxtimes} = \{\hat{\boxplus}, \hat{\boxminus}, \hat{\boxtimes}, \hat{\boxdiv}\}$, and $\hat{\epsilon}$, we obtain **AI abstract arithmetic**, and
- $(\hat{x}_f, \hat{x}_r), (\hat{y}_f, \hat{y}_r), \hat{\boxtimes} = \{\hat{\boxplus}, \hat{\boxminus}, \hat{\boxtimes}, \hat{\boxdiv}\}$, and $\hat{\epsilon}$, we obtain **EAI abstract arithmetic**,

where

$$\begin{cases} \bar{\epsilon} = \hat{\epsilon} = [b^{-fp}/2, b^{-fp}/2] \\ \hat{\epsilon} = (b^{-fp}/2)\epsilon_r \text{ with a fresh noise symbol } \epsilon_r \end{cases}$$

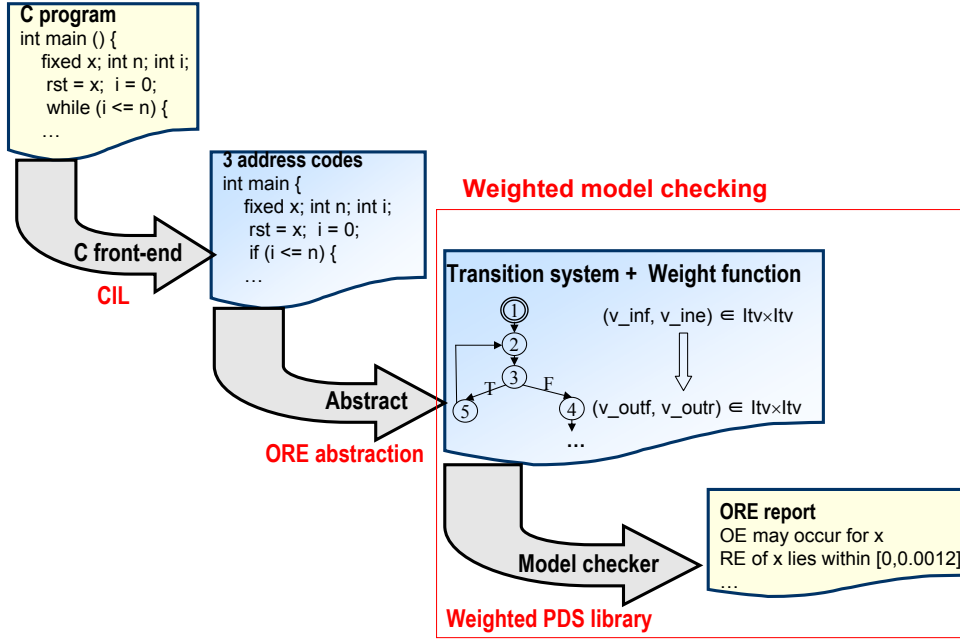


Fig. 5 ORE analysis as weighted model checking

To illustrate how to create EAI abstract numbers and compute EAI abstract arithmetic, we consider the following example:

Example 15 The EAI abstract numbers $(\hat{x}_f, \hat{x}_r), (\hat{y}_f, \hat{y}_r)$ for variables x, y , respectively, in Example 1 are

$$\begin{cases} \hat{x}_f = [1, 1] \overline{+} [2, 2]\varepsilon_1 \\ \hat{x}_r = [2^{-5}, 2^{-5}]\varepsilon_3 \\ \hat{y}_f = [0, 0] \overline{+} [10, 10]\varepsilon_2 \\ \hat{y}_r = [2^{-5}, 2^{-5}]\varepsilon_4 \end{cases}$$

since $fp = 5$, the REs of x and y are in $[-2^{-fp-1}, 2^{-fp-1}] = [-2^{-5}, 2^{-5}]$ ($= [-0.03125, 0.03125]$) and $\delta_x = 2^{-6} - 2^{-11} + 2^{-43} \approx 2^{-6}$.

Then, at the line 2 ($rst = x * x$), $(\widehat{rst}_f, \widehat{rst}_r) = (\hat{x}_f, \hat{x}_r) \hat{\boxtimes} (\hat{x}_f, \hat{x}_r)$ is projected as follows:

$$\begin{aligned} \widehat{rst}_f &= \hat{x}_f \hat{\times} \hat{x}_f \\ &= ([1, 1] \overline{+} [2, 2]\varepsilon_1) \hat{\times} ([1, 1] \overline{+} [2, 2]\varepsilon_1) \\ &= [1, 1] \overline{+} [0, 8]\varepsilon_1 \end{aligned}$$

$$\begin{aligned} \widehat{rst}_r &= 2 \hat{\times} \hat{x}_f \hat{\times} \hat{x}_r \hat{+} \hat{x}_r \hat{\times} \hat{x}_r \hat{+} \delta_x \\ &= [-0.031250, 0.031250] \overline{+} [-0.123091, 0.123091]\varepsilon_1 \\ &\quad \overline{+} [0.059615, 0.065385]\varepsilon_2 \end{aligned}$$

Abstract comparison operations

Instead of nondeterministic transitions at a conditional branch, the conditional expression can often be evaluated by using abstract environment. This is useful in avoiding unnecessary execution paths. The abstract comparison operations are defined by using ORE comparisons as follows:

Definition 28 Replacing $(x_f, x_r), (y_f, y_r)$, in the definition of \tilde{x} in ORE comparisons (Definition 6) with

- $(\bar{x}_f, \bar{x}_r), (\bar{y}_f, \bar{y}_r)$, we obtain **CI abstract comparison operations** $\bar{\Delta} = \{\bar{\leq}, \bar{=}\}$,
- $(\tilde{x}_f, \tilde{x}_r), (\tilde{y}_f, \tilde{y}_r)$, we obtain **AI abstract comparison operations** $\tilde{\Delta} = \{\tilde{\leq}, \tilde{=}\}$, and
- $(\hat{x}_f, \hat{x}_r), (\hat{y}_f, \hat{y}_r)$, we obtain **EAI abstract comparison operations** $\hat{\Delta} = \{\hat{\leq}, \hat{=}\}$.

The following example illustrates how to evaluate EAI abstract comparison $\hat{\leq}$.

Example 16 Use (\hat{x}_f, \hat{x}_r) , and (\hat{y}_f, \hat{y}_r) as in Example 17. $(\hat{x}_f, \hat{x}_r) \hat{\leq} 0$ is evaluated as follows:

$$\begin{aligned} -\bar{x}_f &= [-1, 3] \\ -\bar{x}_r &= [-2^{-5}, 2^{-5}] \end{aligned}$$

Since $(\bar{x}_f, \bar{x}_r) < 0$ is **unknown**, we can conclude that $(\hat{x}_f, \hat{x}_r) \hat{\leq} (\hat{y}_f, \hat{y}_r)$ is **unknown**.

Meet operation

At the meet of two paths in a program, we need to combine the results that are generated from these paths. The result of the meet must bind all input abstract values. We first consider how to compute the union of two ranges:

Definition 29 The unions of ranges are:

- *CI*: $[x_l, x_h] \overline{\cup} [y_l, y_h] = [\min(x_l, y_l), \max(x_h, y_h)]$.

- *AI*: $(u_0 + \sum_{i=1}^n u_i \varepsilon_i) \dot{\cup} (v_0 + \sum_{i=1}^n v_i \varepsilon_i) = (\frac{u_0 + v_0}{2} + \frac{|u_0 - v_0|}{2} \varepsilon_{n+1} + \sum_{i=1}^n t_i \varepsilon_i)$ where $\varepsilon_{n+1} \in [-1, 1]$ is a new noise symbol and, for each i ,

$$t_i = \begin{cases} u_i & \text{if } |u_i| > |v_i|, \\ v_i & \text{otherwise.} \end{cases}$$

- *EAI*: $(\bar{u}_0 \bar{+} \sum_{i=1}^n \bar{u}_i \varepsilon_i) \hat{\cup} (\bar{v}_0 \bar{+} \sum_{i=1}^n \bar{v}_i \varepsilon_i) = (\bar{u}_0 \bar{\cup} \bar{v}_0) \bar{+} \sum_{i=1}^n (\bar{u}_i \bar{\cup} \bar{v}_i) \varepsilon_i$.

Then, the result of meet operation is a pair of the union of fixed point ranges and the union of roundoff error ranges.

Definition 30 The meets in abstract values are:

- *CI meet*: $(\bar{x}_f, \bar{x}_r) \sqcap (\bar{y}_f, \bar{y}_r) = (\bar{x}_f \bar{\cup} \bar{y}_f, \bar{x}_r \bar{\cup} \bar{y}_r)$
- *AI meet*: $(\ddot{x}_f, \ddot{x}_r) \dot{\cup} (\ddot{y}_f, \ddot{y}_r) = (\ddot{x}_f \dot{\cup} \ddot{y}_f, \ddot{x}_r \dot{\cup} \ddot{y}_r)$
- *EAI meet*: $(\hat{x}_f, \hat{x}_r) \hat{\cup} (\hat{y}_f, \hat{y}_r) = (\hat{x}_f \hat{\cup} \hat{y}_f, \hat{x}_r \hat{\cup} \hat{y}_r)$

$\sqcap \in \{\bar{\cup}, \dot{\cup}, \hat{\cup}\}$ is extended to $\Phi_{\perp} \in \{\bar{\Phi}_{\perp}, \dot{\Phi}_{\perp}, \hat{\Phi}_{\perp}\}$ by $\perp \sqcap (x_f, x_r) = (x_f, x_r) \perp = (x_f, x_r)$.

4.3.2 Weighted transition system for ORE analysis

In preprocessing phase, C programs are transformed into CIL, a three address code language. We replace each array with fixed length by a set of variables that correspond to locations in an array. Thus, our target instructions are restricted as follows.

- **Assignment**: “ $x = y \circ z$ ” with $\circ \in \{+, -, *, /\}$.
- **Conditional instruction**: “if $x \circ y$ then s ” where s is an instruction and $\circ \in \{<, <=, >, >=, =, !=\}$. If the condition $(x \circ y)$ is false, s is not visited; otherwise, s is visited.
- **Control instruction**: “return loc ”, “goto loc ”, “break”, “continue”. Control moves to the specified location, and the values of variables do not change.
- **While Loop**: “while $x \circ y$ { $body$ }” with $\circ \in \{<, <=, >, >=, =, !=\}$. $body$ is repeated as long as the condition $(x \circ y)$ holds. Inside $body$, “break” will exit from the loop.

In preprocessing phase, the bounded while loops are unfolded as a sequence of conditional instructions, and we obtain an *acyclic* CIL program.

The weight function is defined as follows:

Definition 31 For an acyclic transition system \mathcal{P} and $I \in AbsEnv$, the weight function $f_{\mathcal{P}, I} : \Delta \rightarrow D_{\mathcal{P}, I}^+$ is given in Table 2.

Then, we obtain the weighted transition system:

$$\mathcal{W} = (\mathcal{P}, \mathcal{S}_{\mathcal{P}, I}^+, f_{\mathcal{P}, I}).$$

We explain how to generate a weighted transition system by Example 1 (in Introduction).

instruction	weight
“ $x = y \circ z$ ”	$(W_i, \{x_o = y_i \boxdot z_i, v_o = v_i v \in Var \setminus \{x\}\})$ where \boxdot is the corresponding abstract arithmetic operation of \circ
“if $x \circ y$ then s ”	$\mathbf{0}^+$ if $x_i \boxdot y_i = \text{false}$; $\mathbf{1}^+$ otherwise, where \boxdot is the corresponding abstract comparison of \circ
Control inst.	$\mathbf{1}^+$

Table 2 Weight function of ORE analysis

```

maintest{
st1.  if (x > 0) {
st2.    rst = x * x;
      else {
st3.    __cil_tmp1 = (Real )3;
          rst = __cil_tmp1 * x;
st4.    rst -= y;
st5.    return (rst);
st6.  }
}

```

Fig. 6 CIL code for Example 1

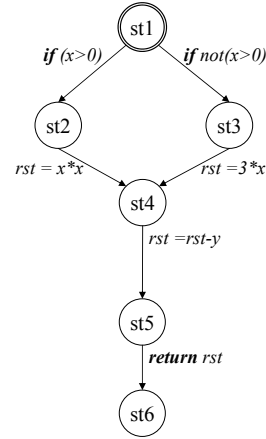


Fig. 7 CFG of three address codes in Fig. 6

Example 17 We use EAI representation type.

The CIL codes of the C program in Fig. 2 are shown in Fig. 6. Let $st1, \dots, st6$ be its locations.

To distinguish variables at each locations, we will denote a variable v at the location sti by $v^{(i)}$. The fixed-point value and RE of v are denoted by $\hat{v}_f^{(i)}, \hat{v}_r^{(i)}$ respectively.

The transition system is $\mathcal{P} = (P, \Delta)$, where $P = \{st1, st2, \dots, st6\}$, Δ is shown in Fig. 7 and f is defined in Table 3.

The initial abstract environment W_{init} at $st1$ is generated from initial range values of variables (given in the topmost comments in Fig. 2), in that:

$$\begin{cases} x^{(init)} = ([1, 1] \bar{+} [2, 2] \varepsilon_1, [2^{-5}, 2^{-5}] \varepsilon_3) \\ y^{(init)} = ([10, 10] \varepsilon_2, [2^{-5}, 2^{-5}] \varepsilon_4) \end{cases}$$

Then, the resulting weighted transition system is $\mathcal{W} = (\mathcal{P}, \mathcal{S}_{\mathcal{P}, W_{init}, f}^+, f)$.

transition	weight
(st1,st2)	$if ((x^{(init)} \succ 0) = false) \mathbf{0}^+ \text{ else } \mathbf{1}^+$
(st1,st3)	$if \text{ not}((x^{(init)} \succ 0) = false) \mathbf{0}^+ \text{ else } \mathbf{1}^+$
(st2,st4)	$(W_i^{(in)}, \{rst^{(2)} = x^{(in)} \hat{\boxtimes} x^{(in)}, v^{(2)} = v^{(in)} \mid v \in Var \setminus \{rst\}\})$
(st3,st4)	$(W_i^{(in)}, \{-cil_tmp1^{(3)} = 3, rst^{(3)} = _cil_tmp1^{(3)} \hat{\boxtimes} x^{(in)}, v^{(3)} = v^{(in)} \mid v \in Var \setminus \{-cil_tmp1, rst\}\})$
(st4,st5)	$(W_i^{(in)}, \{rst^{(4)} = rst^{(in)} \hat{\boxminus} y^{(in)}, v^{(4)} = v^{(in)} \mid v \in Var \setminus \{rst\}\})$
(st5,st6)	$\mathbf{1}^+$

Table 3 Weight function for a CIL code in Example 17

Since the abstraction is an overapproximation, we conclude soundness of ORE analysis.

Theorem 2 *An ORE analysis on a subclass of acyclic CIL programs is sound.*

4.3.3 ORE Analysis Example

We continue to explain how ORE analysis works by the example in Fig. 2.

Example 18 Let the input EAI $x^{(1)} = (\hat{x}_f^{(1)}, \hat{x}_r^{(1)})$ and $y^{(1)} = (\hat{y}_f^{(1)}, \hat{y}_r^{(1)})$ as shown in Example 17.

St1 is a conditional branch. Since the initial range of x is $[-1, 3]$, CANA cannot decide the condition $x > 0$. Thus, it traces both *st2* and *st3*, and later merges their results.

At *st2*, the RE of **rst** is computed by multiplication $\hat{\boxtimes}$ as:

$$\widehat{rst}_r^{(2)} = [-0.031250, 0.031250] \mp [-0.123091, 0.123091]_{\varepsilon_1} \mp [0.059615, 0.065385]_{\varepsilon_3}$$

and at *st3*,

$$\widehat{rst}_r^{(3)} = 3 \hat{\times} \hat{x}_r = [0.093750, 0.093750]_{\varepsilon_3}$$

And they are merged as:

$$\begin{aligned} \widehat{rst}_r &= [-0.031250, 0.031250] \mp [-0.123091, 0.123091]_{\varepsilon_1} \mp \\ &\quad ([0.059615, 0.065385] \cup [0.093750, 0.093750])_{\varepsilon_3} \\ &= [-0.031250, 0.031250] \mp [-0.123091, 0.123091]_{\varepsilon_1} \\ &\quad \mp ([0.059615, 0.093750])_{\varepsilon_3} \end{aligned}$$

At *st4*, by subtraction $\hat{\boxminus}$, we get the RE of **rst**:

$$\widehat{rst}_r^{(4)} = [-0.031250, 0.031250] \mp [-0.123091, 0.123091]_{\varepsilon_1} \mp [0.059615, 0.09375]_{\varepsilon_3} \mp [-0.031250, -0.031250]_{\varepsilon_4}$$

RE bound of **rst** is by $[-0.279341, 0.279341]$ by replacing each ε_i with $[-1, 1]$ in \hat{r} .

By nature of analysis, this analysis overapproximates REs. It occurs at the conditional branch (line 1) and the multiplication. For instance, at *st2*, $\hat{\delta}$ (in $\hat{\boxtimes}$) is approximated with $[-0.031250, 0.031250]$.

The example above shows that the overapproximations will occur at (1) nonlinear operations; (2) undecided conditional branch. In case the conditional branch is decided, analysis will not overapproximate it. Below, how the treatment of conditional branches in ORE analysis will improve the overapproximation.

Example 19 In Example 18, if we reduce the initial range of x_f to $[1, 3]$, the condition $x > 0$ is decided to be true at *st1*, and CANA ensures that *st3* will not be executed. Then, at *st4*, by subtraction $\hat{\boxminus}$,

$$\begin{aligned} \hat{r} &= [-0.031250, 0.031250] \mp [-0.123091, 0.123091]_{\varepsilon_1} \\ &\quad \mp [0.059615, 0.065385]_{\varepsilon_3} \\ &\quad \mp [-0.031250, -0.031250]_{\varepsilon_4} \end{aligned}$$

RE of **rst** is bounded by $\bar{r} = [-0.250976, 0.250976]$ by replacing each ε_i with $[-1, 1]$ in \hat{r} .

This result is more precise than that in Example 18.

4.4 Experiments

Implementation

We have implemented our analysis framework in a tool *C ANalyzer* (CANA). CANA uses two libraries: CIL library³ and WPDS library⁴.

- CIL (C Intermediate Language) is a high-level representation that permit source-to-source transformation of C programs. CIL is used to generate three address codes, information about variables, and the CFG of a C program.
- WPDS (Weighted Pushdown System) is a library, which provides functions to the sets of forward- or backward- reachable configurations in a weighted pushdown system. Since we exclude procedure calls and unbounded loops, we adopt WPDS only for weighted finite acyclic state transition systems

The inputs of CANA are subclass of ANSI C programs and initial ranges of variables. The outputs of CANA are roundoff error ranges of variables at each point of the program, and warning about overflow errors (if they occur). CANA has six main modules (Fig. 5.4) as follows:

³ <http://hal.cs.berkeley.edu/cil/>

⁴ <http://www.fmi.uni-stuttgart.de/szs/tools/wpds/>

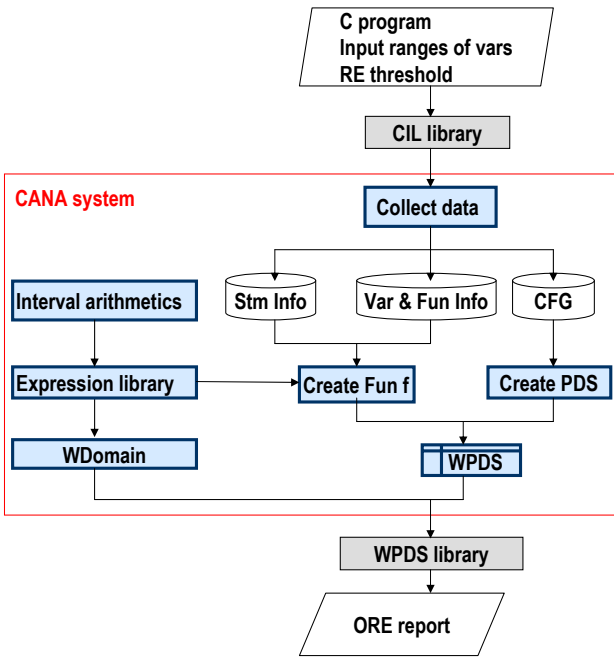


Fig. 8 CANA system

1. *Collect data* module generates information required for the analysis, including: statement information (Stm Info), (2) variable and function information (Var and Func Info), and (3) CFG of C program.
2. *Range arithmetics* module includes three types of interval arithmetics: CI arithmetic, AI arithmetic, and EAI arithmetic.
3. *Evaluate exps* module evaluates the abstract values of expressions based on types of range arithmetics.
4. *Create PDS* module generates transition system from control flow graph of a C program.
5. *Create Fun f* module assigns a weight to each transition.
6. *WDomain* module defines two operations: \otimes and \oplus .

Experimental results

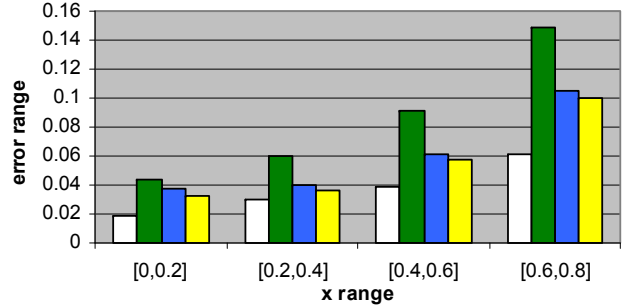
We have implemented in CANA three types of interval representations: CI, AI, and EAI. CANA can analyze programs that have nested loops 64×64 .

In order to compare the efficiency of EAI arithmetic to CI and AI arithmetics, we analyzed source codes of three examples:

1. a program computes a polynomial of degree 5
2. a program computes the sine function
3. a tiny fragment, which frequently appears in the mpeg decoder reference algorithm.

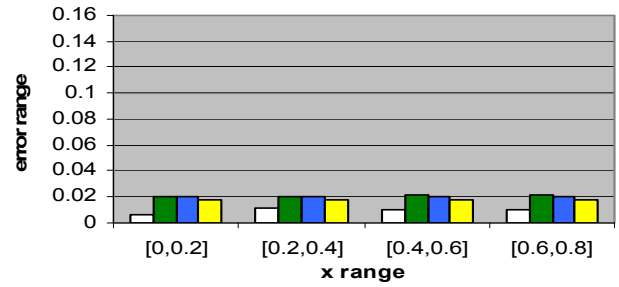
Fig. 9 shows experimental results of analyzing these programs (on PC with Intel(R) Core(TM) Duo CPU

[x]	[0,0.2]	[0.2,0.4]	[0.4,0.6]	[0.6,0.8]
real err	0.01909	0.03	0.03891	0.061588
CI	0.04373	0.0599	0.09172	0.148848
AI	0.0377	0.0406	0.06179	0.104468
EAI	0.03232	0.0356	0.05763	0.100454



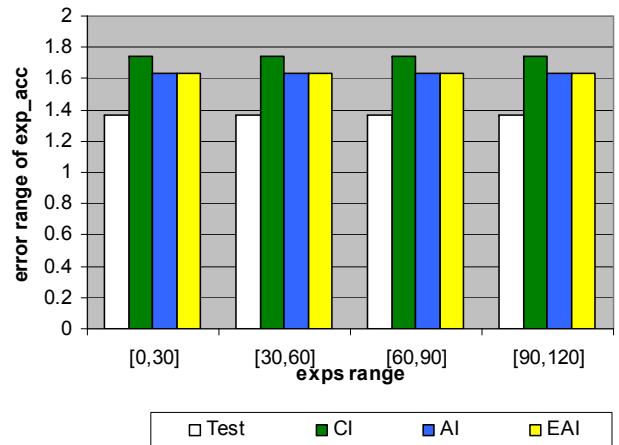
a. The analysis result of P5(x)

[x]	[0,0.2]	[0.2,0.4]	[0.4,0.6]	[0.6,0.8]
real err	0.00647	0.0108	0.01049	0.010206
CI	0.0204	0.0208	0.02139	0.022198
AI	0.02035	0.0204	0.02029	0.019976
EAI	0.01759	0.0176	0.01749	0.017184



b. The analysis result of Sin(x)

[exps]	[0,30]	[30,60]	[60,90]	[90,120]
real err	1.36309	1.3631	1.36309	1.363086
CI	1.73766	1.7377	1.73766	1.737656
AI	1.63453	1.6345	1.63453	1.634531
EAI	1.63453	1.6345	1.63453	1.634531



c. The analysis result of pMpeg(exps)

Fig. 9 The experimental results

1.66GHz, 1.5Gb of memory). Fig. 9a shows results of analyzing the program that computes $P5(x) = 1 - x + 3x^2 - 2x^3 + x^4 - 5x^5$, where the fraction part $fp = 8$. The *true err* row is the width of roundoff error ranges which are found by testing; the CI, AI, and EAI rows are the widths of roundoff error ranges computed by CI arithmetic, AI arithmetic, and EAI arithmetic, respectively. Our experiments show that EAI is more precise than CI and is comparable to AI. We got similar results for the program that computes sine of x shown in Fig. 9b, and the fragment of the mpeg decoder $pMpeg(exps)$ in Fig. 9c. All tests run in less than 2 seconds.

5 Detecting REs based on Counterexample-guided Narrowing

Static analysis is useful in proving safety properties of ORE problem. But it requires overapproximation in both propagating REs and control flows like conditional branches. Hence, it may return spurious counterexamples. Fortunately, in our setting, the floating to fixed-point conversion, we can compute the exact RE, whereas there are in general no ways to compute exact real numbers. Though testing can return exact REs, it cannot cover all possible inputs. If we are lucky, witness of large REs would be eventually found, yet most of them may be missed. Another challenging problem is how to reduce the number of test cases if the input domain is large and the input parameters are many.

A popular approach to deal with spurious counterexamples is *counterexample-guided abstraction refinement* (CEGAR) [6]. Inspired by CEGAR, we combine testing and RE analysis.

This section proposes an approach for detecting REs of C programs, which combines static analysis and testing, and make them refine each other. We call this combination *counterexample-guided narrowing*. First, we apply an overflow and roundoff errors analysis from our previous work [44] which returns an overapproximation of REs as an Extended affine interval (EAI). Fortunately, an EAI represents the relations between the input value and the RE of the output. These relations can be used to clarify: variables are irrelevant to REs of the results, variables affect the REs the most, and the ranges of inputs are most likely to cause the maximum RE. These observations effectively narrow the focus of test data generation. Second, in case testing does not find a witness of RE violation, the analysis may overapproximate too much. Further, the narrower the input ranges are, the more precise the analysis result will be. Therefore, with a “divide and conquer” refinement strategy, we can check the most suspicious part first.

Throughout the section, we focus only on roundoff errors. We assume that ORE analyzer (CANA) does not detect any overflow errors.

5.1 Observation on RE Analysis

The inputs of an RE analysis consist of

- a C program (with m -input variables) to be analyzed (base $b = 2$),
- a fixed-point format (sp, ip, fp) ,
- an RE threshold $\theta (> 0)$, and
- a pair of a fixed-point range $[l_i, h_i]$ and an RE range $[l_{m+i}, h_{m+i}]$ with $-2^{-fp-1} < l_{m+i} \leq h_{m+i} < 2^{-fp-1}$ for each i -th input variable.

We will fix the last three elements as an environment of the RE analysis. We call the Cartesian product $D = [l_1, h_1] \times \dots \times [l_{2m}, h_{2m}]$ an *input domain*.

Throughout the RE analysis, all ranges are represented as EAIs with $2m$ noise symbols (where ε_i and ε_{m+i} correspond to noise symbols of the fixed-point part and the RE of values of the i -th input variable). Thus, we coerce input CIs to EAIs by EAI coercion. In the context of the RE analysis, we denote: Input domain $D = [l_1, h_1] \times \dots \times [l_{2m}, h_{2m}]$ is $(\hat{v}_1, \dots, \hat{v}_{2m})$ with $\hat{v}_i = (\frac{l_i+h_i}{2}) \mp (\frac{h_i-l_i}{2})\varepsilon_i$.

As notational convention, the analysis result is denoted by an EAI:

$$\hat{r} = \bar{r}_0 \mp \sum_{i=1}^{2m} \bar{r}_i \varepsilon_i$$

The analysis result \hat{r} shows extra information about the effects of inputs on \hat{r} , since EAI coercions of input ranges and \hat{r} share common noise symbols. If violations are found (i.e., EAI projection of \hat{r} exceeds $[-\theta, \theta]$), we need to check whether they are spurious by testing. Fortunately, we have useful observations below, which will optimize test data generation and testing.

- **RE bound for each test case:** Assume that the valuation of noise symbols for the test case $t = (t_1, \dots, t_{2m})$ is $(\lambda_1, \dots, \lambda_{2m})$, i.e.,

$$\lambda_i = \begin{cases} 0 & \text{if } l_i = h_i \\ \frac{2t_i - (l_i + h_i)}{(h_i - l_i)} & \text{otherwise} \end{cases}$$

Then, the RE for the input t is bounded by the valuation of \hat{r} with $(\lambda_1, \dots, \lambda_{2m})$.

- **Irrelevant noise symbol:** If the coefficient of a noise symbol ε_k is $\bar{r}_k = [0, 0]$, the noise symbol ε_k will not affect the RE of result.
- **Sensitivity of noise symbols:** If $|\bar{r}_k| \leq |\bar{r}_h|$, the noise symbol ε_k affects \hat{r} more than ε_h .

The following example will demonstrate how they affect the testing phases:

Example 20 In the analysis result of Example 18:

- For the test case $t = (x_f, y_f, x_r, y_r) = (1, 5, 0, 0)$, the valuation of noise symbols is $(\varepsilon_1, \varepsilon_2, \varepsilon_3, \varepsilon_4) = (0, 0.5, 0, 0)$. The RE bound of t is

$$[-0.031250, 0.031250] \subseteq [-0.26, 0.26]$$

Therefore, t is not a counterexample.

- Since $\bar{r}_2 = [0, 0]$, ε_2 is an *irrelevant noise symbol*. Hence, v_2 (or fixed-point part of y) does not affect RE of rst .
- We have $|\bar{r}_1| = 0.123091 = \max\{|\bar{r}_1|, \dots, |\bar{r}_4|\}$, thus ε_1 is the most sensitive noise symbol. Hence, v_1 (or fixed-point part of x) affect the RE of rst the most.

The analysis result is helpful to optimize test phase, includes:

- **Pre-evaluate test case:** if RE bound of a test case t lies in the RE threshold bound, we need not execute test for t .
- **Reduce input ranges:** which reduces segments in each input range if corresponding REs are subsumed. Especially, for an irrelevant noise symbol.
- **Choice of the number of ticks,** which takes more ticks in the input ranges of more sensitive noise symbols.

5.2 Refining Test Data Generation

5.2.1 Test Data Generation

For an input domain $D = [l_1, h_1] \times \dots \times [l_{2m}, h_{2m}]$, a basic strategy of test data generation is to divide the input domain into meshes and select one test case from each mesh.

Definition 32 For an interval $[l, h]$, and $k > 1$, the k -ticks of $[l, h]$ starting from c are $\{c, c + \Delta, \dots, c + (k - 1)\Delta\}$ where $\Delta = \frac{h-l}{k}$ and c is a number lies in $[l, l + \Delta]$.

Remark 3 This c intends a randomly generated flicker; since refinement loops will repeat, we would like to avoid generating the same test cases.

Example 21 Let us consider the C program as Fig. 2.

- For $x_f \in [-1, 3]$, let $k_1 = 10$, then $\Delta_1 = (3 - (-1))/10 = 0.4$. Let $c_1 = -0.8$, we got the set of ticks $X_f = \{-0.8, -0.4, \dots, 3.8\}$
- For $x_r \in [-2^{-5}, 2^{-5}]$, let $k_2 = 10$, then $\Delta_2 = (2^{-5} - (-2^{-5}))/10 = 0.00625$. Let $c_2 = -0.03$, we got the set of ticks $X_r = \{-0.03, -0.02375, \dots, 0.0265\}$

- For $y_f \in [-10, 10]$, let $k_3 = 10$, then $\Delta_3 = (10 - (-10))/10 = 2$. Let $c_3 = -9$, we got the set of ticks $Y_f = \{-9, -7, \dots, 9\}$
- For $y_r \in [-2^{-5}, 2^{-5}]$, let $k_4 = 10$, then $\Delta_4 = (2^{-5} - (-2^{-5}))/10 = 0.00625$. Let $c_4 = -0.028$, we got the set of ticks $Y_r = \{-0.028, -0.02175, \dots, 0.0285\}$

Hence, the set of test data is $T = X_f \times Y_f \times X_r \times Y_r$. E.g., For test case $t_1 = (-0.8, -9, -0.03, -0.028)$, the input of fixed-point program is $(x, y) = (-0.8, -9)$; the input of floating-point program is $(x, y) = (-0.83, -9.028)$.

For an input domain $D = [l_1, h_1] \times \dots \times [l_{2m}, h_{2m}]$, all combinations of k_i -ticks of $[l_i, h_i]$ for $i \leq 2m$ are the set of test data. Then, we execute a program in two ways: with floating-point arithmetic and fixed-point arithmetic. The difference between them is a true RE. However, the number of test data grows with the power of the $2m$ -th degree.

5.2.2 Range Reduction

For two ranges $[l_1, h_1], [l_2, h_2]$, we denote $[l_1, h_1] \leq [l_2, h_2]$ if $u \leq v$ for each $u \in [l_1, h_1]$ and $v \in [l_2, h_2]$ (i.e., $h_1 \leq l_2$). Reducing input range is executed based on the observation as the following lemma.

Lemma 2 Assume $0 \notin [u, v]$. Then,

- $[-v, -u] \leq [u, v][-\frac{u}{v}, \frac{u}{v}] \leq [u, v]$ if $0 < u \leq v$
- $[u, v] \leq [u, v][-\frac{v}{u}, \frac{v}{u}] \leq [-v, u]$ if $u \leq v < 0$

This lemma means that if $0 \notin [u_i, v_i] = \bar{r}_i$ we can ignore test data with corresponding noise symbol ε_i in $[-\frac{u_i}{v_i}, \frac{u_i}{v_i}]$ (resp. $[-\frac{v_i}{u_i}, \frac{v_i}{u_i}]$) for $0 < u_i \leq v_i$ (resp. $u_i \leq v_i < 0$). The reason for this is that the true REs for test data with ε_i in $[-\frac{u_i}{v_i}, \frac{u_i}{v_i}]$ (resp. $[-\frac{v_i}{u_i}, \frac{v_i}{u_i}]$) are bounded by the valuations when a noise symbol ε_i is either 1 or -1 whatever a true coefficient has a value in $[u_i, v_i]$.

In a special case $u_i = v_i = 0$, a valuation of ε_i does not matter. Thus, only a valuation with 0 is considered.

From observations above we reduce the input domain D to two input domains D_{max} (which contain an input that causes the maximum RE) and D_{min} (which contain an input that causes the minimum RE) without losing opportunities to find test cases that cause violation of REs.

Definition 33 For an input domain

$D = (\frac{l_1+h_1}{2} + \frac{h_1-l_1}{2}\varepsilon_1) \times \dots \times (\frac{l_{2m}+h_{2m}}{2} + \frac{h_{2m}-l_{2m}}{2}\varepsilon_{2m})$ and the analysis result

$$\hat{r} = [u_0, v_0] \mp \sum_{i=1}^{2m} [u_i, v_i] \varepsilon_i,$$

the subdomain D_{max} of D is

$$\left(\frac{l_1 + h_1}{2} + \frac{h_1 - l_1}{2}\bar{\varepsilon}_1\right) \times \dots \times \left(\frac{l_{2m} + h_{2m}}{2} + \frac{h_{2m} - l_{2m}}{2}\bar{\varepsilon}_{2m}\right)$$

where

$$\bar{\varepsilon}_i = \begin{cases} \left[\frac{|u_i|}{|v_i|}, 1\right] & \text{if } 0 < u_i \leq v_i \\ [-1, -\frac{|v_i|}{|u_i|}] & \text{if } u_i \leq v_i < 0 \\ [0, 0] & \text{if } u_i = v_i = 0 \\ [-1, 1] & \text{otherwise} \end{cases}$$

and the subdomain D_{min} is

$$\left(\frac{l_1 + h_1}{2} + \frac{h_1 - l_1}{2}\underline{\varepsilon}_1\right) \times \dots \times \left(\frac{l_{2m} + h_{2m}}{2} + \frac{h_{2m} - l_{2m}}{2}\underline{\varepsilon}_{2m}\right)$$

where

$$\underline{\varepsilon}_i = \begin{cases} [-1, -\frac{|u_i|}{|v_i|}] & \text{if } 0 < u_i \leq v_i \\ \left[\frac{|v_i|}{|u_i|}, 1\right] & \text{if } u_i \leq v_i < 0 \\ [0, 0] & \text{if } u_i = v_i = 0 \\ [-1, 1] & \text{otherwise} \end{cases}$$

An EAI $\hat{r} = \bar{r}_0 \bar{\top} \sum_{i=1}^{2m} \bar{r}_i \varepsilon_i$ is maximum (resp. minimum) if each of elements \bar{r}_0 and $\bar{r}_i \varepsilon_i$ is maximized (resp. minimized). Thus the maximal (resp. minimal) RE will occur among valuations of ε_i in $[u/v, 1]$ (resp. $[-1, -u/v]$) if $0 < u \leq v$, and in $[-1, -v/u]$ (resp. $[v/u, 1]$) if $u \leq v < 0$. Therefore, from Lemma 2, we obtain the next theorem.

Theorem 3 *If there exists a counterexample in D , then there exists a counterexample in $D_{max} \cup D_{min}$.*

Example 22 By observation of analysis result \hat{r} in Example 18, we can find D_{max} of the input domain D is

$$([1, 1] \bar{\top} [2, 2]\bar{\varepsilon}_1) \times ([10, 10]\bar{\varepsilon}_2) \times ([2^{-5}, 2^{-5}]\bar{\varepsilon}_3) \times ([2^{-5}, 2^{-5}]\bar{\varepsilon}_4)$$

with $\bar{\varepsilon}_1 = [-1, 1]$, $\bar{\varepsilon}_2 = [0, 0]$, $\bar{\varepsilon}_3 = [0.63589, 1]$, and $\bar{\varepsilon}_4 = [-1, -1]$. Hence, D_{max} is projected to

$$[-1, 3] \times [0, 0] \times [0.019872, 0.03125] \times [-0.03125, -0.03125].$$

Hence, we can conclude that the input with $y_f = 0$ and $y_r \simeq -0.03125$ will cause the maximum RE.

5.2.3 More Ticks for more Sensitive Noise Symbols

We need a strategy to setting ticks for each initial ranges in input domain. A bad strategy of setting ticks will create several test cases which cause similar REs and all of them lie within RE threshold bound. Testing over these test cases are not needed.

Analysis result \hat{r} shows the effects of noise symbols to the REs. A larger coefficient of a noise symbol causes

stronger effect on the RE of the result. For example, the variable corresponding to dominant noise symbol will strongly affect REs, in other words, the changing of this variable causes the changing of RE the most. Therefore, setting more ticks on the initial range of this variables can lead the test cases to various REs. Based on this observation, our basic idea is, in input domains D_{min}, D_{max} , the initial range of variables which are predicated strongly affect REs will be set more ticks than other initial ranges. The strategy of setting number of ticks is then depending on coefficients of noise symbols in analysis result as follows:

Definition 34 Let $\hat{r} = \bar{r}_0 \bar{\top} \sum_{i=0}^{2m} \varepsilon_i$ be analysis result of input domain $D' = [u_1, v_1] \times \dots \times [u_{2m}, v_{2m}]$.

For $\sigma > 0$, a *tick frequency* t_i (wrt σ) for the input interval $[u_i, v_i]$ is

$$t_i = \begin{cases} \lceil \frac{2 \times \bar{r}_i}{\sigma} \rceil & \text{if } u_i < v_i \\ 1 & \text{if } u_i = v_i \end{cases}$$

Here, $\lceil x \rceil$ denotes the round up of x .

Example 23 For \hat{r} in Example 18 and D_{max} in Example 22, the tick frequency t_1, t_2, t_3 , and t_4 wrt $\sigma = 0.01$ (of $\bar{\varepsilon}_1, \bar{\varepsilon}_2, \bar{\varepsilon}_3$, and $\bar{\varepsilon}_4$, respectively) are,

$$\begin{aligned} t_1 &= \lceil \frac{2 \times 0.123091}{0.01} \rceil = 25 \\ t_2 &= 1 \\ t_3 &= \lceil \frac{2 \times 0.09375}{0.01} \rceil = 19 \\ t_4 &= 1 \end{aligned}$$

Thus, the number of test cases is $25 \times 1 \times 19 \times 1 = 475$, and the RE found by testing 475 test cases is 0.219720.

5.3 Refinement of Analysis by Narrowing Input Domains

An analysis may report spurious counterexamples. Fortunately, our RE analysis becomes more precise if an input domain becomes narrower. There are two reasons that make input domain decomposition reduces the overapproximations:

- a smaller input domain is more likely to be deterministic on conditional branches, and
- a smaller input domain is more likely makes EAI arithmetic more precise.

Our “divide and conquer” strategy has two phases:

- Reduce an input domain D to D_{max} and D_{min} (Definition 33)
- Divide the input ranges (in D_{max} and D_{min}) of the most sensitive noise symbol ε_k into two ranges.

Definition 35 Let $D_{max} = [l_1, h_1] \times \dots \times [l_{2m}, h_{2m}]$ be an input domain and ε_k be the most sensitive noise symbol. Then:

$$\begin{aligned} - D_{max}^1 &= D_{max}|_{\bar{v}_k=[l_k, \frac{l_k+h_k}{2}]} \\ - D_{max}^2 &= D_{max}|_{\bar{v}_k=[\frac{l_k+h_k}{2}, h_k]} \end{aligned}$$

where \bar{v}_k is the k -th element of D_{max}^1 (D_{max}^2).

For D_{min} , we also have a similar partition strategy.

The next round of the RE analysis will be performed for input domains D_{max}^1 and D_{max}^2 . Our early experience shows that often one of analysis results of D_{max}^1 and D_{max}^2 lie in the RE threshold bound. Thus this simple strategy becomes quite effective.

Example 24 From Example 20, the most sensitive noise symbol is ε_1 .

Form Example 22, the new input domain D_{max} is

$$[-1, 3] \times [0, 0] \times [0.019872, 0.03125] \times [-0.03125, -0.03125]$$

and ε_1 is the most sensitive symbol (Example 20).

Hence, we will divide the initial range of v_1 ($[-1, 3]$) into two new subranges $[-1, 1]$ and $[1, 3]$ and we get:

$$\begin{aligned} D_{max}^1 &= D_{max}|_{\bar{v}_1=[-1,1]} \\ &= [-1, 1] \times [0, 0] \times [0.019872, 0.03125] \times \\ &\quad [-0.03125, -0.03125] \\ D_{max}^2 &= D_{max}|_{\bar{v}_1=[1,3]} \\ &= [1, 3] \times [0, 0] \times [0.019872, 0.03125] \times \\ &\quad [-0.03125, -0.03125] \end{aligned}$$

RE analyses on two domains D_{max}^1 and D_{max}^2 report that:

- the REs of all input in D_{max}^1 lie in $[-0.22, 0.22]$
- the REs of all input in D_{max}^2 lie in $[-0.25, 0.25]$

Hence, we can conclude the REs of all input in D_{max} lie in RE threshold bound $[-0.26, 0.26]$.

Similarly, we get the REs of all input in D_{min} also lie in the RE threshold bound, and we can conclude the program is “safe”. Note that, before decomposition, it was $[-0.28, 0.28]$ (Example 18), which exceeds the RE threshold bound (Fig. 10). If we reduce RE threshold bound θ in $0.219720 < \theta < 0.22$, both D_{max}^1 and D_{max}^2 are not enough. In such a case, we will investigate subdomain which has larger RE found testing first in the later rounds.

Combining Analysis and Testing Algorithm

Algorithm 5.3 shows the algorithm combining the analysis and testing. Function $analyze(P_{fl}, D)$ analyzes the program P_{fl} with input domain D , and return the overapproximate RE in EAI form \hat{r} . Function $reduceMax(\hat{r}, D)$ reduce domain D to D_{max} . Function $reduceMin(\hat{r}, D)$

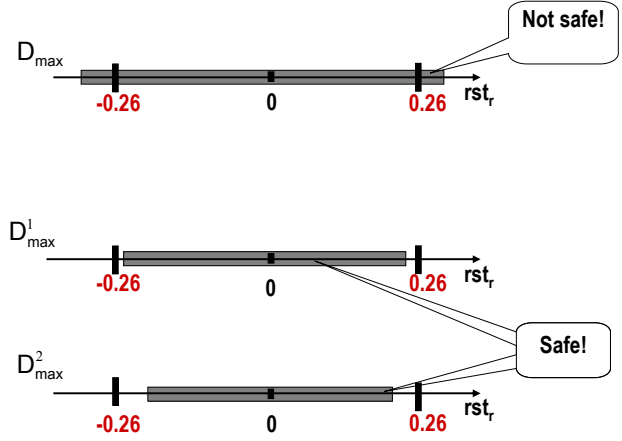


Fig. 10 Effects of decomposition of D_{max} to D_{max}^1, D_{max}^2

reduce domain D to D_{min} . Function $gentest(D, \hat{r})$ generate a set of test cases T of domain D using analysis result \hat{r} . Function $test(P_{fl}, P_{fx}, t)$ executes both two programs P_{fl}, P_{fx} with test case t and return the difference between result of these two program, called set of test results D_R . Function $divide(D, \hat{r})$ divides domain D into two new subdomains D_1, D_2 based on analysis result \hat{r} .

Input: P_{fl}, P_{fx} , initial ranges of variables, RE threshold θ

Output: Return “safe” or counterexample or unknown

Initial list of subdomains $L_D = [D]$, a set of test data

$T = \emptyset$, a set of RE $D_R = \emptyset$;

while $length(L_D) < 10$ **do**

 pop one element D from L_D ;

$\hat{r} = analyze(P_{fl}, D)$;

if ($\hat{r} \subseteq [-\theta, \theta]$) **then**

 | continue ;

end

$D_{max} = reduceMax(\hat{r}, D)$;

$D_{min} = reduceMin(\hat{r}, D)$;

 push D_{min} into L_D ;

$r_{max} = \max\{|test(P_{fl}, P_{fx}, t)| | t \in T\}$;

 Let $t_{max} \in T$ such that $r_{max} = test(P_{fl}, P_{fx}, t_{max})$;

if ($r_{max} > \theta$) **then**

 | **return** counterexample t_{max} ;

end

$\{D_1, D_2\} = divide(D_{max}, \hat{r})$;

if $t_{max} \in D_1$ **then**

 | push D_2 into L_D ;

 | push D_1 into L_D ;

else

 | push D_1 into L_D ;

 | push D_2 into L_D ;

end

end

return unknown;

Algorithm 1: Combining analysis and testing

5.4 Implementation and Experiments

CANAT implementation

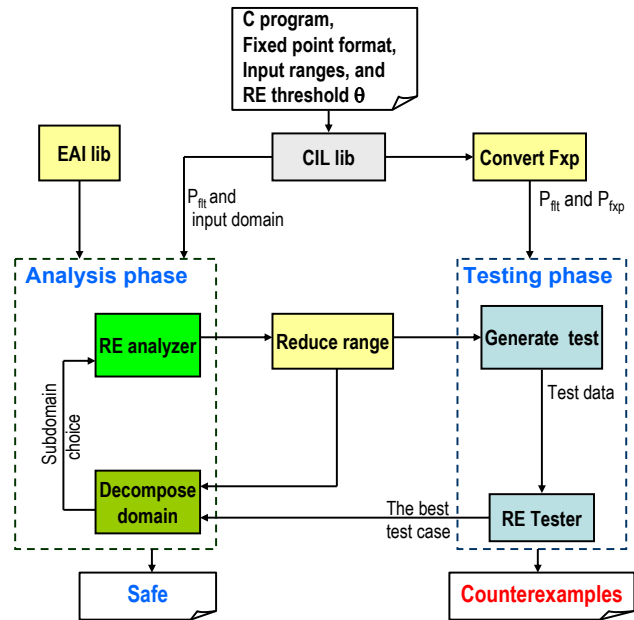


Fig. 11 CANAT system

Fig. 11 shows the construction of CANAT. CANAT uses three external tools:

- CIL library⁵ as a preprocessor,
- Weighted PDS library⁶ as a backend weighted model checking engine, and
- CANA [44] as an RE analyzer.

CANAT has 4 main modules as follows:

1. *Reduce range* module clarify: (1) the subdomains of input domain that contain maximum (or minimum) REs, (2) the choice of ticks for testing generation.
2. *Decompose domain* module divides the reduced domain (e.g., D_{max}, D_{min}) into two subdomains.
3. *Generate Test* module generates set of test data T based on information obtained from *Reduce range* module.
4. *RE Tester* module automatically generates two programs corresponding to input C program, one uses fixed-point arithmetic, while the other uses floating-point arithmetic. Then, these two programs are executed with the set of test data T . The testing REs are the differences between the results of these two programs.

⁵ <http://hal.cs.berkeley.edu/cil/>

⁶ <http://www.fmi.uni-stuttgart.de/szs/tools/wpds/>

Experimental results

Table 4 shows the results of checking 4 programs (on PC with Intel(R) Xeon(TM)CPU 3.60GHz, 3.37Gb of memory). *The first column* shows the names of 4 programs.

1. **P2** from Fig. 1: We set the initial range $[-1, 3] \times [-10, 10]$, $fp \in \{7, 8, 9, 10\}$, and $\theta \in \{0.001 + 0.002i \mid 0 \leq i \leq 9\}$.
2. **P5** that computes $P5(x) = 1 - x + 3x^2 - 2x^3 + x^4 - 5x^5$: We set the initial range $[0, 1]$, $fp \in \{7, 8, 9, 10\}$, and $\theta \in \{0.001 + 0.01i \mid 0 \leq i \leq 9\}$.
3. **Sine** that computes the sine function using Taylor expansion up to degree 21: We set the initial range $[0, 1]$, $fp \in \{7, 8, 9, 10\}$, and $\theta \in \{0.001 + 0.005i \mid 0 \leq i \leq 9\}$.
4. **pMpeg** that consists of one 8×8 loop taken from the mpeg4 decoder reference algorithm: We set the initial range $[0, 30]$, $fp \in \{7, 8, 9, 10\}$, and $\theta \in \{0.001 + 0.05i \mid 0 \leq i \leq 9\}$.

We compare experimental results by

1. ORE analysis by CANA (CANA) followed by Randomly generated testing (**Random test**), and
2. Counter example guided narrowing loop by repeating ORE analyses (**CANAT ananalysis**) and refined testing (**CANAT test**).

Both try 40 settings (i.e., different numbers of digits in fixed-point numbers, different RE thresholds) for each program. For fair comparison, we make the total numbers of test cases to be the same for each setting; **Random test** generates 200 test cases, and (**CANAT test**) generates 20 test cases for each refinement loop, which is repeated 10 times.

%Checked columns show that how many percentage (among 40 settings) is detected either to be safe or ORE violation. Although the experiment remains a toy, it shows clear improvement of testing and analysis.

6 Related Work

The ORE problems are one of the central issues in the numerical analysis [17, 13]. There are lots of works on mathematical reasoning to estimate OREs [17, 16], and there is a well-known methodology for the precise addition of floating numbers, which cancels the effect of REs [27]. It is extended to the precise multiplication [46], and recently verified numerical computation is evolving.

Our focus is more on static detections of OREs of programs, and we omit huge references of these areas, which are beyond scope of the paper.

Input program	CANA	Random test	Time (s)	%Checked	CANAT analysis	CANAT test	CANAT time(s)	%Checked of CANAT
P2	15	11	7	65.00%	20	18	13	95.00%
P5	9	15	14	60.00%	12	19	24	77.50%
Sine	19	7	37	65.00%	21	8	81	72.50%
pMpeg	11	11	65	55.00%	11	19	121	75.00%

Table 4 Experimental results of CANAT

Range representations of real numbers

Due to REs, we need to evaluate values of real numbers by some representation of ranges. They are classically classified into:

- Interval, which is the Cartesian product of one dimensional intervals [2, 40].
- Octagon, which is surrounded by either vertical, horizontal, or diagonal lines [39].
- Polyhedra, which is represented as the conjunction of linear inequalities [8]. Recently, its refinement SubPolyhedra was proposed [29] by reducing deduction rules among linear inequalities, yet preserving expressiveness.

We are more focus on intervals, and we call a range described by a pair of the lowest and the highest value by a classical interval (CI). CI is generalized to allow swapping of boundaries [18].

By introducing noise symbols, which preserve dependency of uncertainty, Affine interval (AI) has been proposed [54, 55]. Later, we will see how AI is applied as overapproximation for ORE analyses.

Extended Affine interval (EAI) has proposed for under approximation, based on the mean value theorem and Kaucher arithmetic [21]. EAI replaces real coefficients of AI with CI coefficients. We apply this idea for overapproximation of ORE analysis.

ORE analysis

For a static analysis, we need a concrete semantics. We obey the semantics of propagation of OREs to [32].

There are three kinds of OREs, caused by:

- real numbers to floating-point numbers conversion,
- real numbers to fixed-point numbers conversion, and
- floating-point numbers to fixed-point numbers conversion.

ORE analysis are mainly investigated for the first and the third.

For the real numbers to floating-point numbers conversion, ORE analysis adapt AI [20, 33] (which introduced widening operators to handle loops), except that the octagon abstract domain is used in [38].

They are implemented and showed experimental results. FLUCTUAT is presented in [20] and [38] showed a case study on an embedded avionics software. The

technique of [33] is further applied on TMS320 C3X assembler programs [34].

These ORE analyses are overapproximation, and easily cause false positives. An ORE analysis with under approximation is proposed based on the mean value theorem and Kaucher arithmetic [21], to sandwich OREs from both sides. However, strictly speaking, this under approximation is for real number variables rather than floating-point number variables.

The floating-point numbers to fixed-point numbers conversion typically appears in hard-wired algorithms and/or embedded systems. Apart from difficulties in hardware encoding difficulties [3, 5, 25, 26, 37, 51, 52], there are strong demand to solve ORE problems.

For the floating-point numbers to fixed-point numbers conversion, Fang, et.al. proposed an ORE analysis based on AI, intended for DSP applications [9–11]. We are facing on the same problem, but with different intervals, EAI. In our implementation, we adapted a sophisticated weighted model checking, whereas they adapt direct bit-vector encoding. For scalability, they also applied probabilistic reduction of the search space.

Thanks to the problem nature, we can examine OREs by testing, since we can compute both floating-point numbers and fixed-point numbers. [56] showed a such testing tool.

We further combined an ORE analysis and testing by a counter example guided narrowing approach, which refines the focus of testing and avoid false positives in an early stage.

Numerical Constraint Solvers

Recently, several tools have been developed as variations of SMT to solve non-linear numerical constrains. For instance,

- iSAT [12], which evaluates non-linear operations to interval constraints by overfapproximation.
- minimt [57], which covers specific irrationals, such as rational numbers and roots of small integers. They are symbolically represented and its bounded search is encoded as CNF.
- a tool for Simulink/Stateflow models [24], which applied a variation of polyhedra, called the *bounded vertex representation* for under-approximation.

Ganai and Ivancic [14] introduced a new method to face with decision problems involving non-linear constraints on bounded integers. Each nonlinear operation is encoded into a Boolean combination of linear arithmetic constraints based on CORDIC algorithm. Then, the linearized formula will be input of a DPLL-style Interval Search Engine that explores various combination of interval bounds using a SMT solver.

To solve an interval constraint, they divide the input ranges to smaller ranges, which is similar to ours. Adding to the difference of target domains (i.e., bounded integers and floating/fixed-point numbers), the differences are, (1) we use EAI instead of CI, (2) combination with testing, and (3) weighted model checking instead of (SMT) solver.

To solve dataflow equations over infinite domains, such as numerical constraints (mostly on integers), several algorithms are proposed in the context of weighted pushdown model checking [19, 15, 30, 41].

The library Apron of numerical abstract domains is also freely available [23].

Refining analyses and testing

As general setting of static detection, recently the refinement loop of analyses and testing is extensively investigated.

Counterexample Guided Abstraction-refinement (CEGAR) [6] is widely applied methodology, in which the initial abstract model typically nondeterministic control structures for conditional branches. When (possible spurious) counter examples are found, symbolic techniques refine the model by hooking more deterministic behavior.

Proofs from Tests [4] presented an algorithm DASH to check if a program satisfies a safety property. It uses only test generation, and it refines and maintains a sound program abstraction as a consequence of failed test generation operations. This enables us a light-weight refinement loop with neither any extra theorem prover calls nor any global may-alias information.

Our methodology of counter example guided narrowing tries to refine the focus of testing based on ORE analysis results. Fortunately, by the nature of EAI, ORE analysis results tell us which input parameter is dominant for REs. By using this information, we can effectively focus on the most problematic point.

7 Conclusion

Motivated by automatically detecting OREs of the hard-wired conversion of DSP encoders/decoders, this paper proposed techniques for automatic detection of overflow

and roundoff errors, caused by the floating-point number to fixed-point number conversion. Our contribution is summarized as follows.

- An extended affine interval (EAI) is newly proposed to overapproximate overflow and roundoff errors. EAI has two main advantages over current methods: EAI is more precise than CI and EAI forms are more compact than AI forms.
- An ORE analysis method based on weighted model checking is proposed and implemented as CANA.
- The counterexample-guided narrowing, in which an analysis and testing refine each other, is applied to the RE problem and implemented as CANAT. Thus, the result will be more precise than either testing or static analysis alone.

We also performed preliminary experiments of CANA and CANAT, which shows optimistic results.

For future works, one obstacle is the scalability of CANA and CANAT. We are optimistic since DSP algorithms (e.g., digital video compression [47]) are often compositional. They typically consist of sequences of computations with fixed-length arrays and bounded loops. Therefore, we can divide the algorithm to small fragments and check each separately.

Second, widening operator design. Currently, we did not introduce widening operators, but first focus on precision of ORE detection. (For instance, the widening operator [20] leads to inevitably lose precision a lot.) Its drawback is that the class of target programs has strong limitation, though it seems enough for the core part of DSP encoders/decoders.

We are also interested in an automatic correction of a program to improve with less OREs. For instance, in [36], there are several techniques to automatically improve numerical precision, such as:

- swapping the order of arithmetic operations,
- explicit shifting of the order of magnitude, and
- symbolic executions.

Our ORE analyzer, CANA, can generate the information about range values of fixed-point numbers and their RE at each point of the program. This information is helpful information to automatic source code correction.

References

1. Aho, A.V., Sethi, R., Ullman, J.D.: Compilers: Principles, Techniques, and Tools. Addison-Wesley (1986)
2. Alefeld, G. and Herzberger, J.: Introduction to Interval Computations. Academic Press, N.Y. (1983)

3. Banerjee, P., Bagchi, D., Haldar, M., Nayak, A., Kim, V., and Uribe, R.: Automatic Conversion of floating-point MATLAB Programs into fixed-point FPGA Based Hardware Design. In *Proc. of the 11th Annual IEEE Symposium on Field-Programmable Custom Computing Machines*, pp. 263, IEEE Computer Society (2003)
4. Beckman, N.E., Nori, A.V., Rajamani, S.K., Simmons, R.J., Tetali, S.D., Thakur, A.V.: Proofs from Tests. In *IEEE Transactions on Software Engineering: Special Issue on the ISSTA 2008 Best Papers*, IEEE Computer Society (2010)
5. Belanovic, P. and Rupp, M.: Automated Floating-Point to Fixed-Point Conversion with the Fixify Environment. In *Proc. of the 16th IEEE International Workshop on Rapid System Prototyping*, pp. 172 - 178, IEEE Computer Society (2007)
6. Clarke, E., Grumberg, O., Jha, S., Lu, Y., and Veith, H.: Counterexample-guided abstraction refinement. In *Proc. of the 12th International Conference on Computer Aided Verification*, pp. 154-169, Springer-Verlag (2000)
7. Cousot, P. and Cousot, R.: Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proc. of the 4th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages*, pp. 238-252, ACM (1977)
8. Cousot, P. and Halbwachs, N.: Automatic discovery of linear restraints among variables of a program. In *Proc. of the 5th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pp.84-96, ACM (1978)
9. Fang, C.F., Rutenbar, R.A., and Chen, T.: Fast, accurate static analysis for fixed-point finite precision effects in DSP designs. In *Proc. of the International Conference on Computer Aided Design*, pp. 275-282, IEEE Computer Society (2003)
10. Fang, C.F.: Probabilistic Interval-Valued Computation: Representing and Reasoning about Uncertainty in DSP and VLSI Design. Ph. D. Dissertation, Department of Electrical and Computer Engineering, Carnegie Mellon University (2005)
11. Fang, C.F., Rutenbar, R.A., Püschel, M. and Chen, T.: Toward efficient static analysis of finite-precision effects in DSP applications via Affine arithmetic modeling. In *Proc. of Design Automation Conference*, pp. 496-501, ACM Press (2003)
12. Franzle, M., Herde, C., Teige, T., Ratschan, S., and Schubert, T.: Efficient solving of large non-linear arithmetic constraint systems with complex boolean structure. In *Journal on Satisfiability, Boolean Modeling and Computation*, Vol 1, pp. 209-236 (2007)
13. Friedman, M. and Kandel, A.: Fundamentals of Computer numerical analysis. CRG Press (1994)
14. Ganai, M.K and Ivancic, F.: Efficient Decision Procedure for Bounded Integer Non-linear Operations using SMT(LIA), In *Proc. of Haifa Verification Conference 2008*, LNCS, vol.5394, pp.68-83, Springer (2009)
15. Gawlitza, T. and Seidl, H.: Precise Fixpoint Computation Through Strategy Iteration. In *Proc. of the European Symposium on Programming 2007*, LNCS, vol. 4421, pp. 300315, Springer (2007)
16. Giraud, L., Langou, J., Rozloznik, M., and Eshof, J.V.D.: Rounding error analysis of the classical Gram-Schmidt orthogonalization process. *Numerische Mathematik*, v.101 n.1, pp.87-100 (2005)
17. Goldberg, D.: What every computer scientist should know about floating-point arithmetic. *ACM Computing Surveys*, pp.5-48, ACM (1991)
18. Goldsztejn, A., Daney, D., Rueher, M., and Taillibert, P.: Modal intervals revisited: a mean-value extension to generalized intervals, In *Proc. of the 1st International Workshop on Quantification in Constraint Programming*, 2005.
19. Gopan, D.: Numeric program analysis techniques with applications to array analysis and library summarization. Ph. D Thesis, University of Wisconsin-Madison (2007)
20. Goubault, E. and Putot, S.: Static analysis of numerical algorithms. In *Proc. of the 13th International Static Analysis Symposium*, pp. 18-34, Springer-Verlag (2006)
21. Goubault, E. and Putot, S.: Under-approximations of computations in real numbers based on generalized affine arithmetic. In *Proc. of the 14th International Static Analysis Symposium*, pp. 137-152, Springer-Verlag (2007)
22. 754-2008 IEEE Standard for Floating-Point Arithmetic at: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4610935. Accessed February 2010.
23. Jeannot, B., and Miné, A.: Apron: A library of numerical abstract domains for static analysis. In *Proc. of the 21st International Conference on Computer Aided Verification*, LNCS, vol.5643, pp. 661-667, Springer-Verlag (2009)
24. Kanade, A., Alur, R., Ivanc, F., Ramesh, S., Sankaranarayanan, S., and Shashidhar, K.C.: Generating and Analyzing Symbolic Traces of Simulink/Stateflow Models. In *Proc. of the 21st International Conference on Computer Aided Verification*, LNCS, vol.5643, pp. 430-445, Springer-Verlag (2009)
25. Keding, H., Hurtgen, F., Willems, M., and Coors, M.: Transformation of Floating-Point into fixed-point Algorithms by Interpolation Applying a Statistical Approach. In *Proc. of the 9th International Conference on Signal Processing Applications and Technology*, pp. 270-276, ACM (1998)
26. Kim, S., Kum, K., and Sung, W.: Fixed-point optimization utility for C and C++ based digital signal processing programs. In *IEEE Transactions on Circuits and Systems II*, vol. 45, no.11, pp. 1455-1464, IEEE Computer Society (1998)
27. Knuth, D.E.: The Art of Computer Programming, Vol. 2: Seminumerical Algorithms. AddisonWesley (2002)
28. Lacey, D., Jones, N.D., Wyk, E.V., and Frederiksen, C.C.: Proving correctness of compiler optimizations by temporal logic. *POPL*, pp.283-294, ACM (2002)
29. Laviron, V. and Logozzo, F.: SubPolyhedra: A (more) scalable approach to infer linear inequalities. *Proc. of the 10th International Conference on Verification, Model Checking, and Abstract Interpretation*, LNCS, vol. 5403, pp. 229-244, Springer-Verlag (2009)
30. Leroux, J. and Sutre, G.: Accelerated data-flow analysis. In *Proc. of the 14th International Static Analysis Symposium*, LNCS, vol. 4634, pp. 184199, Springer (2007)
31. Loh, E. and Walster, G.W.: Rump's example revisited. *Reliable Computing*, vol.8 (3), pp.245-248 (2002)
32. Martel, M.: Semantics of roundoff error propagation in finite precision calculations. In *Higher-Order and Symbolic Computation*, v19(1), pp.7-30, Springer Netherlands (2006)
33. Martel, M.: Static analysis of the numerical stability of loops. In *Proc. of the 9th International Static Analysis Symposium*, LNCS, vol. 2477, pp. 133-150, Springer-Verlag (2002)
34. Martel, M.: Validation of assembler programs for DSPs: A static analyzer. In *Proc. of ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering*, pp. 8-13, ACM Press (2004)
35. Martel, M.: Semantics-Based Transformation of Arithmetic Expressions. In *Proc. of the 14th International Static Analysis Symposium*, LNCS, vol.4634 pp.298-314, Springer-Verlag (2007)
36. Martel, M.: Program transformation for numerical precision. In *Proc. of ACM SIGPLAN Workshop on Partial Evaluation and Program Manipulation 2009* pp. 101-109, ACM Press (2009)
37. Menard, D., Chillet, D., Charot, F., and Sentieys, O.: Automatic floating-point to fixed-point conversion for DSP code

- generation. In *Proc. of the 2002 international conference on Compilers, architecture, and synthesis for embedded systems*, pp. 270-276, ACM (2002)
38. Mine, A.: Relational abstract domains for the detection of floating-point run-time errors. In *Proc. of the European Symposium on Programming (ESOP 2004)*, Vol. 2986, pp. 317, Springer (2004)
 39. Mine, A.: The octagon abstract domain. *Higher-Order and Symbolic Computation*, Vol.19 (1), pp.31-100, 2006.
 40. Moore, R.E.: *Interval Analysis*. Prentice-Hall (1966)
 41. Mller-Olm, M. and Seidl, H.: Precise interprocedural analysis through linear algebra. In *Proc. of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pp. 330-341, ACM (2004)
 42. Necula, G.C., McPeak, S., Rahul, S.P., and Weimer, W.: CIL: Intermediate Language and Tools for Analysis and Transformation of C Programs. In *Proc. of the 11th International Conference on Compiler Construction*, pp. 213-228, Springer-Verlag (2002)
 43. Nielson, F., Nielson, H. R., Hankin, C.: *Principles of program analysis*. Springer (1998)
 44. Ngoc, D.T.B. and Ogawa, M.: Roundoff and Overflow Error Analysis via Model Checking. In *Proc. of the 7th International Conference on Software Engineering and Formal Methods*, pp.105-114, IEEE Computer Society (2009)
 45. Ngoc, D.T.B. and Ogawa, M.: Checking Roundoff Errors based on Counterexample-Guided Narrowing. In *Proc. of the 25th IEEE/ACM International Conference on Automated Software Engineering*, to appear as a short paper (2010)
 46. Ogita, T., Rump, S.M., and Oishi, S.: Accurate sum and dot product. *SIAM Sci. Comp*, vol.26 (6), pp.1955-1988 (2005)
 47. Peter, S.: *Digital Video Compression*. mcgraw-Hill (2004)
 48. Reps, T., Schwoon, S., Jha, S., and Melski D.: Weighted pushdown systems and their application to interprocedural dataflow analysis. In *Proc. of the 10th International Static Analysis Symposium*, pp. 189-213, Springer-Verlag (2003)
 49. Sagiv, M., Reps, T., and Horwitz, S.: Precise interprocedural dataflow analysis with applications to constant propagation. In *Theoretical Computer Science*, v.167 n.1-2, p.131-170 (1996)
 50. Schmidt, D.A. and Steffen, B.: Program analysis as model checking of abstract interpretations. In *Proc. of the 5th International Conference on Static Analysis Symposium*, pp. 351-380, Springe-Verlag (1998)
 51. Shi, C. and Brodersen, R.W.: An automated floating-point to fixed-point conversion methodology. In *Proc. IEEE International Conference on Acoust., Speech, and Signal Processing*, Vol. 2, pp. 529-532 (2003)
 52. Shi, C.: *Floating-point to Fixed-point Conversion*. Ph. D. Thesis, University of California - Berkeley (2004)
 53. Steffen, B.: Data Flow Analysis as Model Checking. In *Proc. of the International Conference on Theoretical Aspects of Computer Software*, p.346-365, Springer-Verlag (1991)
 54. Stolfi, J.: *Self-Validated Numerical Methods and Applications*. Ph. D. Dissertation, Computer Science Department, Stanford University (1997)
 55. Stolfi, J. and Figueiredo, L.H. de: An introduction to affine arithmetic. In *Tendencias em Matematica Aplicada e Computacional*, Vol.3(4), pp. 297-312 (2003)
 56. Wijaya, S. and Cantoni, A.: A Java Simulation Tool for Fixed-Point System Design. In *Proc. of the 2nd International Conference on Simulation Tools and Techniques*, pp. 1-10, ACM (2009)
 57. Zankl, H. and Middeldorp, A.: Satisfiability of non-linear (ir)rational arithmetic. In *Proc. of the 16th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning*, to appear (2010).