

# Certified Roundoff Error Bounds Using Semidefinite Programming

VICTOR MAGRON, CNRS Verimag

GEORGE CONSTANTINIDES and ALASTAIR DONALDSON, Imperial College London

Roundoff errors cannot be avoided when implementing numerical programs with finite precision. The ability to reason about rounding is especially important if one wants to explore a range of potential representations, for instance, for FPGAs or custom hardware implementations. This problem becomes challenging when the program does not employ solely linear operations as non-linearities are inherent to many interesting computational problems in real-world applications.

Existing solutions to reasoning possibly lead to either inaccurate bounds or high analysis time in the presence of nonlinear correlations between variables. Furthermore, while it is easy to implement a straightforward method such as interval arithmetic, sophisticated techniques are less straightforward to implement in a formal setting. Thus there is a need for methods that output certificates that can be formally validated inside a proof assistant.

We present a framework to provide upper bounds on absolute roundoff errors of floating-point nonlinear programs. This framework is based on optimization techniques employing semidefinite programming and sums of squares certificates, which can be checked inside the Coq theorem prover to provide formal roundoff error bounds for polynomial programs. Our tool covers a wide range of nonlinear programs, including polynomials and transcendental operations as well as conditional statements. We illustrate the efficiency and precision of this tool on non-trivial programs coming from biology, optimization, and space control. Our tool produces more accurate error bounds for 23% of all programs and yields better performance in 66% of all programs.

**CCS Concepts:** • **Theory of computation** → **Automated reasoning**; **Numeric approximation algorithms**; **Semidefinite programming**; **Convex optimization**

**Additional Key Words and Phrases:** Correlation sparsity pattern, floating-point arithmetic, formal verification, polynomial optimization, proof assistant, roundoff error, semidefinite programming, transcendental functions

## ACM Reference Format:

Victor Magron, George Constantinides, and Alastair Donaldson. 2017. Certified roundoff error bounds using semidefinite programming. *ACM Trans. Math. Softw.* 43, 4, Article 34 (January 2017), 31 pages.

DOI: <http://dx.doi.org/10.1145/3015465>

## 1. INTRODUCTION

Constructing numerical programs that perform accurate computation turns out to be difficult, due to finite numerical precision of implementations such as floating-point

---

This work was partly funded by the Engineering and Physical Sciences Research Council (EPSRC) “Challenging Engineering” Grant (EP/I020457/1) and European Research Council (ERC) “STATOR” Grant Agreement nr. 306595.

Authors’ addresses: V. Magron, Verimag, Bâtiment Imag, Office 2053 700 Avenue Centrale, 38401 Saint-Martin d’Hères France; email: [victor.magron@imag.fr](mailto:victor.magron@imag.fr); G. Constantinides, Imperial College, Circuits and Systems, South Kensington Campus London SW7 2AZ, UK; email: [george.constantinides@imperial.ac.uk](mailto:george.constantinides@imperial.ac.uk); A. Donaldson, Imperial College London, Room 422, Huxley Building, 180 Queen’s Gate, London SW7 2AZ UK; email: [alastair.donaldson@imperial.ac.uk](mailto:alastair.donaldson@imperial.ac.uk).

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2017 ACM 0098-3500/2017/01-ART34 \$15.00

DOI: <http://dx.doi.org/10.1145/3015465>

or fixed-point representations. Finite-precision numbers induce roundoff errors, and knowledge of the range of these roundoff errors is required to fulfill safety criteria of critical programs, as typically arising in modern embedded systems such as aircraft controllers. Such a knowledge can be used in general for developing accurate numerical software but is also particularly relevant when considering migration of algorithms onto hardware (e.g., FPGAs). The advantage of architectures based on FPGAs is that they allow more flexible choices rather than choosing either for IEEE standard single or double precision. Indeed, in this case, we benefit from a more flexible number representation while still ensuring guaranteed bounds on the program output.

To obtain lower bounds on roundoff errors, one can rely on testing approaches, such as meta-heuristic search [Borges et al. 2012] or under-approximation tools (e.g., `s3fp` [Chiang et al. 2014]). Here we are interested in efficiently handling the complementary over-approximation problem, namely to obtain precise upper bounds on the error. This problem boils down to finding tight abstractions of linearities or non-linearities while being able to bound the resulting approximations in an efficient way. For computer programs consisting of linear operations, automatic error analysis can be obtained with well-studied optimization techniques based on satisfiability modulo theories (SAT/SMT) solvers [Haller et al. 2012] and affine arithmetic [Delmas et al. 2009]. However, non-linear operations are key to many interesting computational problems arising in physics, biology, controller implementations, and global optimization. Recently, two promising frameworks have been designed to provide upper bounds for roundoff errors of nonlinear programs. The corresponding algorithms rely on Taylor-interval methods [Solovyev et al. 2015], implemented in the `FPTaylor` tool, and on combining SMT with interval arithmetic [Darulova and Kuncak 2014], implemented in the `Rosa` real compiler.

The complexity of the mathematics underlying techniques for nonlinear reasoning, and the intricacies associated with constructing an efficient implementation, are such that a means for independent formal validation of results is particularly desirable. The `Rosa` tool is based on theoretical results that should provide sound over-approximations of error bounds. While `Rosa` relies on an SMT solver capable of generating unsatisfiability proof witnesses (thus allowing independent soundness checking), it does not verify formally these certificates inside a proof assistant. To the best of our knowledge, the `FPTaylor` and `GAPPA` software are the only academic tools that can produce formal proof certificates. For `FPTaylor`, this is based on the framework developed in Solovyev and Hales [2013] to verify nonlinear inequalities in `HOL-LIGHT` [Harrison 1996] using Taylor-interval methods. However, most of computation performed in the informal optimization procedure ends up being redone inside the `HOL-LIGHT` proof assistant, yielding a formal verification that may be computationally demanding.

The aim of this work is to provide a formal framework to perform automated precision analysis of computer programs that manipulate finite-precision data using nonlinear operators. For such programs, guarantees can be provided with certified programming techniques. Semidefinite programming (SDP) is relevant to a wide range of mathematical fields, including combinatorial optimization, control theory, and matrix completion. In 2001, Lasserre introduced a hierarchy of SDP relaxations [Lasserre 2001] for approximating polynomial infima. Our method to bound the error is a decision procedure based on a specialized variant of the Lasserre hierarchy [Lasserre 2006]. The procedure relies on SDP to provide sparse sum-of-squares decompositions of nonnegative polynomials. Our framework handles polynomial program analysis (involving the operations  $+$ ,  $\times$ ,  $-$ ) as well as extensions to the more general class of semialgebraic and transcendental programs (involving  $\sqrt{\cdot}$ ,  $\min$ ,  $\max$ ,  $\arctan$ ,  $\exp$ ), following the approximation scheme described in Magron et al. [2015a].

### 1.1. Overview of our Method

We present an overview of our method and of the capabilities of related techniques, using an example. Consider a program implementing the following polynomial expression  $f$ :

$$\begin{aligned} f(\mathbf{x}) := & x_2 \times x_5 + x_3 \times x_6 - x_2 \times x_3 - x_5 \times x_6 \\ & + x_1 \times (-x_1 + x_2 + x_3 - x_4 + x_5 + x_6), \end{aligned}$$

where the six-variable vector  $\mathbf{x} := (x_1, x_2, x_3, x_4, x_5, x_6)$  is the input of the program. For this example, assume that the set  $\mathbf{X}$  of possible input values is a product of closed intervals:  $\mathbf{X} = [4.00, 6.36]^6$ . This function  $f$  together with the set  $\mathbf{X}$  appear in many inequalities arising from the proof of the Kepler Conjecture [Hales 2006], yielding challenging global optimization problems.

The polynomial expression  $f$  is obtained by performing 15 basic operations (1 negation, 3 subtractions, 6 additions, and 5 multiplications). When executing this program with a set of floating-point numbers  $\hat{\mathbf{x}} := (\hat{x}_1, \hat{x}_2, \hat{x}_3, \hat{x}_4, \hat{x}_5, \hat{x}_6) \in \mathbf{X}$ , one actually computes a floating-point result  $\hat{f}$ , where all operations  $+$ ,  $-$ ,  $\times$  are replaced by the respectively associated floating-point operations  $\oplus$ ,  $\ominus$ ,  $\otimes$ . The results of these operations comply with IEEE 754 standard arithmetic [IEEE 2008] (see relevant background in Section 2.1). Here, for the sake of clarity, we do not consider real input variables but we do it later while performing detailed comparison (see Section 4). For instance (in the absence of underflow), one can write  $\hat{x}_2 \otimes \hat{x}_5 = (x_2 \times x_5)(1 + e_1)$  by introducing an error variable  $e_1$  such that  $-\epsilon \leq e_1 \leq \epsilon$ , where the bound  $\epsilon$  is the machine precision (e.g.,  $\epsilon = 2^{-24}$  for single precision). One would like to bound the absolute roundoff error  $|r(\mathbf{x}, \mathbf{e})| := |\hat{f}(\mathbf{x}, \mathbf{e}) - f(\mathbf{x})|$  over all possible input variables  $\mathbf{x} \in \mathbf{X}$  and error variables  $e_1, \dots, e_{15} \in [-\epsilon, \epsilon]$ . Let us define  $\mathbf{E} := [-\epsilon, \epsilon]^{15}$  and  $\mathbf{K} := \mathbf{X} \times \mathbf{E}$ . Then our bound problem can be cast as finding the maximum  $r^*$  of  $|r|$  over  $\mathbf{K}$ , yielding the following nonlinear optimization problem:

$$\begin{aligned} r^* := & \max_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} |r(\mathbf{x}, \mathbf{e})| \\ = & \max\{-\min_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} r(\mathbf{x}, \mathbf{e}), \max_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} r(\mathbf{x}, \mathbf{e})\}. \end{aligned} \quad (1)$$

One can directly try to solve these two polynomial optimization problems using classical SDP relaxations [Lasserre 2001]. As in Solovyev et al. [2015], one can also decompose the error term  $r$  as the sum of a term  $l(\mathbf{x}, \mathbf{e})$ , which is affine w.r.t.  $\mathbf{e}$ , and a nonlinear term  $h(\mathbf{x}, \mathbf{e}) := r(\mathbf{x}, \mathbf{e}) - l(\mathbf{x}, \mathbf{e})$ . Then the triangular inequality yields

$$r^* \leq \max_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} |l(\mathbf{x}, \mathbf{e})| + \max_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} |h(\mathbf{x}, \mathbf{e})|. \quad (2)$$

It follows for this example that  $l(\mathbf{x}, \mathbf{e}) = x_2x_5e_1 + x_3x_6e_2 + (x_2x_5 + x_3x_6)e_3 + \dots + f(\mathbf{x})e_{15} = \sum_{i=1}^{15} s_i(\mathbf{x})e_i$ , with  $s_1(\mathbf{x}) := x_2x_5$ ,  $s_2(\mathbf{x}) := x_3x_6$ ,  $\dots$ ,  $s_{15}(\mathbf{x}) := f(\mathbf{x})$ . The *Symbolic Taylor expansions* method [Solovyev et al. 2015] consists of using a simple branch and bound algorithm based on interval arithmetic to compute a rigorous interval enclosure of each polynomial  $s_i$ ,  $i = 1, \dots, 15$ , over  $\mathbf{X}$  and finally obtain an upper bound of  $|l| + |h|$  over  $\mathbf{K}$ . In contrast, our method uses sparse semidefinite relaxations for polynomial optimization (derived from Lasserre [2006]) to bound  $l$  and basic interval arithmetic as in Solovyev et al. [2015] to bound  $|h|$  (i.e., we use interval arithmetic to bound second-order error terms in the multivariate Taylor expansion of  $r$  w.r.t.  $\mathbf{e}$ ).

The following comparison results have been obtained on an Intel Core i7-5600U CPU (2.60GHz). All execution times have been computed by averaging over five runs.

—A direct attempt to solve the two polynomial problems occurring in Equation (1) fails as the SDP solver (in our case SDPA [Yamashita et al. 2010]) runs out of memory.

- Using our method implemented in the `Real2Float` tool, one obtains an upper bound of  $760\epsilon$  for  $|l| + |h|$  over  $\mathbf{K}$  in 0.15s. This bound is provided together with a certificate that can be formally checked inside the Coq proof assistant in 0.20 seconds.
- After normalizing the polynomial expression and using basic interval arithmetic, one obtains 8 times more quickly a coarser bound of  $922\epsilon$ .
- Symbolic Taylor expansions implemented in `FPTaylor` [Solovyev et al. 2015] provide a more precise bound of  $721\epsilon$  but 28 times slower than with our implementation. Formal verification of this bound inside the `HOL-LIGHT` proof assistant takes 27.7s and is 139 times slower than proof checking with `Real2Float` inside Coq. One can obtain an even more precise bound of  $528\epsilon$  (but 37 times slower than with our implementation) by turning on the improved rounding model of `FPTaylor` and bound the number of branch and bound iterations with 10,000. The drawback of this bound is that it cannot be formally verified.
- Finally, a slightly coarser bound of  $762\epsilon$  is obtained with the Rosa real compiler [Darulova and Kuncak 2014] but 19 times slower than with our implementation and we cannot get formal verification of this bound.

## 1.2. Related Works

SMT solvers allow to analyze programs with various semantics or specifications but are limited for the manipulation of problems involving nonlinear arithmetic. Several solvers, including Z3 [De Moura and Bjørner 2008], provide partial support for the IEEE floating-point standard [Rümmer and Wahl 2010]. They suffer from a lack of scalability when used for roundoff error analysis in isolation (as emphasized in Darulova and Kuncak [2014]) but can be integrated into existing frameworks, for example, `FPHILE` [Paganelli and Ahrendt 2013]. The procedure in Gao et al. [2013] can solve SMT problems over the real numbers, using interval constraint propagation but has not yet been applied to quantification of roundoff error. The Rosa tool [Darulova and Kuncak 2014] provides a way to compile functional `SCALA` programs involving semialgebraic functions and conditional statements. The tool uses affine arithmetic to provide sound over-approximations of roundoff errors, allowing for generation of finite precision implementations that fulfill the required precision given as input by the user. This tool thus relies on abstract interpretation but bounds of the affine expressions are provided through an optimization procedure based on SMT. In our case, we use the same rounding model but provide approximations that are affine w.r.t. the additional error variables and nonlinear w.r.t. the input variables. Instead of using SMT, we bound the resulting expressions with optimization techniques based on semidefinite programming. Abstract interpretation [Cousot and Cousot 1977] has been extensively used in the context of static analysis to provide sound over-approximations, called *abstractions*, of the sets of values taken by program variables. The effects of variable assignments, guards and conditional branching statements are handled with several domain-specific operators (e.g. inclusion, meet, and join). Well-studied abstract domains include intervals [Moore 1962] as well as more complicated frameworks based on affine arithmetic [Stolfi and de Figueiredo 2003], octogons [Miné 2006], zonotopes [Ghorbal et al. 2010], polyhedra [Chen et al. 2008], and interval polyhedra [Chen et al. 2009], some of them being implemented inside a tool called `APRON` [Jeannet and Miné 2009]. Abstract domains provide sound over-approximations of program expressions and allow upper bounds on roundoff error to be computed. The `GAPPA` tool [Daumas and Melquiond 2010] relies on interval abstract domains with an extension to affine domains [Linderman et al. 2010] to reason about roundoff errors. As demonstrated in Solovyev et al. [2015], the bounds obtained inside `GAPPA` are often coarser than other methods. Formal guarantees can be provided as `GAPPA` benefits from an interface with Coq while making use of interval libraries [Melquiond 2012] relying on formalized floating-points [Boldo and Melquiond

Table I. Comparison of Roundoff Error Tools w.r.t. Expressiveness

Feature	Real2Float	Rosa	FPTaylor	GAPPA	FLUCTUAT
Basic FP operations/formats	✓	✓	✓	✓	✓
Special values ( $\pm\infty$ , NaN)					✓
Improved rounding model			✓	✓	✓
Input uncertainties	✓	✓	✓	✓	✓
Transcendental functions	✓		✓		
Discontinuity errors	✓	✓			✓
Proof certificates	✓		✓	✓	

2011]. The static analysis commercial tool FLUCTUAT (with a free academic version) relies on affine abstract domains [Blanchet et al. 2003] and techniques that are very similar to the ones in Rosa, including interval subdivision. This tool does not perform optimization but uses forward computation to analyze floating-point programs written in C. Furthermore, FLUCTUAT also has a procedure for discontinuity errors [Ghorbal et al. 2010]. GAPPA and FLUCTUAT tools use a different rounding model (also available as an option inside FPTaylor) based on a piecewise constant absolute error bound. This is more precise than the simple rounding model used in our framework but requires (possibly) extensive use of a branch and bound algorithm as each interval has to be subdivided in intervals  $[2^n, 2^{n+1}]$  for several values of the integer  $n$ . In Solovyev et al. [2015], the authors provide a table (Table 1) comparing relevant features of FPTaylor with three other tools (Rosa, GAPPA, and FLUCTUAT) performing roundoff error estimation. In a similar fashion, we summarize the main features related to our tool Real2Float and the same four above-mentioned tools used for our further benchmark comparisons w.r.t. their expressiveness in Table I.

Computing sound bounds of nonlinear expressions is mandatory to perform formal analysis of finite precision implementations and can be performed with various optimization tools. In the polynomial case, alternative approaches to semidefinite relaxations are based on decomposition in the multivariate Bernstein basis. Formal verification of bounds obtained with this decomposition has been investigated by Muñoz and Narkawicz [Muñoz and Narkawicz 2013] in the PVS theorem prover. We are not aware of any work based on these techniques that can quantify roundoff errors. Another decomposition of nonnegative polynomials into sums of squares (SOS) certificates consists in using the Krivine [1964] and Handelman [1988] representation and boils down to solving linear programming (LP) relaxations. In our case, we use a different representation, leading to solve SDP relaxations. The Krivine-Handelman representation has been used in Boland and Constantinides [2010] to compute roundoff error bounds. LP relaxations often provide coarser bounds than SDP relaxations, and it has been proven in Lasserre [2009] that generically finite convergence does not occur for convex problems, at the exception of the linear case. The work in Roux [2015] focuses on formalization of roundoff errors bounds related to positive definiteness verification. Branch and bound methods with Taylor models [Berz and Makino 2009] are not restricted to polynomial systems and have been formalized Solovyev and Hales [2013] to solve nonlinear inequalities occurring in the proof of Kepler Conjecture. Symbolic Taylor Expansions [Solovyev et al. 2015] have been implemented in the FPTaylor tool to compute formal bounds of roundoff errors for programs involving both polynomial and transcendental functions.

### 1.3. Contributions

Our key contributions can be summarized as follows:

- We present an optimization algorithm providing sound over-approximations for roundoff errors of floating-point nonlinear programs. This algorithm is based on



sparse sums of squares programming [Lasserre 2006]. In comparison with other methods, our algorithm allows us to obtain tighter upper bounds, while overcoming scalability and numerical issues inherent in SDP solvers [Todd 2001]. Our algorithm can currently handle programs implementing polynomial functions but also involving non-polynomial components, including either semialgebraic or transcendental operations (e.g.,  $/$ ,  $\sqrt{\cdot}$ ,  $\arctan$ ,  $\exp$ ), as well as conditional statements. Programs containing iterative or while loops are not currently supported.

- Our framework is fully implemented in the `Real2Float` tool. Among several features, the tool can optionally perform formal verification of roundoff error bounds for polynomial programs, inside the Coq proof assistant [Coq 2016]. The most recent software release of `Real2Float` provides OCAML [2015] and Coq libraries and is freely available.<sup>1</sup> Our implementation tool is built on top of the `NLCertify` verification system [Magron 2014]. Precision and efficiency of the tool are evaluated on several benchmarks coming from the existing literature. Numerical experiments demonstrate that our method competes well with recent approaches relying on Taylor-interval approximations [Solovyev et al. 2015] or combining SMT solvers with affine arithmetic [Darulova and Kuncak 2014]. We also compared our tool with `GAPPA` [Daumas and Melquiond 2010] and `FLUCTUAT` [Delmas et al. 2009].

The article is organized as follows. In Section 2, we present mandatory background on roundoff errors due to finite precision arithmetic before describing our nonlinear program semantics (Section 2.1). Then we recall how to perform certified polynomial optimization based on semidefinite programming (Section 2.2) and how to obtain formal bounds while checking the certificates inside the Coq proof assistant (Section 2.3). Section 3 contains the main contribution of the article, namely how to compute tight over-approximations for roundoff errors of nonlinear programs with sparse semidefinite relaxations. Finally, Section 4 is devoted to the evaluation of our nonlinear verification tool `Real2Float` on benchmarks arising from control systems, optimization, physics, and biology, as well as comparisons with the tools `FPTaylor`, `Rosa`, `GAPPA`, and `FLUCTUAT`.

## 2. PRELIMINARIES

### 2.1. Program Semantics and Floating-Point Numbers

We support conditional code without procedure calls or loops. Despite these restrictions, we can consider a wide range of nonlinear programs while assuming that important numerical calculations can be expressed in a loop-free manner. Our programs are encoded in an ML-like language:

```
let box_prog       $x_1 \dots x_n = [(a_1, b_1); \dots; (a_n, b_n)];;$ 
let obj_prog       $x_1 \dots x_n = [(f(\mathbf{x}), \epsilon_{\text{Real2Float}})];;$ 
let cstr_prog       $x_1 \dots x_n = [g_1(\mathbf{x}); \dots; g_k(\mathbf{x})];;$ 
let uncert_prog    $x_1 \dots x_n = [u_1; \dots; u_n];;$ 
```

Here, the first line encodes interval floating-point bound constraints for input variables, namely  $\mathbf{x} := (x_1, \dots, x_n) \in [a_1, b_1] \times \dots \times [a_n, b_n]$ . The second line provides the function  $f(\mathbf{x})$  as well as the total roundoff error bound  $\epsilon_{\text{Real2Float}}$ . Then, one encodes polynomial nonnegativity constraints over the input variables, namely  $g_1(\mathbf{x}) \geq 0, \dots, g_k(\mathbf{x}) \geq 0$ . Finally, the last line allows the user to specify a numerical constant  $u_i$  to associate a given uncertainty to the variable  $x_i$ , for each  $i = 1, \dots, n$ .

The type of numerical constants is denoted by `C`. In our current implementation, the user can choose either 64-bit floating-point or arbitrary-size rational numbers. This

<sup>1</sup>[forge.ocamlcore.org/frs/?group\\_id=351](https://forge.ocamlcore.org/frs/?group_id=351).

type  $C$  is used for the terms  $\epsilon_{\text{Real2Float}}, u_1, \dots, u_n, a_1, \dots, a_n, b_1, \dots, b_n$ . The inductive type of polynomial expressions with coefficients in  $C$  is  $\text{pExprC}$  defined as follows:

```
type pexprC = Pc of C | Px of positive
| Psub of pexprC * pexprC | Pneg of pexprC
| Padd of pexprC * pexprC
| Pmul of pexprC * pexprC
```

The constructor  $Px$  takes a positive integer as argument to represent either an input or local variable. The inductive type  $\text{nlexpr}$  of nonlinear expressions (such as  $f(\mathbf{x})$ ) is defined as follows:

```
type nlexpr =
| Pol of pexprC | Neg of nlexpr
| Add of nlexpr * nlexpr
| Mul of nlexpr * nlexpr
| Sub of nlexpr * nlexpr
| Div of nlexpr * nlexpr | Sqrt of nlexpr
| Transc of transc * nlexpr
| IfThenElse of pexprC * nlexpr * nlexpr
| Let of positive * nlexpr * nlexpr
```

The type  $\text{transc}$  corresponds to the dictionary  $\mathcal{D}$  of special functions. In our case  $\mathcal{D} := \{\exp, \log, \cos, \sin, \tan, \arccos, \arcsin, \arctan\}$ . For instance, the term  $\text{Transc}(\exp, f(\mathbf{x}))$  represents the program implementing  $\exp(f(\mathbf{x}))$ .

Given a polynomial expression  $p$  and two nonlinear expressions  $f$  and  $g$ , the term  $\text{IfThenElse}(p(\mathbf{x}), f(\mathbf{x}), g(\mathbf{x}))$  represents the conditional program implementing<sup>2</sup> if  $(p(\mathbf{x}) \geq 0)$   $f(\mathbf{x})$  else  $g(\mathbf{x})$ . The constructor  $\text{Let}$  allows us to define local variables in an ML fashion, for example,  $\text{let } t_1 = 331.4 + 0.6 * T \text{ in } -t_1 * v / ((t_1 + u) * (t_1 + u))$  (part of the *doppler1* program considered in Section 4).

Finally, one obtains rounded nonlinear expressions using a recursive procedure  $\text{round}$ , defined according to Equation (3) and Equation (4). Rounded expressions are supported inside conditions. When an uncertainty  $u_i$  is specified for an input variable  $x_i$ , the corresponding rounded expression is given by  $x_i(1 + e)$ , with  $|e| \leq u_i$ , the uncertainty  $u_i$  being a relative error.

We adopt the standard practice [Higham 2002] to approximate a real number  $x$  with its closest floating-point representation  $\hat{x} = x(1 + e)$ , where  $|e|$  is less than the machine precision  $\epsilon$ . In the sequel, we neglect both overflow and denormal range values. The operator  $\hat{\cdot}$  is called the rounding operator and can be selected among rounding to nearest, rounding toward zero (respectively  $\pm\infty$ ). In the sequel, we assume rounding to nearest. The scientific notation of a binary (respectively, decimal) floating-point number  $\hat{x}$  is a triple  $(s, \text{sig}, \text{exp})$  consisting of a sign bit  $s$ , a *significand*  $\text{sig} \in [1, 2)$  (respectively  $[1, 10)$ ), and an *exponent*  $\text{exp}$ , yielding numerical evaluation  $(-1)^s \text{sig} 2^{\text{exp}}$  (respectively  $(-1)^s \text{sig} 10^{\text{exp}}$ ).

The value of  $\epsilon$  actually gives the upper bound on the relative floating-point error and is equal to  $2^{-\text{prec}}$ , where  $\text{prec}$  is called the *precision*, referring to the number of significant bits used. For single precision floating-point, one has  $\text{prec} = 24$ . For double (respectively, quadruple) precision, one has  $\text{prec} = 53$  (respectively,  $\text{prec} = 113$ ). Let  $\mathbb{R}$  denote the set of real numbers and  $\mathbb{F}$  the set of binary floating-point numbers. For each real-valued operation  $\text{bop}_{\mathbb{R}} \in \{+, -, \times, /\}$ , the result of the corresponding

<sup>2</sup>Our general framework could theoretically handle nested if statements. However, our current implementation is limited to programs involving single conditional statements at the code toplevel.

floating-point operation  $\text{bop}_{\mathbb{F}} \in \{\oplus, \ominus, \otimes, \oslash\}$  satisfies the following when complying with IEEE 754 standard arithmetic [IEEE 2008] (without overflow, underflow, and denormal occurrences):

$$\text{bop}_{\mathbb{F}}(\hat{x}, \hat{y}) = \text{bop}_{\mathbb{R}}(\hat{x}, \hat{y})(1 + e), \quad |e| \leq \epsilon = 2^{-\text{prec}}. \quad (3)$$

Other operations include special functions taken from a *dictionary*  $\mathcal{D}$ , containing the unary functions  $\tan$ ,  $\arctan$ ,  $\cos$ ,  $\arccos$ ,  $\sin$ ,  $\arcsin$ ,  $\exp$ ,  $\log$ ,  $(\cdot)^r$  with  $r \in \mathbb{R} \setminus \{0\}$ . For  $f_{\mathbb{R}} \in \mathcal{D}$ , the corresponding floating-point evaluation satisfies

$$f_{\mathbb{F}}(\hat{x}) = f_{\mathbb{R}}(\hat{x})(1 + e), \quad |e| \leq \epsilon(f_{\mathbb{R}}). \quad (4)$$

The value of the relative error bound  $\epsilon(f_{\mathbb{R}})$  differs from the machine precision  $\epsilon$  in Equation (3) and has to be properly adjusted on a per-operator basis. We refer the interested reader to Bingham and Leslie-Hurd [2014] for relative error bound verification of transcendental functions (see also Harrison [2000] for formalization in HOL-LIGHT).

## 2.2. SDP Relaxations for Polynomial Optimization

The sums of squares method involves approximation of polynomial inequality constraints by SOS equality constraints. Here we recall mandatory background about SOS. We apply this method in Section 3 to solve the problems of Equation (1) when the nonlinear function  $r$  is a polynomial. In the sequel, let us denote by  $n$  the number of initial variables of the polynomial optimization problem and by  $k$  the number of optimization constraints.

**2.2.1. Sums of Squares Certificates and SDP.** First, we recall basic facts about generation of SOS certificates for polynomial optimization, using semidefinite programming, which can be found in texts such as Lasserre [2001]. Denote by  $\mathbb{R}[\mathbf{x}]$  the vector space of polynomials and by  $\mathbb{R}_{2d}[\mathbf{x}]$  the restriction of  $\mathbb{R}[\mathbf{x}]$  to polynomials of degree at most  $2d$ . Let us define the set of SOS polynomials:

$$\Sigma[\mathbf{x}] := \left\{ \sum_i q_i^2, \text{ with } q_i \in \mathbb{R}[\mathbf{x}] \right\}, \quad (5)$$

as well as its restriction  $\Sigma_{2d}[\mathbf{x}] := \Sigma[\mathbf{x}] \cap \mathbb{R}_{2d}[\mathbf{x}]$  to polynomials of degree at most  $2d$ . For instance, the following bivariate polynomial  $\sigma(\mathbf{x}) := 1 + (x_1^2 - x_2^2)^2$  lies in  $\Sigma_4[\mathbf{x}] \subseteq \mathbb{R}_4[\mathbf{x}]$ .

Optimization methods based on SOS use the implication  $r \in \Sigma[\mathbf{x}] \Rightarrow \forall \mathbf{x} \in \mathbb{R}^n, r(\mathbf{x}) \geq 0$ , that is, the inclusion of  $\Sigma[\mathbf{x}]$  in the set of nonnegative polynomials.

The underlying reason for using SOS polynomials is that optimizing over positive polynomials is NP Hard [Laurent 2009]. Thus, one would like to replace such positivity constraints by more tractable ones, and in particular the SOS decompositions admitted by positive polynomials provide a suitable alternative: When fixing the degree of such decompositions, the resulting relaxed problem becomes more tractable.

Given  $r \in \mathbb{R}[\mathbf{x}]$ , one considers the following polynomial minimization problem:

$$r^* := \inf_{\mathbf{x} \in \mathbb{R}^n} \{r(\mathbf{x}) : \mathbf{x} \in \mathbf{K}\}, \quad (6)$$

where the set of constraints  $\mathbf{K} \subseteq \mathbb{R}^n$  is defined by

$$\mathbf{K} := \{\mathbf{x} \in \mathbb{R}^n : g_1(\mathbf{x}) \geq 0, \dots, g_k(\mathbf{x}) \geq 0\},$$

for polynomial functions  $g_1, \dots, g_k$ . The set  $\mathbf{K}$  is called a *basic semialgebraic* set. Membership of semialgebraic sets is ensured by satisfying conjunctions of polynomial non-negativity constraints.

**Remark 2.1.** When the input variables satisfy interval constraints  $\mathbf{x} \in [a_1, b_1] \times \dots \times [a_n, b_n]$ , then one can easily show that there exists some integer  $M > 0$  such



that  $M - \sum_{i=1}^n x_i^2 \geq 0$ . In the sequel, we assume that this nonnegativity constraint appears explicitly in the definition of  $\mathbf{K}$ . Such an assumption is mandatory to prove the convergence of semidefinite relaxations recalled in Theorem 2.3.

In general, the objective function  $r$  and the set of constraints  $\mathbf{K}$  can be nonconvex, which makes Problem (6) difficult to solve in practice. One can rewrite Problem (6) as the equivalent maximization problem:

$$r^* := \sup_{\mu \in \mathbb{R}} \{ \mu : r(\mathbf{x}) - \mu \geq 0, \forall \mathbf{x} \in \mathbf{K} \}. \quad (7)$$

Now we outline how to handle the nonnegativity constraint  $r - \mu \geq 0$ . Given a non-negative polynomial  $p \in \mathbb{R}[\mathbf{x}]$ , the existence of an SOS decomposition  $p = \sum_i q_i^2$  valid over  $\mathbb{R}^n$ , is equivalent to the existence of a symmetric real matrix  $\mathbf{Q}$ , a solution of the following linear matrix feasibility problem:

$$r(\mathbf{x}) = \mathbf{m}_d(\mathbf{x})^\top \mathbf{Q} \mathbf{m}_d(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^n, \quad (8)$$

where  $\mathbf{m}_d(\mathbf{x}) := (1, x_1, \dots, x_n, x_1^2, x_1x_2, \dots, x_n^d)$  and the matrix  $\mathbf{Q}$  has only nonnegative eigenvalues. Such a matrix  $\mathbf{Q}$  is called *positive semidefinite*. The vector  $\mathbf{m}_d$  (respectively, matrix  $\mathbf{Q}$ ) has a size (respectively, dimension) equal to  $s_n^d := \binom{n+d}{d}$ . Problem (8) can be handled with SDP solvers, such as MOSEK [Andersen and Andersen 2000] or SDPA [Yamashita et al. 2010] (see Vandenberghe and Boyd [1994] for specific background about SDP). Then, one computes the “LDL” decomposition  $\mathbf{Q} = \mathbf{L}^\top \mathbf{D} \mathbf{L}$  (a variant of the classical Cholesky decomposition), where  $\mathbf{L}$  is a lower triangular matrix and  $\mathbf{D}$  is a diagonal matrix. Finally, one obtains  $r(\mathbf{x}) = (\mathbf{L} \mathbf{m}_d(\mathbf{x}))^\top \mathbf{D} (\mathbf{L} \mathbf{m}_d(\mathbf{x})) = \sum_{i=1}^{s_n^d} q_i(\mathbf{x})^2$ . Such a decomposition is called an SOS *certificate*.

*Example 2.2.* Let us define  $r(\mathbf{x}) := \frac{1}{4} + x_1^4 - 2x_1^2x_2^2 + x_2^4$ . With  $\mathbf{m}_2(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$ , one solves the linear matrix feasibility problem  $r(\mathbf{x}) = \mathbf{m}_2(\mathbf{x})^\top \mathbf{Q} \mathbf{m}_2(\mathbf{x})$ . One can show that the solution writes  $\mathbf{Q} = \mathbf{L}^\top \mathbf{D} \mathbf{L}$  for a  $6 \times 6$  matrix  $\mathbf{L}$  and a diagonal matrix  $\mathbf{D}$  with entries  $(\frac{1}{2}, 0, 0, 1, 0, 0)$ , yielding the SOS decomposition:  $r(\mathbf{x}) = (\frac{1}{2})^2 + (x_1^2 - x_2^2)^2$ . This is enough to prove that  $p$  is nonnegative.

**2.2.2. Dense SDP Relaxations for Polynomial Optimization.** In order to solve our goal problem (Problem (1)), we are trying to solve Problem (6), recast as Problem (7). We first explain how to obtain tractable approximations of this difficult problem. Define  $g_0 := 1$ . The hierarchy of SDP relaxations developed by Lasserre [2001] provides lower bounds of  $r^*$ , through solving the optimization problems  $(\mathbf{P}_d)$ :

$$(\mathbf{P}_d) : \begin{cases} p_d^* := \sup_{\sigma_j, \mu} \mu, \\ \text{s.t. } r(\mathbf{x}) - \mu = \sum_{j=0}^k \sigma_j(\mathbf{x}) g_j(\mathbf{x}), \forall \mathbf{x}, \\ \mu \in \mathbb{R}, \sigma_j \in \Sigma[\mathbf{x}], \quad j = 0, \dots, k, \\ \deg(\sigma_j g_j) \leq 2d, \quad j = 0, \dots, k. \end{cases}$$

One can solve  $(\mathbf{P}_d)$  with SDP optimization to find a tuple  $(\mu, \sigma_0, \dots, \sigma_k)$ , which allows us to prove that  $g_1(\mathbf{x}) \geq 0 \wedge \dots \wedge g_k(\mathbf{x}) \geq 0 \Rightarrow r(\mathbf{x}) - \mu \geq 0$ .

The next theorem is a consequence of the assumption mentioned in Remark 2.1.

**THEOREM 2.3 (LASSERRE [2001]).** *Let  $p_d^*$  be the optimal value of the SDP relaxation  $(\mathbf{P}_d)$ . Then, the sequence of optimal values  $(p_d^*)_{d \in \mathbb{N}}$  is nondecreasing and converges to  $r^*$ .*

The number of SDP variables (i.e., the number of variables of the semidefinite relaxation  $(\mathbf{P}_d)$ ) grows polynomially with the integer  $d$ , called the *relaxation order*. Indeed,

at fixed number of variables  $n$ , the relaxation  $(\mathbf{P}_d)$  involves  $O((2d)^n)$  SDP variables and  $(k+1)$  linear matrix inequalities (LMIs) of size  $O(d^n)$ . When  $d$  increases, then more accurate lower bounds of  $r^*$  can be obtained, at an increasing computational cost. At fixed  $d$ , the relaxation  $(\mathbf{P}_d)$  involves  $O(n^{2d})$  SDP variables and  $(d+1)$  linear matrix inequalities (LMIs) of size  $O(n^d)$ .

*Example 2.4.* Consider the polynomial  $f$  mentioned in Section 1:  $f(\mathbf{x}) := x_2x_5 + x_3x_6 - x_2x_3 - x_5x_6 + x_1(-x_1 + x_2 + x_3 - x_4 + x_5 + x_6)$  and the set  $\mathbf{K} := [4, 6.36]^6$ . The set  $\mathbf{K}$  can be equivalently rewritten as:

$$\mathbf{K} := \{\mathbf{x} \in \mathbb{R}^n : g_1(\mathbf{x}) \geq 0, \dots, g_7(\mathbf{x}) \geq 0\},$$

with  $g_i(\mathbf{x}) := (6.36 - x_i)(x_i - 4)$  for each  $i = 1, \dots, 6$  and  $g_7(\mathbf{x}) := 243 - \sum_{i=0}^6 x_i^2$ . Here the constant  $M = 243$  is chosen so  $M \geq 6 \times 6.36^2$  and the assumption in Remark 2.1 is fulfilled. The number of initial variables of the optimization problem  $r^* := \inf_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x}) : \mathbf{x} \in \mathbf{K}\}$  is  $n = 6$  and the number of optimization constraints is  $k = 7$ . For  $d = 1$ , the dense SDP relaxation  $(\mathbf{P}_1)$  involves  $\binom{n+2d}{2d} = \binom{6+2}{2} = 28$  variables and provides a lower bound  $p_1^* = 20.755$  for  $r^*$ . The dense SDP relaxation  $(\mathbf{P}_2)$  involves  $\binom{6+4}{4} = 210$  variables and provides a tighter lower bound of  $p_2^* = 20.8608$  for  $r^*$ .

*2.2.3. Exploiting Sparsity.* Here we recall how to exploit the structured sparsity of the problem to replace one SDP problem  $(\mathbf{P}_d)$  by an SDP problem  $(\mathbf{S}_d)$  of size  $O(\kappa^{2d})$  where  $\kappa$  is the average size of the maximal cliques of the correlation sparsity pattern (csp) of the polynomial variables (see Waki et al. [2006] and Lasserre [2006] for more details). We now present these notions as well as the formulation of sparse SDP relaxations  $(\mathbf{S}_d)$ .

We denote by  $\mathbb{N}^n$  the set of  $n$ -tuple of nonnegative integers. The support of a polynomial  $r(\mathbf{x}) := \sum_{\alpha \in \mathbb{N}^n} r_\alpha \mathbf{x}^\alpha$  is defined as  $\text{supp}(r) := \{\alpha \in \mathbb{N}^n : r_\alpha \neq 0\}$ . For instance the support of  $r(\mathbf{x}) := \frac{1}{4} + x_1^4 - 2x_1^2x_2^2 + x_2^4$  is  $\text{supp}(p) = \{(0, 0), (4, 0), (2, 2), (0, 4)\}$ .

Let  $F_j$  be the index set of variables that are involved in the polynomial  $g_j$ , for each  $j = 1, \dots, k$ . The correlative sparsity is represented by the  $n \times n$  correlation sparsity pattern matrix (csp matrix)  $\mathbf{R}$  defined by

$$\mathbf{R}(i, j) := \begin{cases} 1 & \text{if } i = j, \\ 1 & \text{if } \exists \alpha \in \text{supp}(f) \text{ such that } \alpha_i, \alpha_j \geq 1, \\ 1 & \text{if } \exists l \in \{1, \dots, k\} \text{ such that } i, j \in F_l, \\ 0 & \text{otherwise.} \end{cases}$$

We define the undirected csp graph  $G(N, E)$  with  $N = \{1, \dots, n\}$  and  $E = \{\{i, j\} : i, j \in N, i < j, \mathbf{R}(i, j) = 1\}$ . Then, let  $C_1, \dots, C_m \subseteq N$  denote the maximal cliques of  $G(N, E)$  and define  $n_j := \#C_j$ , for each  $j = 1, \dots, m$ .

*Remark 2.5.* Assuming that the set  $\mathbf{K}$  is as in Remark 2.1, one replaces the constraint  $M - \sum_{i=1}^n x_i^2 \geq 0$  by the  $m$  redundant additional constraints,

$$g_{k+j} := n_j M^2 - \sum_{i \in C_j} x_i^2 \geq 0, \quad j = 1, \dots, m, \quad (9)$$

set  $k' = k + m$  and define the compact semialgebraic set,

$$\mathbf{K}' := \{\mathbf{x} \in \mathbb{R}^n : g_1(\mathbf{x}) \geq 0, \dots, g_{k'}(\mathbf{x}) \geq 0\},$$

and modify Problem (6) into the following optimization problem:

$$r^* := \inf_{\mathbf{x} \in \mathbb{R}^n} \{r(\mathbf{x}) : \mathbf{x} \in \mathbf{K}'\}. \quad (10)$$

For each  $j = 1, \dots, m$ , we note  $\mathbb{R}_{2d}[\mathbf{x}, C_j]$  the set of polynomials of  $\mathbb{R}_{2d}[\mathbf{x}]$  that involve the variables  $(x_i)_{i \in C_j}$ . We denote  $\Sigma[\mathbf{x}, C_j] := \Sigma[\mathbf{x}] \cap \mathbb{R}_{2d}[\mathbf{x}, C_j]$ . Similarly, we define

$\Sigma[\mathbf{x}, F_j]$  for each  $j = 1, \dots, k'$ . The following program is the sparse variant of the SDP program  $(\mathbf{P}_d)$ :

$$(\mathbf{S}_d) : \begin{cases} r_d^* := \sup_{\mu, \sigma_j} \mu, \\ \text{s.t. } r(\mathbf{x}) - \mu = \sum_{j=0}^{k'} \sigma_j(\mathbf{x})g_j(\mathbf{x}), \forall \mathbf{x}, \\ \mu \in \mathbb{R}, \sigma_0 \in \Sigma[\mathbf{x}, C_j], \\ \sigma_j \in \Sigma[\mathbf{x}, F_j], j = 1, \dots, k', \\ \deg(\sigma_j g_j) \leq 2d, j = 0, \dots, k', \end{cases}$$

where  $\sigma_0 \in \Sigma[\mathbf{x}, C_j]$  if and only if there exist  $\sigma^1 \in \Sigma[\mathbf{x}, C_1], \dots, \sigma^m \in \Sigma[\mathbf{x}, C_m]$  such that  $\sigma_0(\mathbf{x}) = \sum_{j=1}^m \sigma^j(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{R}^n$ .

The number of SDP variables of the relaxation  $(\mathbf{S}_d)$  is  $\sum_{j=1}^m \binom{n_j+2d}{2d}$ . At fixed  $d$ , it yields an SDP problem with  $O(\kappa^{2d})$  variables, where  $\kappa := \frac{1}{m} \sum_{j=1}^m n_j$  is the average size of the cliques  $C_1, \dots, C_m$ . Moreover, the cliques  $C_1, \dots, C_m$  satisfy the running intersection property.

**Definition 2.6 (RIP).** Let  $m \in \mathbb{N}_0$  and  $I_1, \dots, I_m$  be subsets of  $\{1, \dots, n\}$ . We say that  $I_1, \dots, I_m$  satisfy the running intersection property (RIP) when for all  $i = 1, \dots, m$ , there exists an integer  $l < i$  such that  $I_i \cap (\cup_{j < i} I_j) \subseteq I_l$ .

This RIP property together with the assumption mentioned in Remark 2.5 allow us to state the sparse variant of Theorem 2.3:

**THEOREM 2.7 (LASSERRE [2006, THEOREM 3.6]).** *Let  $r_d^*$  be the optimal value of the sparse SDP relaxation  $(\mathbf{S}_d)$ . Then the sequence  $(r_d^*)_{d \in \mathbb{N}}$  is nondecreasing and converges to  $r^*$ .*

The interested reader can find more details in Waki et al. [2006] about additional ways to exploit sparsity in order to derive analogous sparse SDP relaxations. We illustrate the benefits of the SDP relaxations  $(\mathbf{S}_d)$  with the following example:

**Example 2.8.** Consider the polynomial  $f$  mentioned in Section 1:  $f(\mathbf{x}) := x_2x_5 + x_3x_6 - x_2x_3 - x_5x_6 + x_1(-x_1 + x_2 + x_3 - x_4 + x_5 + x_6)$ . Here  $n = 6, d = 2, N = \{1, \dots, 6\}$ . The  $6 \times 6$  correlative sparsity matrix  $\mathbf{R}$  is

$$\mathbf{R} = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 & 1 \end{pmatrix}.$$

The csp graph  $G$  associated to  $\mathbf{R}$  is depicted in Figure 1. The maximal cliques of  $G$  are  $C_1 := \{1, 4\}$ ,  $C_2 := \{1, 2, 3\}$ ,  $C_3 := \{1, 2, 5\}$ ,  $C_4 := \{1, 5, 6\}$ , and  $C_5 := \{1, 3, 6\}$ . For  $d = 2$ , the dense SDP relaxation  $(\mathbf{P}_2)$  involves  $\binom{6+4}{4} = 210$  variables against  $\binom{2+4}{4} + 4\binom{3+4}{4} = 155$  for the sparse variant  $(\mathbf{S}_2)$ . The dense SDP relaxation  $(\mathbf{P}_3)$  involves 924 variables against 364 for the sparse variant  $(\mathbf{S}_3)$ . This difference becomes significant while considering that the time complexity of semidefinite programming is polynomial w.r.t. the number of variables with an exponent greater than 3 (see Ben-Tal and Nemirovski [2001, Chapter 4] for more details).

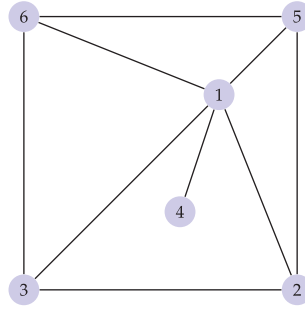


Fig. 1. Correlative sparsity pattern graph for the variables of  $f$  from Example 2.8.

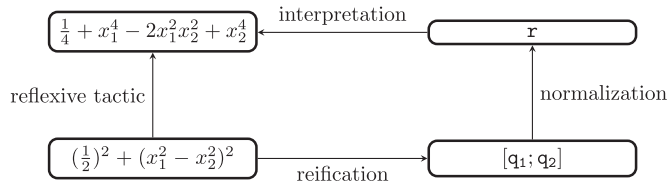


Fig. 2. An illustration of computational reflection.

### 2.3. Computer Proofs for Polynomial Optimization

Here, we briefly recall some existing features of the Coq proof assistant to handle formal polynomial optimization when using SDP relaxations. The advantage of such relaxations is that they provide SOS certificates, which can be formally checked *a posteriori*. For more details on Coq, we recommend the documentation available in Bertot and Castéran [2004]. Given a polynomial  $r$  and a set of constraints  $\mathbf{K}$ , one can obtain a lower bound on  $r$  by solving any instance of Problem  $(\mathbf{P}_d)$ . Then, one can verify formally the correctness of the lower bound  $r_d^*$ , using the SOS certificate output  $\sigma_0, \dots, \sigma_k$ . Indeed it is enough to prove the polynomial equality  $r(\mathbf{x}) - r_d^* = \sum_{j=0}^k \sigma_j(\mathbf{x})g_j(\mathbf{x})$  inside Coq. Such equalities can be efficiently proved using Coq's ring tactic [Grégoire and Mahboubi 2005] via the mechanism of computational reflection [Boutin 1997]. Any polynomial of type `pexprC` (see Section 2.1) can be normalized to a unique polynomial of type `polC` (see Grégoire and Mahboubi [2005] for more details on the constructors of this type). For the sake of clarity, let us consider the unconstrained case, that is,  $\mathbf{K} = \mathbb{R}^n$ . One encodes an SOS certificate  $\sigma_0(\mathbf{x}) = \sum_{i=1}^m q_i^2$  with the sequence of polynomials  $[q_1; \dots; q_m]$ , each  $q_i$  being of type `polC`. To prove the equality  $r = \sigma_0$ , our version of the ring tactic normalizes both  $r$  and the sequence  $[q_1; \dots; q_m]$  and compares the two normalization results. This mechanism is illustrated in Figure 2 with the polynomial  $r(\mathbf{x}) := \frac{1}{4} + x_1^4 - 2x_1^2x_2^2 + x_2^4$  (see Example 2.2) being encoded by  $r$  and the polynomials  $1/2$  and  $x_1^2 - x_2^2$  being encoded respectively by  $q_1$  and  $q_2$ .

In the general case, this computational step is done through a `checker_sos` procedure that returns a Boolean value. If this value is true, then one applies a correctness lemma, whose conclusion yields the nonnegativity of  $r - r_d^*$  over  $\mathbf{K}$ . In practice, the SDP solvers are implemented in floating-point arithmetic, thus the above equality between  $r - r_d^*$  and the SOS certificate does not hold. However, following Remark 2.1, each variable lies in a closed interval, and thus one can bound the remainder polynomial  $\epsilon(\mathbf{x}) := r(\mathbf{x}) - r_d^* - \sum_{j=0}^k \sigma_j(\mathbf{x})g_j(\mathbf{x})$  using basic interval arithmetic, so the lower bound  $\epsilon^*$  of  $\epsilon$  yields the valid inequality:  $\forall \mathbf{x} \in \mathbf{K}, r(\mathbf{x}) \geq r_d^* + \epsilon^*$ . For more explanation, we refer the interested reader to the formal framework [Magron et al. 2015b, Section 2.3]. Note that this formal verification remains valid when considering the sparse variant  $(\mathbf{S}_d)$ .

**Input:** input variables  $\mathbf{x}$ , input constraints  $\mathbf{X}$ , nonlinear expression  $f$ , rounded expression  $\hat{f}$ , error variables  $\mathbf{e}$ , error constraints  $\mathbf{E}$ , relaxation order  $d$

**Output:** interval enclosure  $I_d$  of the error  $\hat{f} - f$  over  $\mathbf{K} := \mathbf{X} \times \mathbf{E}$

- 1: Define the absolute error  $r(\mathbf{x}, \mathbf{e}) := \hat{f}(\mathbf{x}, \mathbf{e}) - f(\mathbf{x})$
- 2: Compute  $l(\mathbf{x}, \mathbf{e}) := r(\mathbf{x}, 0) + \sum_{j=1}^m \frac{\partial r(\mathbf{x}, \mathbf{e})}{\partial e_j}(\mathbf{x}, 0) e_j$
- 3: Define  $h := r - l$
- 4: Compute bounds for  $h$ :  $I^h := \text{ia\_bound}(h, \mathbf{K})$
- 5: Compute bounds for  $l$ :  $I_d^l := \text{sdp\_bound}(l, \mathbf{K}, d)$
- 6: **return**  $I_d := I_d^l + I^h$

Fig. 3. bound: our algorithm to compute roundoff errors bounds of nonlinear programs.

### 3. GUARANTEED ROUNDOFF ERROR BOUNDS USING SDP RELAXATIONS

In this section, we present our new algorithm, relying on sparse SDP relaxations, to bound roundoff errors of nonlinear programs. After stating our general algorithm (Section 3.1), we detail how this procedure can handle polynomial programs (Section 3.2). Extensions to the non-polynomial case are presented in Section 3.3.

#### 3.1. The General Optimization Framework

Here we consider a given program that implements a nonlinear transcendental expression  $f$  with input variables  $\mathbf{x}$  satisfying a set of constraints  $\mathbf{X}$ . We assume that  $\mathbf{X}$  is included in a box (i.e., a product of closed intervals)  $[\mathbf{a}, \mathbf{b}] := [a_1, b_1] \times \dots \times [a_n, b_n]$  and that  $\mathbf{X}$  is encoded as follows:

$$\mathbf{X} := \{ \mathbf{x} \in \mathbb{R}^n : g_1(\mathbf{x}) \geq 0, \dots, g_k(\mathbf{x}) \geq 0 \},$$

for polynomial functions  $g_1, \dots, g_k$ . Then, we denote by  $\hat{f}(\mathbf{x}, \mathbf{e})$  the rounded expression of  $f$  after applying the round procedure (see Section 2.1), introducing additional error variables  $\mathbf{e}$ .

The algorithm bound, depicted in Figure 3, takes as input  $\mathbf{x}, \mathbf{X}, f, \hat{f}, \mathbf{e}$  as well as the set  $\mathbf{E}$  of bound constraints over  $\mathbf{e}$ . Here we assume that our program implementing  $f$  does not involve conditional statements (this case will be discussed later in Section 3.3). For a given machine  $\epsilon$ , one has  $\mathbf{E} := [-\epsilon, \epsilon]^m$ , with  $m$  being the number of error variables. This algorithm actually relies on the sparse SDP optimization procedure ( $\mathbf{S}_d$ ) (see Section 2.2 for more details), thus bound also takes as input a relaxation order  $d \in \mathbb{N}$ . The algorithm provides as output an interval enclosure  $I_d$  of the error  $\hat{f}(\mathbf{x}, \mathbf{e}) - f(\mathbf{x})$  over  $\mathbf{K}$ . From this interval  $I_d := [\underline{f}_d, \overline{f}_d]$ , one can compute  $f_d := \max\{-\underline{f}_d, \overline{f}_d\}$ , which is a sound upper bound of the maximal absolute error  $r^* := \max_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} |\hat{f}(\mathbf{x}, \mathbf{e}) - f(\mathbf{x})|$ .

After defining the absolute roundoff error  $r := \hat{f} - f$  (Line 1), one decomposes  $r$  as the sum of an expression  $l$  that is affine w.r.t. the error variable  $\mathbf{e}$  and a remainder  $h$ . One way to obtain  $l$  is to compute the vector of partial derivatives of  $r$  w.r.t.  $\mathbf{e}$  evaluated at  $(\mathbf{x}, 0)$  and, finally, to take the inner product of this vector and  $\mathbf{e}$  (Line 2). Then, the idea is to compute a precise bound of  $l$  and a coarse bound of  $h$ . The underlying reason is that  $h$  involves error term products of degree greater than 2 (e.g.,  $e_1 e_2$ ), yielding an interval enclosure  $I^h$  of *a priori* much smaller width, compared to the interval enclosure  $I^l$  of  $l$ . One obtains  $I^h$  using the procedure `ia_bound` implementing basic interval arithmetic (Line 4) to bound the remainder of the multivariate Taylor expansion of  $r$  w.r.t.  $\mathbf{e}$ , expressed as a combination of the second-order derivatives (similar as in Solovyev et al. [2015]). The main algorithm presented in Figure 3 is very similar to the algorithm of `FPTaylor` [Solovyev et al. 2015], except that SDP-based techniques are used instead of the global optimization procedure from Solovyev et al. [2015]. Note that overflow and denormal are neglected here but one could handle them, as in Solovyev et al. [2015],



while adding additional error variables and discarding the related terms using naive interval arithmetic.

### 3.2. Polynomial Programs

We first describe our `sdp_bound` optimization algorithm when implementing polynomial programs. In this case, `sdp_bound` calls an auxiliary procedure `sdp_poly`. The bound of  $l$  is provided through solving two sparse SDP instances of Problem  $(\mathbf{S}_d)$  at relaxation order  $d$ . We now give more explanation about the `sdp_poly` procedure.

We can map each input variable  $x_i$  to the integer  $i$ , for all  $i = 1, \dots, n$ , as well as each error variable  $e_j$  to  $n + j$  for all  $j = 1, \dots, m$ . Then we define the sets  $C_1 := \{1, \dots, n, n + 1\}, \dots, C_m := \{1, \dots, n, n + m\}$ . Here we take advantage of the correlation sparsity pattern of  $l$  by using  $m$  distinct sets of cardinality  $n + 1$  rather than a single one of cardinality  $n + m$ , that is, the total number of variables. After writing  $l(\mathbf{x}, \mathbf{e}) = r(\mathbf{x}, 0) + \sum_{j=1}^m \frac{\partial r(\mathbf{x}, \mathbf{e})}{\partial e_j}(\mathbf{x}, 0) e_j$  and noticing that  $r(\mathbf{x}, 0) = \hat{f}(\mathbf{x}, 0) - f(\mathbf{x}) = 0$ , one can scale the optimization problems by writing

$$l(\mathbf{x}, \mathbf{e}) = \sum_{j=1}^m s_j(\mathbf{x}) e_j = \epsilon \sum_{j=1}^m s_j(\mathbf{x}) \frac{e_j}{\epsilon}, \quad (11)$$

with  $s_j(\mathbf{x}) := \frac{\partial r(\mathbf{x}, \mathbf{e})}{\partial e_j}(\mathbf{x}, 0)$  for all  $j = 1, \dots, m$ . Replacing  $\mathbf{e}$  by  $\mathbf{e}/\epsilon$  leads to computing an interval enclosure of  $l/\epsilon$  over  $\mathbf{K}' := \mathbf{X} \times [-1, 1]^m$ . Recall that from Remark 2.1, there exists an integer  $M > 0$  such that  $M - \sum_{i=1}^n x_i^2 \geq 0$ , as the input variables satisfy box constraints. Moreover, to fulfill the assumption of Remark 2.5, one encodes  $\mathbf{K}'$  as follows:

$$\mathbf{K}' := \{(\mathbf{x}, \mathbf{e}) \in \mathbb{R}^{n+m} : g_1(\mathbf{x}) \geq 0, \dots, g_k(\mathbf{x}) \geq 0, \\ g_{k+1}(\mathbf{x}, e_1) \geq 0, \dots, g_{k+m}(\mathbf{x}, e_m) \geq 0\},$$

with  $g_{k+j}(\mathbf{x}, e_j) := M + 1 - \sum_{i=1}^n x_i^2 - e_j^2$  for all  $j = 1, \dots, m$ . The index set of variables involved in  $g_j$  is  $F_j := N = \{1, \dots, n\}$  for all  $j = 1, \dots, k$ . The index set of variables involved in  $g_{k+j}$  is  $F_{k+j} := C_j$  for all  $j = 1, \dots, m$ .

Then one can compute a lower bound of the minimum of  $l'(\mathbf{x}, \mathbf{e}) := l(\mathbf{x}, \mathbf{e})/\epsilon = \sum_{j=1}^m s_j(\mathbf{x}) e_j$  over  $\mathbf{K}'$  by solving the following optimization problem:

$$\left\{ \begin{array}{l} \underline{l}'_d := \sup_{\mu, \sigma_j} \mu, \\ \text{s.t. } l' - \mu = \sigma_0 + \sum_{j=1}^{k+m} \sigma_j g_j, \\ \mu \in \mathbb{R}, \sigma_0 \in \sum_{j=1}^m \Sigma[(\mathbf{x}, \mathbf{e}), C_j], \\ \sigma_j \in \Sigma[(\mathbf{x}, \mathbf{e}), F_j], \quad j = 1, \dots, k + m, \\ \deg(\sigma_j g_j) \leq 2d, \quad j = 1, \dots, k + m. \end{array} \right. \quad (12)$$

A feasible solution of Problem (12) ensures the existence of  $\sigma^1 \in \Sigma[(\mathbf{x}, e_1)], \dots, \sigma^m \in \Sigma[(\mathbf{x}, e_m)]$  such that  $\sigma_0 = \sum_{j=0}^m \sigma^j$ , allowing the following reformulation:

$$\left\{ \begin{array}{l} \underline{l}'_d := \sup_{\mu, \sigma_j} \mu, \\ \text{s.t. } l' - \mu = \sum_{j=1}^m \sigma^j + \sum_{j=1}^{k+m} \sigma_j g_j, \\ \mu \in \mathbb{R}, \sigma_j \in \Sigma[\mathbf{x}], \quad j = 1, \dots, m, \\ \sigma^j \in \Sigma[(\mathbf{x}, e_j)], \deg(\sigma^j) \leq 2d, \quad j = 1, \dots, m, \\ \deg(\sigma_j g_j) \leq 2d, \quad j = 1, \dots, k + m. \end{array} \right. \quad (13)$$

An upper bound  $\bar{l}'_d$  can be obtained by replacing  $\sup$  with  $\inf$  and  $l' - \mu$  by  $\mu - l'$  in Problem (13). Our optimization procedure `sdp_poly` computes the lower bound  $\underline{l}'_d$  as well as an upper bound  $\bar{l}'_d$  of  $l'$  over  $\mathbf{K}'$  then returns the interval  $I'_d := [\epsilon \underline{l}'_d, \epsilon \bar{l}'_d]$ , which is a sound enclosure of the values of  $l$  over  $\mathbf{K}$ .

We emphasize two advantages of the decomposition  $r := l + h$  and, more precisely, of the linear dependency of  $l$  w.r.t.  $\mathbf{e}$ : scalability and robustness to SDP numerical issues. First, no computation is required to determine the correlation sparsity pattern of  $l$  in comparison to the general case. Thus, it becomes much easier to handle the optimization of  $l$  with the sparse SDP Problem (13) rather than with the corresponding instance of the dense relaxation ( $\mathbf{P}_d$ ). While the latter involves  $\binom{n+m+2d}{2d}$  SDP variables, the former involves only  $m \binom{n+1+2d}{2d}$  variables, ensuring the scalability of our framework. In addition, the linear dependency of  $l$  w.r.t.  $\mathbf{e}$  allows us to scale the error variables and optimize over a set of variables lying in  $\mathbf{K}' := \mathbf{X} \times [-1, 1]^m$ . It ensures that the range of input variables does not significantly differ from the range of error variables. This condition is mandatory while considering SDP relaxations because most SDP solvers (e.g., MOSEK [Andersen and Andersen 2000]) are implemented using double precision floating-point. It is impossible to optimize  $l$  over  $\mathbf{K}$  (rather than  $l'$  over  $\mathbf{K}'$ ) when the maximal value  $\epsilon$  of error variables is less than  $2^{-53}$ , due to the fact that SDP solvers would treat each error variable term as 0 and, consequently,  $l$  as the zero polynomial. Thus, this decomposition insures our framework against numerical issues related to finite-precision implementation of SDP solvers.

Let us define the interval enclosure  $I^l := [\underline{l}, \bar{l}]$ , with  $\underline{l} := \inf_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} l(\mathbf{x}, \mathbf{e})$  and  $\bar{l} := \sup_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} l(\mathbf{x}, \mathbf{e})$ . The next lemma states that one can approximate  $I^l$  as closely as desired using the `sdp_poly` procedure.

**LEMMA 3.1 (CONVERGENCE OF THE `sdp_poly` PROCEDURE).** *Let  $I^l_d$  be the interval enclosure returned by the procedure `sdp_poly`( $l, \mathbf{K}, d$ ). The sequence  $(I^l_d)_{d \in \mathbb{N}}$  converges to  $I^l$ .*

**PROOF.** It is sufficient to show the similar convergence result for  $l' = l/\epsilon$ , as it implies the convergence for  $l$  by a scaling argument. The sets  $C_1, \dots, C_m$  satisfy the RIP property (see Definition 2.6). Moreover, the encoding of  $\mathbf{K}'$  satisfies the assumption mentioned in Remark 2.5. Thus, Theorem 2.7 implies that the sequence of lower bounds  $(\underline{l}'_d)_{d \in \mathbb{N}}$  converges to  $\underline{l}' := \inf_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}'} l'(\mathbf{x}, \mathbf{e})$ . Similarly, the sequence of upper bounds converge to  $\bar{l}'$ , yielding the desired result.  $\square$

Lemma 3.1 guarantees asymptotic convergence to the exact enclosure of  $l$  when the relaxation order  $d$  tends to infinity. However, it is more reasonable in practice to keep this order as small as possible to obtain tractable SDP relaxations. Hence, we generically solve each instance of Problem (13) at the minimal relaxation order, that is,  $d_0 := \max\{\lceil \deg l / 2 \rceil, \max_{1 \leq j \leq k+m} \lceil \deg(g_j) / 2 \rceil\}$ .

### 3.3. Non-Polynomial and Conditional Programs

Other classes of programs do not only involve polynomials but also semialgebraic and transcendental functions as well as conditional statements. Such programs are of particular interest as they often occur in real-world applications such as biology modeling, space control, or global optimization. We present how the general optimization procedure `sdp_bound` can be extended to these nonlinear programs.

**3.3.1. Semialgebraic Programs.** Here we assume that the function  $l$  is semialgebraic, that is, it involves non-polynomial components such as divisions or square roots. Following Lasserre and Putinar [2010], we explain how to transform the optimization problem  $\inf_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}} l(\mathbf{x}, \mathbf{e})$  into a polynomial optimization problem and then use the sparse SDP

**Input:** input variables  $\mathbf{y}$ , input constraints  $\mathbf{K}$ , semialgebraic expression  $f$   
**Output:** variables  $\mathbf{y}_{\text{poly}}$ , constraints  $\mathbf{K}_{\text{poly}}$ , polynomial expression  $f_{\text{poly}}$

```

1:  $I := \text{ia\_bound}(f, \mathbf{K})$ 
2: if  $f = \text{Pol}(p)$  then  $\mathbf{y}_{\text{poly}} := \mathbf{y}$ ,  $\mathbf{K}_{\text{poly}} := \mathbf{K}$ ,  $f_{\text{poly}} := p$ 
3: else if  $f = \text{Div}(g, h)$  then
4:    $\mathbf{y}_g, \mathbf{K}_g, g_{\text{poly}} := \text{lift}(\mathbf{y}, \mathbf{K}, g)$ 
5:    $\mathbf{y}_h, \mathbf{K}_h, h_{\text{poly}} := \text{lift}(\mathbf{y}, \mathbf{K}, h)$ 
6:    $\mathbf{y}_{\text{poly}} := (\mathbf{y}_g, \mathbf{y}_h, x)$     $f_{\text{poly}} := x$ 
7:    $\mathbf{K}_{\text{poly}} := \{\mathbf{y}_{\text{poly}} \in \mathbf{K}_g \times \mathbf{K}_h \times I : x h_{\text{poly}} = g_{\text{poly}}\}$ 
8: else if  $f = \text{Sqrt}(g)$  then
9:    $\mathbf{y}_g, \mathbf{K}_g, g_{\text{poly}} := \text{lift}(\mathbf{y}, \mathbf{K}, g)$ 
10:   $\mathbf{y}_{\text{poly}} := (\mathbf{y}_g, x)$     $f_{\text{poly}} := x$ 
11:   $\mathbf{K}_{\text{poly}} := \{\mathbf{y}_{\text{poly}} \in \mathbf{K}_g \times I : x^2 = g_{\text{poly}}\}$ 
12: ...
13: end
14: return  $\mathbf{y}_{\text{poly}}, \mathbf{K}_{\text{poly}}, f_{\text{poly}}$ 

```

Fig. 4. lift: a recursive procedure to reduce semialgebraic problems to polynomial problems.

program (13). One way to perform this reformulation consists of introducing lifting variables to represent non-polynomial operations. We first illustrate the extension to semialgebraic programs with an example.

*Example 3.2.* Let us consider the program implementing the rational function  $f : [0, 1] \rightarrow \mathbb{R}$  defined by  $f(x_1) := \frac{x_1}{1+x_1}$ . Applying the rounding procedure (with machine  $\epsilon$ ) yields  $\hat{f}(x_1, \mathbf{e}) := \frac{x_1(1+e_2)}{(1+x_1)(1+e_1)}$  and the decomposition  $r(x_1, \mathbf{e}) := \hat{f}(x_1, \mathbf{e}) - f(x_1) = l(x_1, \mathbf{e}) + h(x_1, \mathbf{e}) = s_1(x_1)e_1 + s_2(x_1)e_2 + h(x_1, \mathbf{e})$ . One has  $s_1(x_1) = \frac{\partial r(x_1, \mathbf{e})}{\partial e_1}(x_1, 0) = -\frac{x_1}{1+x_1}$  and  $s_2(x_1) = -s_1(x_1)$ .

Let  $\mathbf{K} := [0, 1] \times [-\epsilon, \epsilon]^2$ . One introduces a lifting variable  $x_2 := \frac{x_1}{1+x_1}$  to handle the division operator and encode the equality constraint  $p(\mathbf{x}) := x_2(1+x_1) - x_1 = 0$  with the two inequality constraints  $p(\mathbf{x}) \geq 0$  and  $-p(\mathbf{x}) \geq 0$ . To ensure the compactness assumption, one bounds  $x_2$  within  $I := [0, 1/2]$ , using basic interval arithmetic.

Let  $\mathbf{K}_{\text{poly}} := \{(\mathbf{x}, \mathbf{e}) \in [0, 1] \times I \times [\epsilon, \epsilon]^2 : p(\mathbf{x}) \geq 0, -p(\mathbf{x}) \geq 0\}$ . Then the rational optimization problem involving  $l$  is equivalent to  $\inf_{(\mathbf{x}, \mathbf{e}) \in \mathbf{K}_{\text{poly}}} x_2(-e_1 + e_2)$ , a polynomial optimization problem that we can handle with the `sdp.poly` procedure, described in Section 3.2.

In the semialgebraic case, `sdp_bound` calls an auxiliary procedure `sdp_sa`. Given input variables  $\mathbf{y} := (\mathbf{x}, \mathbf{e})$ , input constraints  $\mathbf{K} := \mathbf{X} \times \mathbf{E}$ , and a semialgebraic function  $l$ , `sdp_sa` first applies a recursive procedure `lift` that returns variables  $\mathbf{y}_{\text{poly}}$ , constraints  $\mathbf{K}_{\text{poly}}$ , and a polynomial  $f_{\text{poly}}$  such that the interval enclosure  $I^l$  of  $l(\mathbf{y})$  over  $\mathbf{K}$  is equal to the interval enclosure of the polynomial  $l_{\text{poly}}(\mathbf{y}_{\text{poly}})$  over  $\mathbf{K}_{\text{poly}}$ . Calling `sdp_sa` yields the interval enclosure  $I_d^l := \text{sdp\_poly}(l_{\text{poly}}, \mathbf{K}_{\text{poly}}, d)$ . We detail the lifting procedure `lift` in Figure 4 for the constructors `Pol` (Line (2)), `Div` (Line (3)), and `Sqrt` (Line (8)). The interval  $I$  obtained through the `ia_bound` procedure (Line (1)) allows us to constrain the additional variable  $x$  to ensure the assumption of Remark 2.5. For the sake of consistency, we omit the other cases (`Neg`, `Add`, `Mul`, and `Sub`) where the procedure is straightforward. For a similar procedure in the context of global optimization, we refer the interested reader to Magron [2013, Chapter 2]. The set of variables  $\mathbf{y}_{\text{poly}}$  can be decomposed as  $(\mathbf{x}_{\text{poly}}, \mathbf{e})$ , where  $\mathbf{x}_{\text{poly}}$  gathers input variables with lifting variables and has a cardinality equal to  $n_{\text{poly}}$ . Then, one easily shows that the sets  $\{1, \dots, n_{\text{poly}}, e_1\}, \dots, \{1, \dots, n_{\text{poly}}, e_m\}$

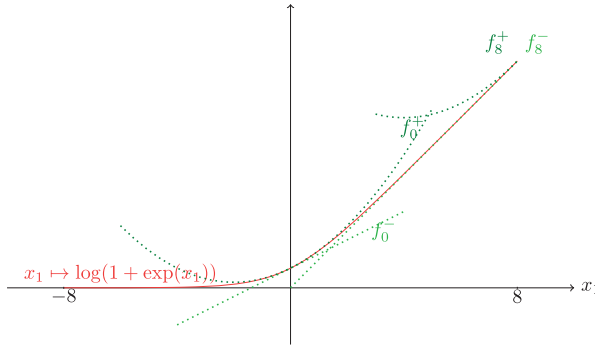


Fig. 5. Semialgebraic Approximations for  $x \mapsto \log(1 + \exp(x))$ :  $\max\{f_0^-(x), f_8^-(x)\} \leq \log(1 + \exp(x)) \leq \min\{f_0^+(x), f_8^+(x)\}$ .

satisfy the RIP, thus ensuring that we solve efficiently the corresponding instances of Problem (13).

**3.3.2. Transcendental Programs.** The above lifting procedure allows to represent exactly the graph of a semialgebraic function with polynomials involving additional (lifting) variables. Here we consider a procedure to approximate transcendental functions with semialgebraic functions. Here we assume that the function  $l$  is transcendental; that is, it involves univariate non-semialgebraic components such as  $\exp$  or  $\sin$ . We use the method presented in Magron et al. [2015a], based on maxplus approximation of semi-convex transcendental functions by quadratic functions. This idea comes from optimal control [McEneaney 2006] and was developed further to represent the value function by a “maxplus linear combination,” which is a supremum of quadratic polynomials at given points  $x_i$ . Given a set of points  $(x_i)$ , we approximate from above and from below every transcendental function  $f_{\mathbb{R}}$  by infima and suprema of finitely many quadratic polynomials ( $f_{x_i}^-$ ) and ( $f_{x_i}^+$ ). Hence, we reduce the problem to semialgebraic optimization problems. We can interpret this method in a geometrical way by thinking of it in terms of “quadratic cuts,” since quadratic inequalities are added to approximate the graph of a transcendental function.

For each univariate transcendental function  $f_{\mathbb{R}}$  in our dictionary set  $\mathcal{D}$ , one assumes that  $f_{\mathbb{R}}$  is twice differentiable, so the univariate function  $g := f_{\mathbb{R}} + \frac{\gamma}{2} |\cdot|^2$  is convex on  $I$  for large-enough  $\gamma > 0$  (for more details, see McEneaney [2006]). It follows that there exists a constant  $\gamma \leq \sup_{x \in I} -f_{\mathbb{R}}''(x)$  such that for all  $x_i \in I$ :

$$\begin{aligned} \forall x \in I, \quad f_{\mathbb{R}}(x) &\geq f_{x_i}^-(x), \\ \text{with } f_{x_i}^- &:= -\frac{\gamma}{2}(x - x_i)^2 + f'_{\mathbb{R}}(x_i)(x - x_i) + f_{\mathbb{R}}(x_i), \end{aligned} \quad (14)$$

implying that for all  $x \in I$ ,  $f_{\mathbb{R}}(x) \geq \max_{x_i \in I} f_{x_i}^-(x)$ . Similarly, one obtains an upper-approximation  $\min_{x_i \in I} f_{x_i}^+(x)$ . Figure 5 provides such approximations for the function  $f_{\mathbb{R}}(x) := \log(1 + \exp(x))$  on the interval  $I := [-8, 8]$ .

For transcendental programs, our procedure `sdp_bound` calls the auxiliary procedure `sdp_transc`. Given input variables  $(\mathbf{x}, \mathbf{e})$ , constraints  $\mathbf{K}$ , and a transcendental function  $l$ , `sdp_transc` first computes a semialgebraic lower (respectively, upper) approximation  $l^-$  (respectively,  $l^+$ ) of  $l$  over  $\mathbf{K}$ . For more details in the context of global optimization, we refer the reader to Magron et al. [2015a]. Then, calling the procedure `sdp_sa` allows us to get interval enclosures of  $l^-$  as well as  $l^+$ . We illustrate the procedure to handle transcendental programs with an example.

**Input:** input variables  $\mathbf{x}$ , input constraints  $\mathbf{X}$ , nonlinear expression  $f$ , rounded expression  $\hat{f}$ , error variables  $\mathbf{e}$ , error constraints  $\mathbf{E}$ , relaxation order  $d$

**Output:** interval enclosure  $I_d$  of the error  $\hat{f} - f$  over  $\mathbf{K} := \mathbf{X} \times \mathbf{E}$

```

1: if  $f = \text{IfThenElse}(p, g, h)$  then
2:    $I_d^p := \text{bound}(\mathbf{x}, \mathbf{X}, p, \hat{p}, \mathbf{e}, \mathbf{E}, d) = [\underline{p_d}, \overline{p_d}]$ 
3:    $\mathbf{X}_1 := \{\mathbf{x} \in \mathbf{X} : 0 \leq p(\mathbf{x}) \leq -\underline{p_d}\}$ 
4:    $\mathbf{X}_2 := \{\mathbf{x} \in \mathbf{X} : -\overline{p_d} \leq p(\mathbf{x}) \leq 0\}$ 
5:    $\mathbf{X}_3 := \{\mathbf{x} \in \mathbf{X} : 0 \leq p(\mathbf{x})\}$ 
6:    $\mathbf{X}_4 := \{\mathbf{x} \in \mathbf{X} : p(\mathbf{x}) \leq 0\}$ 
7:    $I_d^1 := \text{bound\_nlprog}(\mathbf{x}, \mathbf{X}_1, g, \hat{h}, \mathbf{e}, \mathbf{E}, d)$ 
8:    $I_d^2 := \text{bound\_nlprog}(\mathbf{x}, \mathbf{X}_2, h, \hat{g}, \mathbf{e}, \mathbf{E}, d)$ 
9:    $I_d^3 := \text{bound\_nlprog}(\mathbf{x}, \mathbf{X}_3, g, \hat{g}, \mathbf{e}, \mathbf{E}, d)$ 
10:   $I_d^4 := \text{bound\_nlprog}(\mathbf{x}, \mathbf{X}_4, h, \hat{h}, \mathbf{e}, \mathbf{E}, d)$ 
11:  return  $I_d := I_d^1 \cup I_d^2 \cup I_d^3 \cup I_d^4$ 
12: else return  $I_d := \text{bound}(\mathbf{x}, \mathbf{X}, f, \hat{f}, \mathbf{e}, \mathbf{E}, d)$ 
13: end

```

Fig. 6. bound\_nlprog: our algorithm to compute roundoff error bounds of programs with conditional statements.

*Example 3.3.* Let us consider the program implementing the transcendental function  $f : [-8, 8] \rightarrow \mathbb{R}$  defined by  $f(x_1) := \log(1 + \exp(x_1))$ . Applying the rounding procedure yields  $\hat{f}(x_1, \mathbf{e}) := \log[(1 + \exp(x_1)(1 + e_1))(1 + e_2)](1 + e_3)$ . Here,  $|e_2|$  is bounded by the machine  $\epsilon$  while  $|e_1|$  (respectively,  $|e_3|$ ) is bounded with an adjusted absolute error  $\epsilon_1 := \epsilon(\exp)$  (respectively,  $\epsilon_3 := \epsilon(\log)$ ). Let  $\mathbf{K} := [-8, 8] \times [-\epsilon_1, \epsilon_1] \times [-\epsilon, \epsilon] \times [-\epsilon_3, \epsilon_3]$ .

One obtains the decomposition  $r(x_1, \mathbf{e}) := \hat{f}(x_1, \mathbf{e}) - f(x_1) = l(x_1, \mathbf{e}) + h(x_1, \mathbf{e}) = s_1(x_1)e_1 + s_2(x_1)e_2 + s_3(x_1)e_3 + h(x_1, \mathbf{e})$ , with  $s_1(x_1) = \frac{\exp(x_1)}{1 + \exp(x_1)}$ ,  $s_2(x_1) = 1$  and  $s_3(x_1) = \log(1 + \exp(x_1)) = f(x_1)$ . Figure 5 provides a lower approximation  $s_3^- := \max\{f_0^-, f_8^-\}$  of  $s_3$  as well as an upper approximation  $s_3^+ := \min\{f_0^+, f_8^+\}$ . One can get similar approximations  $s_1^-$  and  $s_1^+$  for  $s_1$ . One first obtains (coarse) interval enclosures  $I_2 = \text{ia\_bound}(s_1, \mathbf{K})$  and  $I_3 = \text{ia\_bound}(s_3, \mathbf{K})$  and one introduces extra variables  $x_2 \in I_2$  and  $x_3 \in I_3$  to represent  $s_1$  and  $s_3$ , respectively. Then the interval enclosure of  $l$  over  $\mathbf{K}$  is equal to the interval enclosure of  $l_{\text{sa}}(\mathbf{x}, \mathbf{e}) := x_2e_1 + e_2 + x_3e_3$  over the set  $\mathbf{K}_{\text{sa}} := \{(x_1, \mathbf{e}) \in \mathbf{K}, (x_2, x_3) \in I_2 \times I_3, s_1^-(x_1) \leq x_2 \leq s_1^+(x_1), s_3^-(x_1) \leq x_3 \leq s_3^+(x_1)\}$ .

**3.3.3. Programs with Conditionals.** Finally, we explain how to extend our bounding procedure to nonlinear programs involving conditionals through the recursive algorithm given in Figure 6. The overall procedure is very similar to the one implemented within the Rosa tool [Darulova and Kuncak 2014, Section 7, Figure 6]. The bound\_nlprog algorithm relies on the bound procedure (see Figure 3 in Section 3.1) to compute roundoff error bounds of programs implementing transcendental functions (Line 12). From Line 1 to Line 11, the algorithm handles the case when the program implements a function  $f$  defined as follows:

$$f(\mathbf{x}) := \begin{cases} g(\mathbf{x}) & \text{if } p(\mathbf{x}) \geq 0, \\ h(\mathbf{x}) & \text{otherwise.} \end{cases}$$

The first branch output is  $g$  while the second one is  $h$ . More sophisticated conditionals, such as “ $p_1(x) \geq 0$  or/and  $p_2(x) \geq 0$ ,” are not handled at the moment, but one could easily extend the current framework to do so. A preliminary step consists of computing the roundoff error enclosure  $I_d^p := [\underline{p_d}, \overline{p_d}]$  (Line 2) for the program implementing the



polynomial  $p$ . Then the procedure computes bounds related to the branches discontinuity error, that is, the maximal value between the four following errors, as follows:

- (Line 7) the error obtained while computing the rounded result  $\hat{h}$  of the second branch instead of computing the exact result  $g$  of the first one, occurring for the set of variables  $(\mathbf{x}, \mathbf{e})$  such that  $\hat{p}(\mathbf{x}, \mathbf{e}) \leq 0 \leq p(\mathbf{x})$ . For scalability and numerical issues, we consider an over-approximation  $\mathbf{X}_1$  (Line 3) of this set, where the variables  $\mathbf{x}$  satisfy the relaxed constraints  $0 \leq p(\mathbf{x}) \leq -\underline{p}_d$ . Note that, in this case, one has  $l(\mathbf{x}, \mathbf{e}) = r(\mathbf{x}, 0) + \sum_{j=1}^m \frac{\partial r(\mathbf{x}, \mathbf{e})}{\partial e_j}(\mathbf{x}, 0) e_j$ , with  $r(\mathbf{x}, 0) = h(\mathbf{x}) - g(\mathbf{x}) \neq 0$ . In general, we expect the magnitude of the partial derivative sum to be very small compared to the one of  $r(\mathbf{x}, 0)$ .
- (Line 8) the error obtained while computing the rounded result  $\hat{g}$  of the first branch instead of computing the exact result  $h$  of the second one, occurring for the set of variables  $(\mathbf{x}, \mathbf{e})$  such that  $p(\mathbf{x}) \leq 0 \leq \hat{p}(\mathbf{x}, \mathbf{e})$ . We also consider an over-approximation  $\mathbf{X}_2$  (Line 4), where the variables  $\mathbf{x}$  satisfy the relaxed constraints  $-\overline{p}_d \leq p(\mathbf{x}) \leq 0$ .
- (Line 9) the roundoff error corresponding to the program implementation of  $g$ .
- (Line 10) the roundoff error corresponding to the program implementation of  $h$ .

**3.3.4. Simplification of Error Terms.** In addition, our algorithm `bound_nlprog` integrates several features to reduce the number of error variables. First, it memorizes all sub-expressions of the nonlinear expression tree to perform common sub-expressions elimination. We can also simplify error term products, thanks to the following lemma.

**LEMMA 3.4** (HIGHAM [2002, LEMMA 3.3]). *Let  $\epsilon$  be the machine precision and assume that, for a given integer  $k$ , one has  $\epsilon < \frac{1}{k}$  and  $\gamma_k := \frac{k\epsilon}{1-k\epsilon}$ . Then, for all  $e_1, \dots, e_k \in [-\epsilon, \epsilon]$ , there exists  $\theta_k$  such that  $\prod_{i=1}^k (1 + e_i) = 1 + \theta_k$  and  $|\theta_k| \leq \gamma_k$ .*

Lemma 3.4 implies that for any  $k$  such that  $\epsilon < \frac{1}{k}$ , one has  $\theta_k \leq (k+1)\epsilon$ . Our algorithm has an option to automatically derive safe over-approximations of the absolute roundoff error while introducing only one variable  $e_1$  (bounded by  $(k+1)\epsilon$ ) instead of  $k$  error variables  $e_1, \dots, e_k$  (bounded by  $\epsilon$ ). The cost of solving the corresponding optimization problem can be significantly reduced but it yields coarser error bounds.

## 4. EXPERIMENTAL EVALUATION

Now we present experimental results obtained by applying our general `bound_nlprog` algorithm (see Section 3, Figure 6) to various examples coming from physics, biology, space control, and optimization. The `bound_nlprog` algorithm is implemented in an open-source tool called `Real2Float`, built on top of the `NLCertify` nonlinear verification package, relying on `OCAML` (Version 4.02.1) and `Coq` (Version 8.4pl5) and interfaced with the SDP solver `SDPA` (Version 7.3.9). The SDP solver output numerical SOS certificates, which are converted into rational SOS using the `ZARITH OCAML` library (Version 1.2), implement arithmetic operations over arbitrary-precision integers. For more details about the installation and usage of `Real2Float`, we refer to the dedicated webpage<sup>3</sup> and the setup instructions.<sup>4</sup> All examples are displayed in Appendix A as the corresponding `Real2Float` input text files and satisfy our nonlinear program semantics (see Section 2.1). All results have been obtained on an Intel Core i7-5600U CPU (2.60GHz). Execution timings have been computed by averaging over five runs.

<sup>3</sup><http://nl-certify.forge.ocamlcore.org/real2float.html>.

<sup>4</sup>See the `README.md` file in the top level directory.

Table II. Comparison Results of Upper and Lower Bounds for Absolute Roundoff Errors among Tools Implementing Either Simple or Advanced Rounding Model. For Each Model, the Winner Results Are Emphasized Using Bold

Benchmark	id	Simple rounding			Improved rounding			lower bound
		Real2Float	Rosa	FPTaylor	FPTaylor	GAPPA	FLUCTUAT	
Programs involving polynomial functions								
rigidBody1	a	5.33e−13	5.08e−13	<b>3.87e−13</b>	<b>2.95e−13</b>	<b>2.95e−13</b>	3.22e−13	2.28e−13
rigidBody2	b	6.48e−11	6.48e−11	<b>5.24e−11</b>	<b>3.61e−11</b>	<b>3.61e−11</b>	3.65e−11	2.19e−11
kepler0	c	1.18e−13	1.16e−13	<b>1.05e−13</b>	<b>7.47e−14</b>	1.12e−13	1.26e−13	2.23e−14
kepler1	d	<b>4.47e−13</b>	6.49e−13	4.49e−13	<b>2.87e−13</b>	4.89e−13	5.57e−13	7.58e−14
kepler2	e	<b>2.09e−12</b>	2.89e−12	2.10e−12	<b>1.58e−12</b>	2.45e−12	2.90e−12	3.03e−13
sineTaylor	f	<b>6.03e−16</b>	9.56e−16	6.75e−16	<b>4.44e−16</b>	8.33e−02	6.86e−16	2.85e−16
sineOrder3	g	1.19e−15	1.11e−15	<b>9.97e−16</b>	7.95e−16	<b>7.62e−16</b>	1.03e−15	3.34e−16
sqroot	h	1.29e−15	8.41e−16	<b>7.13e−16</b>	<b>5.02e−16</b>	5.37e−16	3.21e−13	4.45e−16
himmilbeau	i	1.43e−12	1.43e−12	<b>1.32e−12</b>	<b>1.01e−12</b>	<b>1.01e−12</b>	<b>1.01e−12</b>	1.47e−13
Programs involving semialgebraic functions								
doppler1	j	7.65e−12	4.92e−13	<b>1.59e−13</b>	<b>1.29e−13</b>	1.82e−13	1.34e−13	7.11e−14
doppler2	k	1.57e−11	1.29e−12	<b>2.90e−13</b>	<b>2.39e−13</b>	3.23e−13	2.53e−13	1.14e−13
doppler3	l	8.55e−12	2.03e−13	<b>8.22e−14</b>	<b>6.96e−14</b>	9.29e−14	7.36e−14	4.27e−14
verhulst	m	4.67e−16	6.82e−16	<b>3.53e−16</b>	<b>2.50e−16</b>	3.18e−16	4.84e−16	2.23e−16
carbonGas	n	2.21e−08	4.64e−08	<b>1.23e−08</b>	<b>7.77e−09</b>	8.85e−09	1.86e−08	4.11e−09
predPrey	o	2.52e−16	2.94e−16	<b>1.89e−16</b>	<b>1.60e−16</b>	1.95e−16	2.45e−16	1.47e−16
turbine1	p	2.45e−11	1.25e−13	<b>2.33e−14</b>	<b>1.67e−14</b>	3.88e−14	6.09e−14	1.07e−14
turbine2	q	2.08e−12	1.76e−13	<b>3.14e−14</b>	<b>2.01e−14</b>	3.97e−14	8.96e−14	1.43e−14
turbine3	r	1.71e−11	8.50e−14	<b>1.70e−14</b>	<b>9.58e−15</b>	9.96e+00	4.90e−14	5.33e−15
jetEngine	s	OoM	1.62e−08	<b>1.50e−11</b>	<b>1.03e−11</b>	1.32e+05	1.82e−11	5.46e−12
Programs implementing polynomial functions with polynomial preconditions								
floudas2_6	t	<b>5.15e−13</b>	5.87e−13	7.88e−13	<b>5.94e−13</b>	5.98e−13	7.45e−13	4.56e−14
floudas3_3	u	5.81e−13	<b>4.05e−13</b>	5.76e−13	4.29e−13	<b>2.65e−13</b>	4.32e−13	1.48e−13
floudas3_4	v	2.78e−15	2.56e−15	<b>2.23e−15</b>	1.78e−15	<b>1.23e−15</b>	2.23e−15	3.80e−16
floudas4_6	w	1.82e−15	1.33e−15	<b>1.23e−15</b>	<b>8.89e−16</b>	<b>8.89e−16</b>	1.12e−15	2.35e−16
floudas4_7	x	<b>1.06e−14</b>	1.31e−14	1.80e−14	1.32e−14	<b>7.44e−15</b>	1.71e−14	7.31e−15
Programs involving conditional statements								
cav10	y	<b>2.91e+00</b>	<b>2.91e+00</b>	—	—	—	<b>1.02e+02</b>	2.90e+00
perin	z	<b>2.01e+00</b>	<b>2.01e+00</b>	—	—	—	<b>4.91e+01</b>	2.00e+00
Programs implementing transcendental functions								
logexp	$\alpha$	2.52e−15	—	<b>2.07e−15</b>	<b>1.99e−15</b>	—	—	1.19e−15
sphere	$\beta$	1.53e−14	—	<b>1.29e−14</b>	<b>8.21e−15</b>	—	—	5.05e−15
hartman3	$\gamma$	2.99e−13	—	<b>1.34e−14</b>	<b>4.97e−15</b>	—	—	1.10e−15
hartman6	$\delta$	5.09e−13	—	<b>2.55e−14</b>	<b>8.19e−15</b>	—	—	2.20e−15

#### 4.1. Benchmark Presentation

For each example, we compared the quality of the roundoff error bounds (Table II) and corresponding execution times (Table III) while running our tool Real2Float, FPTaylor (version from May 2016 [Solovyev et al. 2015]), Rosa (version from May 2014 used in Darulova and Kuncak [2014]), GAPPA (version 1.2.0 [Daumas and Melquiond 2010]), and FLUCTUAT (version 3.1370 [Delmas et al. 2009]).

To ensure fair comparison, our initial choice was to focus on tools providing certificates and using the same rounding model (FPTaylor or Rosa, which relies on an SMT solver theoretically able to output satisfiability certificates). However, for the sake of completeness, we have also compared Real2Float with GAPPA [Daumas and Melquiond 2010] and FLUCTUAT [Delmas et al. 2009]. For all tools, we use default parameters, we use the default number of subdivisions in FLUCTUAT and for GAPPA, and we provide the simplest user-provided hints we could think of.

Table III. Comparison of Execution Times (in Seconds) for Absolute Roundoff Error Bounds among Tools Implementing Either Simple or Advanced Rounding Model. For Each Model, the Winner Results Are Emphasized Using Bold

Benchmark	id	Simple rounding			Improved rounding		
		Real2Float	Rosa	FPTaylor	FPTaylor	GAPPA	FLUCTUAT
rigidBody1	a	0.58	<b>0.13</b>	1.84	0.41	<b>0.10</b>	0.29
rigidBody2	b	<b>0.26</b>	2.17	3.01	0.46	<b>0.15</b>	0.52
kepler0	c	<b>0.22</b>	3.78	4.93	6.96	<b>0.44</b>	0.70
kepler1	d	17.6	63.1	<b>9.33</b>	4.90	0.72	<b>0.37</b>
kepler2	e	<b>16.5</b>	106	19.1	13.8	<b>1.58</b>	2.04
sineTaylor	f	<b>1.05</b>	3.50	2.91	<b>0.11</b>	0.16	0.55
sineOrder3	g	<b>0.40</b>	0.48	1.90	0.40	<b>0.06</b>	0.22
sqroot	h	<b>0.14</b>	0.77	2.70	0.44	<b>0.19</b>	0.40
himmilbeau	i	<b>0.20</b>	2.51	3.28	0.49	<b>0.09</b>	1.00
doppler1	j	6.80	6.35	<b>6.13</b>	1.34	<b>0.08</b>	0.77
doppler2	k	6.96	<b>6.54</b>	6.88	1.57	<b>0.08</b>	0.79
doppler3	l	6.84	<b>6.37</b>	9.13	1.50	<b>0.07</b>	0.78
verhulst	m	<b>0.51</b>	1.36	1.37	0.40	<b>0.05</b>	0.25
carbonGas	n	<b>0.83</b>	6.59	3.73	0.52	<b>0.15</b>	2.03
predPrey	o	<b>0.87</b>	4.12	1.78	0.48	<b>0.04</b>	5.14
turbine1	p	72.2	<b>3.09</b>	4.38	0.53	<b>0.21</b>	5.79
turbine2	q	4.72	7.75	<b>3.25</b>	0.60	<b>0.12</b>	4.76
turbine3	r	74.5	4.57	<b>3.46</b>	0.61	<b>0.20</b>	5.84
jetEngine	s	OoM	125	<b>9.79</b>	3.06	<b>0.31</b>	31.2
floudas2_6	t	<b>2.49</b>	159	15.9	14.3	<b>2.35</b>	26.8
floudas3_3	u	<b>0.45</b>	13.9	5.64	13.9	<b>0.76</b>	6.41
floudas3_4	v	<b>0.09</b>	0.49	1.47	1.27	<b>0.07</b>	1.01
floudas4_6	w	<b>0.07</b>	1.20	0.91	0.37	<b>0.05</b>	1.69
floudas4_7	x	<b>0.13</b>	21.8	1.64	0.39	<b>0.07</b>	0.53
cav10	y	<b>0.23</b>	0.59	—	—	—	1.26
perin	z	<b>0.49</b>	2.74	—	—	—	1.19
logexp	$\alpha$	<b>1.05</b>	—	1.10	<b>0.39</b>	—	—
sphere	$\beta$	<b>0.05</b>	—	2.04	<b>3.69</b>	—	—
hartman3	$\gamma$	<b>2.02</b>	—	32.5	<b>27.8</b>	—	—
hartman6	$\delta$	<b>119</b>	—	364	<b>259</b>	—	—

A head-to-head comparison is not straightforward here due to differences in the approaches: GAPPA uses an improved rounding model based on a piecewise constant absolute error bound (see Section 1.2 for more details), and FLUCTUAT does not produce output certificates. We also performed further experiments while turning on the improved rounding model of FPTaylor (which is the same as in GAPPA and FLUCTUAT).

A given program implements a nonlinear function  $f(\mathbf{x})$ , involving variables  $\mathbf{x}$  lying in a set  $\mathbf{X}$  contained in a box  $[\mathbf{a}, \mathbf{b}]$ . Applying our rounding model on  $f$  yields the nonlinear expression  $\hat{f}(\mathbf{x}, \mathbf{e})$ , involving additional error variables  $\mathbf{e}$  lying in a set  $\mathbf{E}$ .

At a given semidefinite relaxation order  $d$ , our tool computes the upper bound  $f_d$  of the absolute roundoff error  $|f - \hat{f}|$  over  $\mathbf{K} := \mathbf{X} \times \mathbf{E}$  and verifies that it is less than a requested number  $\epsilon_{\text{Real2Float}}^+$ . As we keep the relaxation order  $d$  as low as possible to ensure tractable SDP programs, it can happen that  $f_d > \epsilon_{\text{Real2Float}}^+$ . The Real2Float tool has the default option to perform box subdivisions when the number of initial variables and maximal polynomial degree are both small. When the option is disabled, the solver does not perform subdivisions and output the error bound  $f_d$ . When enabled, we subdivide a randomly chosen interval of the box  $[\mathbf{a}, \mathbf{b}]$  in two halves to obtain two sub-sets  $\mathbf{X}_1$  and  $\mathbf{X}_2$ , fulfilling  $\mathbf{X} := \mathbf{X}_1 \cup \mathbf{X}_2$ , and apply the bound\_nprog algorithm on both sub-sets either until we succeed to certify that  $\epsilon_{\text{Real2Float}}^+$  is a sound upper bound of the roundoff error or until the maximal number of branch and bound iterations

is reached. For each benchmark, an error bound  $\epsilon_{\text{Real2Float}}$  is automatically computed while setting  $\epsilon_{\text{Real2Float}}^+ = 0$ .

The number  $\epsilon_{\text{Real2Float}}$  is compared with the upper bounds computed by two other tools implementing simple rounding models: FPTaylor, which relies on Taylor Symbolic expansions [Solovyev et al. 2015], and Rosa, which relies on SMT and affine arithmetic [Darulova and Kuncak 2014].

For comparison purposes, we also executed each program using random inputs, following the approach used in the Rosa article by Darulova and Kuncak [2014]. Specifically, we executed each program on  $10^7$  random inputs satisfying the input restrictions. The results from these random samples provide lower bounds on the absolute error. For consistency of comparison, the error bounds computed with FPTaylor correspond to the procedure FPT. (a) (see Solovyev et al. [2015]) using the same simplified rounding model as the one described in Equation (3), also used in Rosa [Darulova and Kuncak 2014]. For the sake of further presentation, we identify with a letter (from a to z and from  $\alpha$  to  $\gamma$ ) each of the 30 nonlinear programs. The programs a-b and f-s are taken from the Rosa article [Darulova and Kuncak 2014] and were used in the FPTaylor article [Solovyev et al. 2015] as well:

- The first nine programs implement polynomial functions: a-b come from physics, c-e are derived from expressions involved in the proof of Kepler Conjecture [Hales 2006], and f-h implement polynomial approximations of the sine and square-root functions. The program i is issued from the global optimization literature and implements the problem *Himmelbeau* in Ali et al. [2005].
- The 10 programs implement semialgebraic functions: j-l and p-s come from physics, m and o from biology, and n from control. All these programs are used to compare FPTaylor and Rosa in Solovyev et al. [2015].
- The five programs t-x come from the global optimization literature and correspond respectively to Problems 2.6, 3.3, 3.4, 4.6, and 4.7 in Floudas and Pardalos [1990]. We selected them as they typically involve nontrivial polynomial preconditions (i.e.,  $\mathbf{X}$  is not a simple box but rather a set defined with conjunction of nonlinear polynomial inequalities).
- The two programs y-z involve conditional statements and come from the static analysis literature. They correspond to the two respective running examples of Ghorbal et al. [2010] (FLUCTUAT's divergence error computation) and Alexandre Maréchal [2014]. The first program y is used in the Rosa article [Darulova and Kuncak 2014] for the analysis of branches discontinuity error.
- The last four programs  $\alpha$ - $\gamma$  involve transcendental functions. The two programs  $\alpha$  and  $\beta$  are used in the FPTaylor article [Solovyev et al. 2015] and correspond respectively to the program `logexp` (see Example 3.3) and the program `sphere` taken from NASA World Wind Java SDK [Nasa 2011]. The two programs  $\gamma$  and  $\delta$ , respectively, implement the functions coming from the optimization problems *Hartman 3* and *Hartman 6* in Ali et al. [2005], involving both sums of exponential functions composed with quadratic polynomials.

*Tool comparison settings.* The five tools Real2Float, Rosa, FPTaylor, GAPPA, and FLUCTUAT can handle programs with input variable uncertainties as well as any floating-point precision, but, for the sake of conciseness, we only compare their performance on programs implemented in double ( $\epsilon = 2^{-53}$ ) precision floating-point. By contrast with preliminary experiments presented in Section 1.1, where we considered floating-point input variables, we run each tool by considering all input variables as real variables, and thus we apply the rounding operation to all of them. For the programs involving transcendental functions, we followed the same procedure as in FPTaylor while adjusting the precision  $\epsilon(f_{\mathbb{R}}) = 1.5\epsilon$  for each special function  $f_{\mathbb{R}} \in \{\sin, \cos, \log, \exp\}$ . Each univariate transcendental function is approximated from below (respectively, from

above) using suprema (respectively, infima) of linear or quadratic polynomials (see Example 3.3 for the case of program `logexp`).

## 4.2. Comparison Results

Comparison results for error bound computation are presented in Table II. Among the tools using simple rounding model (`Real2Float`, `Rosa`, and `FPTaylor`), our `Real2Float` tool computes the tightest upper bounds for 7 (respectively, 4) of 30 benchmarks when comparing with `Rosa` and `FPTaylor` (respectively, all tools). For all programs `j-s` involving semialgebraic functions and the four programs  $\alpha\text{-}\delta$  involving transcendental functions, `FPTaylor` computes the tightest bounds, when comparing with `Real2Float` and `Rosa`.<sup>5</sup> One current limitation of `Real2Float` is its ability to manipulate symbolic expressions, for example, computing rational function derivatives or yielding reduction to the same denominator. In particular, the analysis of program `s` aborted after running out of memory (meaning of the symbol `OoM`). The `FPTaylor` tool is better suited to handle programs that exhibit such rational functions and also includes an interface with the `MAXIMA` computer algebra system [Maxima 2013] to perform symbolic simplifications.

We mention that the interested reader can find more detailed experimental comparisons between the three tools implementing improved rounding models (`FPTaylor`, `FLUCTUAT`, and `GAPPA`) in Solovyev et al. [2015, Section 5.2]. Note that `FPTaylor` provides the tightest bounds for 24 of 30 benchmarks. The `GAPPA` software provides the tightest bounds for 8 of 30 benchmarks while being almost always faster than other tools. As emphasized in Solovyev et al. [2015], `GAPPA` can sometimes compute more precise bounds with more advanced user-provided hints. Note that `FLUCTUAT` provides sometimes tighter bounds than `GAPPA` for most rational functions while being always slower. It would be worth implementing the same improved rounding model in `Real2Float` to perform numerical comparisons.

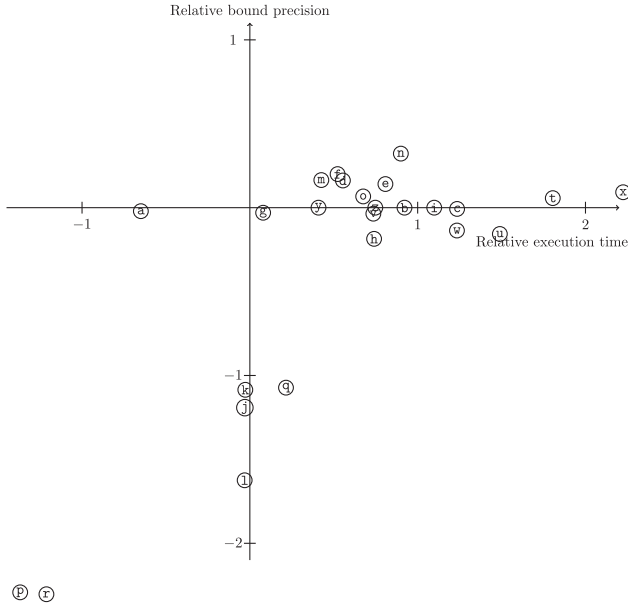
To the best of our knowledge, `Real2Float` is the only academic tool that is able to handle the general class of programs involving either transcendental functions or conditional statements. The `FPTaylor` (respectively, `Rosa`) tool does not currently handle conditionals (respectively, transcendental functions), as meant by the symbol `-` in the corresponding column entries. However, an interface bridging the `FPTaylor` and `Rosa` tools would provide each other with the relevant missing features. These error bound comparison results together with their corresponding execution timings (given in Table III) are used to plot the data points shown in Figure 7. For each benchmark identified by `id`, let  $t_{\text{Real2Float}}$  (in third column of Table III) refer to the execution time of `Real2Float` to obtain the corresponding upper bound  $\epsilon_{\text{Real2Float}}$  (in third column of Table II).

Now, let us define the execution times  $t_{\text{Rosa}}$ ,  $t_{\text{FPTaylor}}$  and the corresponding error bounds  $\epsilon_{\text{Rosa}}$ ,  $\epsilon_{\text{FPTaylor}}$ . Then the  $x$ -axis coordinate of the point  $\textcircled{d}$  (respectively,  $\textcircled{id}$ ) displayed in Figure 7(a) (respectively, 7(b)) corresponds to the logarithm of the ratio between the execution time of `Rosa` (respectively, `FPTaylor`) and `Real2Float`, that is,  $\log \frac{t_{\text{Rosa}}}{t_{\text{Real2Float}}}$  (respectively,  $\log \frac{t_{\text{FPTaylor}}}{t_{\text{Real2Float}}}$ ). Similarly, the  $y$ -axis coordinate of the point  $\textcircled{d}$  (respectively,  $\textcircled{id}$ ) is  $\log \frac{\epsilon_{\text{Rosa}}}{\epsilon_{\text{Real2Float}}}$  (respectively,  $\log \frac{\epsilon_{\text{FPTaylor}}}{\epsilon_{\text{Real2Float}}}$ ).

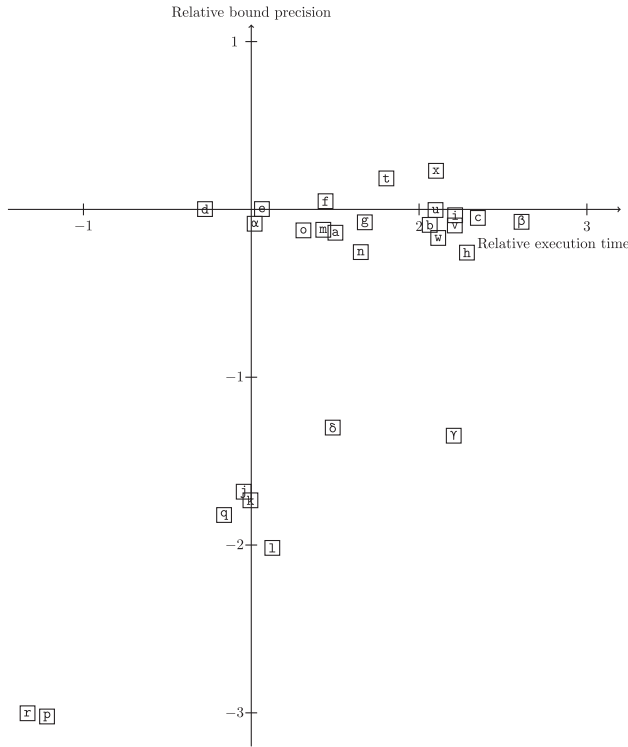
The axes of the coordinate system within Figure 7 divide the plane into four quadrants: the nonnegative quadrant  $(+, +)$  contains data points referring to programs for which `Real2Float` computes the tighter bounds in less time, the second one  $(+, -)$  contains points referring to programs for which `Real2Float` is faster but less accurate, the non-positive quadrant  $(-, -)$  for which `Real2Float` is slower and computes coarser bounds, and the last one  $(-, +)$  for which `Real2Float` is slower but more accurate.

<sup>5</sup>The running execution times of `Rosa` may change with more recent versions.





(a) Comparison with Rosa.



(b) Comparison with FPTaylor.

Fig. 7. Comparisons of execution times and upper bounds of roundoff errors obtained with Rosa and FPTaylor, relatively to Real2Float.

Table IV. Comparisons of Informal and Formal Execution Times to Certify Roundoff Error Bounds Obtained with Real2Float, FPTaylor, and GAPPA

Benchmark	id	Informal timings		Formal timings		Speedup	
		$t_{\text{Real2Float}}$	$t_{\text{FPTaylor}}$	$t_{\text{Coq}}$	$t_{\text{HOL-LIGHT}}$	Informal	Formal
rigidBody1	a	0.58	1.84	0.36	10.2	3.18	12.9
rigidBody2	b	0.26	3.01	4.81	32.3	11.6	6.97
kepler0	c	0.22	4.93	0.29	45.5	22.5	98.9
kepler1	d	17.6	9.33	449	90.5	0.54	0.22
kepler2	e	16.5	19.1	297	274	1.16	0.94
sineTaylor	f	1.05	2.91	7.54	42.1	2.78	5.24
sineOrder3	g	0.40	1.90	0.34	10.4	4.75	16.7
sroot	h	0.14	2.70	0.80	16.8	19.3	20.8
himmilbeau	i	0.20	3.28	0.89	32.2	16.4	32.6

On the quadrant  $(+, -)$  of Figure 7(b), one can see that Real2Float computes bounds that are less accurate than FPTaylor on semialgebraic and transcendental programs but faster for most of them. The quadrant  $(-, -)$  indicates that Rosa and FPTaylor are more precise and efficient than Real2Float on the three programs p-r. The presence of 18 plots on the nonnegative quadrant  $(+, +)$  of Figures 7(a) and 7(b) confirms that Real2Float does not compromise efficiency at the expense of accuracy, in particular for programs implementing polynomials with nontrivial polynomial preconditions.

For each program implementing polynomials, our tool has an option to provide formal guarantees for the corresponding roundoff error bound  $\epsilon_{\text{Real2Float}}$ . Using the formal mechanism described in Section 2.3, Real2Float formally checks inside Coq the SOS certificates generated by the SDP solver for interval enclosures of linear error terms  $l$ . The FPTaylor (respectively, GAPPA) software has a similar option to provide formal scripts that can be checked inside the HOL-LIGHT (respectively, Coq) proof assistant for programs involving polynomial and transcendental functions. For formal verification, our tool is limited compared with FPTaylor and GAPPA, as it cannot handle non-polynomial programs.

Next, we describe the formal proof results obtained while verifying the bounds for the nine polynomial programs a-i. Table IV allows us to compare the execution times of Real2Float, FPTaylor, and GAPPA required to analyze the nine programs in both informal (i.e., without verification inside Coq or HOL-LIGHT) and formal settings. Comparing Real2Float with FPTaylor, the latter performs better than the former to analyze formally the two programs d-e but yields coarser bounds. Our formal procedure is less computationally efficient when the degree (respectively, number of variables) of the SOS polynomials grows, as for these two programs. The informal speedup ratio is greater than 3 for the five programs a, c, and g-i. The table shows that the speedup ratio in the formal setting is higher than in the informal setting for these five programs. The explanation could be that Coq is inherently faster than HOL-LIGHT at checking computations while delegating expensive ones to a virtual machine (being part of Coq's trusted base). Either SOS or Taylor methods could work with the two proof assistants HOL-LIGHT and Coq (with natural changes in execution time), and it would be meaningful to compare similar methods in a given proof assistant (SOS techniques in HOL-LIGHT or Taylor methods in Coq), but both are not yet fully available. The GAPPA tool performs better to analyze the two programs d-e but is less accurate. This is in contrast with all other programs where the formal verification with Real2Float is faster than with GAPPA while providing coarser bounds (except program f where GAPPA yields pessimistic results). Our formal verification framework is a work in progress as it only allows to check correctness of SOS certificates for polynomial programs. The step from formal verification of the step translating the inequality to a statement of explicit syntactic form "Forall  $x$ ,  $|\text{rounded}(f(x)) - f(x)| \leq e$ " is missing and requires more software engineering. Thus, the timings presented in Table IV for Real2Float do not include formal

computation of the symbolic second-order derivatives of  $r$  w.r.t.  $\mathbf{e}$  as well as the cost of bounding them using formal interval arithmetic.

At the moment (and in contrast to our method), the Rosa tool does not formally verify the output bound provided by the SMT solver, but such a feature could be embedded through an interface with the `smtcoq` framework [Armand et al. 2011]. This latter tool allows the proof witness generated by an SMT solver to be formally (and independently) re-checked inside Coq. The `smtcoq` framework uses tactics based on *computational reflection* to enable this re-checking to be performed efficiently.

## 5. CONCLUSION AND PERSPECTIVES

Our verification framework allows us to over-approximate roundoff errors occurring while executing nonlinear programs implemented with finite precision. The framework relies on semidefinite optimization, ensuring certified approximations. Our approach extends to medium-size nonlinear problems, due to automatic detection of the correlation sparsity pattern of input variables and roundoff error variables. Experimental results indicate that the optimization algorithm implemented in our `Real2Float` software package can often produce bounds of quality similar to the ones provided by the competitive solvers Rosa and `FPTaylor`, while saving a significant amount of CPU time. In addition, `Real2Float` produces sums of squares certificates that guarantee the correctness of these upper bounds and can be efficiently verified inside the Coq proof assistant.

This work yields several directions for further research investigation. First, we intend to increase the size of graspable instances by exploiting the SDP relaxations specifically tuned to the case when a program implements the sum of many rational functions [Bugarin et al. 2015]. Symmetry patterns of certain program sub-classes could be tackled with the SDP hierarchies from Riemer et al. [2013]. We could also provide roundoff error bounds for more general programs, involving either finite or infinite conditional loops and additional comparisons with the results described in Darulova and Kuncak [2016]. A preliminary mandatory step is to be able to generate inductive invariants with SDP relaxations. Another interesting direction would be to apply the method used in Lasserre [2011] to derive sequences of lower roundoff error bounds together with SDP-based certificates. On the formal proof side, we could benefit from floating-point/interval arithmetic libraries available inside Coq, first to improve the efficiency of the formal polynomial checker, currently relying on exact arithmetic, and then to extend the formal verification to non-polynomial programs. The method implemented in `FPTaylor` happens to be efficient and precise to analyze various programs, and it would be interesting to design a procedure combining `FPTaylor` with our tool on specific subsets of input constraints. Long-term research perspectives include theoretical study of why/when SOS performs better as well as satisfactory complexity results about SOS certificates (w.r.t. size of polynomials). Finally, we plan to combine this optimization framework with the procedure in Gao and Constantinides [2015] to improve the automatic reordering of arithmetic expressions, allowing more efficient optimization of FPGA implementations.

## NONLINEAR PROGRAM BENCHMARKS

```
let box_rigidbody1 x1 x2 x3 = [(-15, 15); (-15, 15); (-15, 15)];;
let obj_rigidbody1 x1 x2 x3 = [(-x1 * x2 - 2 * x2 * x3 - x1 - x3, 0)];;

let box_rigidbody2 x1 x2 x3 = [(-15, 15); (-15, 15); (-15, 15)];;
let obj_rigidbody2 x1 x2 x3 =
[(2 * x1 * x2 * x3 + 3 * x3 * x3 - x2 * x1 * x2 * x3 + 3 * x3 * x3 - x2, 0)];;

let box_kepler0 x1 x2 x3 x4 x5 x6 = [(4, 6.36); (4, 6.36); (4, 6.36); (4, 6.36); (4, 6.36); (4, 6.36)];;
let obj_kepler0 x1 x2 x3 x4 x5 x6 =
[(x2 * x5 + x3 * x6 - x2 * x3 - x5 * x6 + x1 * (-x1 + x2 + x3 - x4 + x5 + x6), 0)];;
```

```

let box_kepler1  $x_1 x_2 x_3 x_4 = [(4, 6.36); (4, 6.36); (4, 6.36); (4, 6.36)];$ ;
let obj_kepler1  $x_1 x_2 x_3 x_4 = [(x_1 * x_4 * (-x_1 + x_2 + x_3 - x_4)$ 
 $+ x_2 * (x_1 - x_2 + x_3 + x_4) + x_3 * (x_1 + x_2 - x_3 + x_4)$ 
 $- x_2 * x_3 * x_4 - x_1 * x_3 - x_1 * x_2 - x_4, 0)];$ ;

let box_kepler2  $x_1 x_2 x_3 x_4 x_5 x_6 = [(4, 6.36); (4, 6.36); (4, 6.36); (4, 6.36); (4, 6.36); (4, 6.36)];$ ;
let obj_kepler2  $x_1 x_2 x_3 x_4 x_5 x_6 = [(x_1 * x_4 * (-x_1 + x_2 + x_3$ 
 $- x_4 + x_5 + x_6) + x_2 * x_5 * (x_1 - x_2 + x_3 + x_4 - x_5 + x_6)$ 
 $+ x_3 * x_6 * (x_1 + x_2 - x_3 + x_4 + x_5 - x_6) - x_2 * x_3 * x_4$ 
 $- x_1 * x_3 * x_5 - x_1 * x_2 * x_6 - x_4 * x_5 * x_6, 0)];$ ;

let box_sineTaylor  $x = [(-1.57079632679, 1.57079632679)];$ ;
let obj_sineTaylor  $x = [(x - (x * x * x)/6.0$ 
 $+(x * x * x * x * x)/120.0$ 
 $-(x * x * x * x * x * x * x)/5040.0, 0)];$ ;

let box_sineOrder3  $z = [(-2, 2)];$ ;
let obj_sineOrder3  $z = [(0.954929658551372 * z - 0.12900613773279798 * (z * z * z), 0)];$ ;

let box_sqrtroot  $y = [(0, 1)];$ ;
let obj_sqrtroot  $y = [(1.0 + 0.5 * y - 0.125 * y * y + 0.0625 * y * y * y - 0.0390625 * y * y * y * y, 0)];$ ;

let box_himmilbeau  $x_1 x_2 = [(-5, 5); (-5, 5)];$ ;
let obj_himmilbeau  $x_1 x_2 = [($ 
 $(x_1 * x_1 + x_2 - 11) * (x_1 * x_1 + x_2 - 11) + (x_1 + x_2 * x_2 - 7) * (x_1 + x_2 * x_2 - 7), 0)];$ ;

let box_doppler1  $u v T = [(-100, 100); (20, 20e3); (-30, 50)];$ ;
let obj_doppler1  $u v T = [(let t_1 = 331.4 + 0.6 * T in -t_1 * v / ((t_1 + u) * (t_1 + u)), 0)];$ ;

let box_doppler2  $u v T = [(-125, 125); (15, 25e3); (-40, 60)];$ ;
let obj_doppler2  $u v T = [(let t_1 = 331.4 + 0.6 * T in -t_1 * v / ((t_1 + u) * (t_1 + u)), 0)];$ ;

let box_doppler3  $u v T = [(-30, 120); (320, 20300); (-50, 30)];$ ;
let obj_doppler3  $u v T = [(let t_1 = 331.4 + 0.6 * T in -t_1 * v / ((t_1 + u) * (t_1 + u)), 0)];$ ;

let box_verhulst  $x = [(0.1, 0.3)];$ ;
let obj_verhulst  $x = [(4 * x / (1 + (x / 1.11)), 0)];$ ;

let box_carbonGas  $v = [(0.1, 0.5)];$ ;
let obj_carbonGas  $v = [(let p = 3.5e7 in let a = 0.401 in$ 
 $let b = 42.7e-6 in let t = 300 in let n = 1000 in$ 
 $(p + a * (n / v) * * 2) * (v - n * b) - 1.3806503e-23 * n * t, 0)];$ ;

let box_predPrey  $x = [(0.1, 0.3)];$ ;
let obj_predPrey  $x = [(4 * x * x / (1 + (x / 1.11) * * 2), 0)];$ ;

let box_turbine1  $v w r = [(-4.5, -0.3); (0.4, 0.9); (3.8, 7.8)];$ ;
let obj_turbine1  $v w r = [(3 + 2 / (r * r) - 0.125 * (3 - 2 * v) * (w * w * r * r) / (1 - v) - 4.5, 0)];$ ;

let box_turbine2  $v w r = [(-4.5, -0.3); (0.4, 0.9); (3.8, 7.8)];$ ;
let obj_turbine2  $v w r = [(6 * v - 0.5 * v * (w * w * r * r) / (1 - v) - 2.5, 0)];$ ;

let box_turbine3  $v w r = [(-4.5, -0.3); (0.4, 0.9); (3.8, 7.8)];$ ;
let obj_turbine3  $v w r = [(3 - 2 / (r * r) - 0.125 * (1 + 2 * v) * (w * w * r * r) / (1 - v) - 0.5, 0)];$ ;

let box_jet  $x_1 x_2 = [(-5, 5); (-20, 5)];$ ;
let obj_jet  $x_1 x_2 = [($ 
 $x_1 + ((2 * x_1 * ((3 * x_1 * x_1 + 2 * x_2 - x_1) / (x_1 * x_1 + 1))$ 
 $* ((3 * x_1 * x_1 + 2 * x_2 - x_1) / (x_1 * x_1 + 1) - 3)$ 
 $+ x_1 * x_1 * (4 * ((3 * x_1 * x_1 + 2 * x_2 - x_1) / (x_1 * x_1 + 1)) - 6))$ 
 $* (x_1 * x_1 + 1) + 3 * x_1 * x_1 * ((3 * x_1 * x_1 + 2 * x_2 - x_1) / (x_1 * x_1 + 1))$ 
 $+ x_1 * x_1 * x_1 + x_1 + 3 * ((3 * x_1 * x_1 + 2 * x_2 - x_1) / (x_1 * x_1 + 1))), 0)];$ ;

let box_floudas2_6  $x_1 x_2 x_3 x_4 x_5 x_6 x_7 x_8 x_9 x_{10} =$ 
 $[(0, 1); (0, 1); (0, 1); (0, 1); (0, 1); (0, 1); (0, 1); (0, 1)];$ ;
let cstr_floudas2_6  $x_1 x_2 x_3 x_4 x_5 x_6 x_7 x_8 x_9 x_{10} = [$ 
 $-4 + 2 * x_1 + 6 * x_2 + 1 * x_3 + 0 * x_4 + 3 * x_5 + 3 * x_6 + 2 * x_7 + 6 * x_8 + 2 * x_9 + 2 * x_{10};$ 
 $22 - (6 * x_1 - 5 * x_2 + 8 * x_3 - 3 * x_4 + 0 * x_5 + 1 * x_6 + 3 * x_7 + 8 * x_8 + 9 * x_9 - 3 * x_{10});$ 
 $-6 - (5 * x_1 + 6 * x_2 + 5 * x_3 + 3 * x_4 + 8 * x_5 - 8 * x_6 + 9 * x_7 + 2 * x_8 + 0 * x_9 - 9 * x_{10});$ 

```

```

-23 - (9 * x1 + 5 * x2 + 0 * x3 - 9 * x4 + 1 * x5 - 8 * x6 + 3 * x7 - 9 * x8 - 9 * x9 - 3 * x10);
-12 - (-8 * x1 + 7 * x2 - 4 * x3 - 5 * x4 - 9 * x5 + 1 * x6 - 7 * x7 - 1 * x8 + 3 * x9 - 2 * x10)];;
let obj_floudas2_6 x1 x2 x3 x4 x5 x6 x7 x8 x9 x10 = [(
48 * x1 + 42 * x2 + 48 * x3 + 45 * x4 + 44 * x5 + 41 * x6 + 47 * x7
+ 42 * x8 + 45 * x9 + 46 * x10
- 50 * (x1 * x1 + x2 * x2 + x3 * x3 + x4 * x4 + x5 * x5
+ x6 * x6 + x7 * x7 + x8 * x8 + x9 * x9 + x10 * x10), 0)];;

let box_floudas3_3 x1 x2 x3 x4 x5 x6 = [(0, 6); (0, 6); (1, 5); (0, 6); (1, 5); (0, 10)];;
-(x2 - 2) ** 2 - (x3 - 1) ** 2 - (x4 - 4) ** 2 - (x5 - 1) ** 2 - (x6 - 4) ** 2, 0)];;
let cstr_floudas3_3 x1 x2 x3 x4 x5 x6 =
[(x3 - 3) ** 2 + x4 - 4; (x5 - 3) ** 2 + x6 - 4;
2 - x1 + 3 * x2; 2 + x1 - x2; 6 - x1 - x2; x1 + x2 - 2];;
let obj_floudas3_3 x1 x2 x3 x4 x5 x6 = [(-25 * (x1 - 2) ** 2
-(x2 - 2) ** 2 - (x3 - 1) ** 2 - (x4 - 4) ** 2
-(x5 - 1) ** 2 - (x6 - 4) ** 2, 0)];;

let box_floudas3_4 x1 x2 x3 = [(0, 2); (0, 2); (0, 3)];;
let cstr_floudas3_4 x1 x2 x3 = [
4 - x1 - x2 - x3; 6 - 3 * x2 - x3;
-0.75 + 2 * x1 - 2 * x3 + 4 * x1 * x1 - 4 * x1 * x2
+ 4 * x1 * x3 + 2 * x2 * x2 - 2 * x2 * x3 + 2 * x3 * x3];;
let obj_floudas3_4 x1 x2 x3 = [(-2 * x1 + x2 - x3, 0)];;

let box_floudas4_6 x1 x2 = [(0, 3); (0, 4)];;
let cstr_floudas4_6 x1 x2 = [
2 * x1 ** 4 - 8 * x1 ** 3 + 8 * x1 * x1 - x2;
4 * x1 ** 4 - 32 * x1 ** 3 + 88 * x1 * x1 - 96 * x1 + 36 - x2];;
let obj_floudas4_6 x1 x2 = [(-x1 - x2, 0)];;

let box_floudas4_7 x1 x2 = [(0, 2); (0, 3)];;
let cstr_floudas4_7 x1 x2 = [-2 * x1 ** 4 + 2 - x2];;
let obj_floudas4_7 x1 x2 = [(-12 * x1 - 7 * x2 + x2 * x2, 0)];;

let box_cav10 x = [(0, 10)];;
let obj_cav10 x = [(if (x * x - x > 0) then x * 0.1 else x * x + 2, 0)];;

let box_perin xy = [(1, 7); (-2, 7)];;
let cstr_perin xy = [x - 1; y + 2; x - y; 5 - y - x];;
let obj_perin xy = [(if (x * x + y * y ≤ 4) then y * x else 0, 0)];;

let box_logexp x = [(-8, 8)];;
let obj_logexp x = [(log(1 + exp(x)), 0)];;

let box_sphere xyz = [(-10, 10); (0, 10); (-1.570796, 1.570796); (-3.14159265, 3.14159265)];;
let obj_sphere xyz = [(x + r * sin(y) * cos(z), 0)];;

let box_hartman3 x1 x2 x3 = [(0, 1); (0, 1); (0, 1)];;
let obj_hartman3 x1 x2 x3 = [(
let e1 = 3.0 * (x1 - 0.3689) ** 2 + 10.0 * (x2 - 0.117) ** 2 + 30.0 * (x3 - 0.2673) ** 2 in
let e2 = 0.1 * (x1 - 0.4699) ** 2 + 10.0 * (x2 - 0.4387) ** 2 + 35.0 * (x3 - 0.747) ** 2 in
let e3 = 3.0 * (x1 - 0.1091) ** 2 + 10.0 * (x2 - 0.8732) ** 2 + 30.0 * (x3 - 0.5547) ** 2 in
let e4 = 0.1 * (x1 - 0.03815) ** 2 + 10.0 * (x2 - 0.5743) ** 2 + 35.0 * (x3 - 0.8828) ** 2 in
- (1.0 * exp(-e1) + 1.2 * exp(-e2) + 3.0 * exp(-e3) + 3.2 * exp(-e4)), 0)];;

let box_hartman6 x1 x2 x3 x4 x5 x6 = [(0, 1); (0, 1); (0, 1); (0, 1); (0, 1); (0, 1)];;
let obj_hartman6 x1 x2 x3 x4 x5 x6 = [(
let e1 = 10.0 * (x1 - 0.1312) ** 2 + 3.0 * (x2 - 0.1696) ** 2 + 17.0 * (x3 - 0.5569) ** 2 + 3.5 * (x4 - 0.0124) ** 2
+ 1.7 * (x5 - 0.8283) ** 2 + 8.0 * (x6 - 0.5886) ** 2 in
let e2 = 0.05 * (x1 - 0.2329) ** 2 + 10.0 * (x2 - 0.4135) ** 2 + 17.0 * (x3 - 0.8307) ** 2 + 0.1 * (x4 - 0.3736) ** 2
+ 8.0 * (x5 - 0.1004) ** 2 + 14.0 * (x6 - 0.9991) ** 2 in
let e3 = 3.0 * (x1 - 0.2348) ** 2 + 3.5 * (x2 - 0.1451) ** 2 + 1.7 * (x3 - 0.3522) ** 2 + 10.0 * (x4 - 0.2883) ** 2
+ 17.0 * (x5 - 0.3047) ** 2 + 8.0 * (x6 - 0.665) ** 2 in
let e4 = 17.0 * (x1 - 0.4047) ** 2 + 8 * (x2 - 0.8828) ** 2 + 0.05 * (x3 - 0.8732) ** 2 + 10.0 * (x4 - 0.5743) ** 2
+ 0.1 * (x5 - 0.1091) ** 2 + 14.0 * (x6 - 0.0381) ** 2 in
- (1.0 * exp(-e1) + 1.2 * exp(-e2) + 3.0 * exp(-e3) + 3.2 * exp(-e4)), 0)];;

```

## ACKNOWLEDGMENTS

The authors acknowledge the help of Alexey Solovyev for his excellent feedback and suggestions. They also thank the three referees for helpful comments that improved this article.



## REFERENCES

- Michaël Périn Alexandre Maréchal. 2014. *Three Linearization Techniques for Multivariate Polynomials in Static Analysis using Convex Polyhedra*. Technical Report TR-2014-7. Verimag Research Report.
- M. Montaz Ali, Charoenchai Khompapatporn, and Zeld B. Zabinsky. 2005. A numerical evaluation of several stochastic algorithms on selected continuous global optimization test problems. *J. Global Optim.* 31, 4 (April 2005), 635–672. DOI: <http://dx.doi.org/10.1007/s10898-004-9972-2>
- ErlingD. Andersen and KnudD. Andersen. 2000. The mosek interior point optimizer for linear programming: An implementation of the homogeneous algorithm. In *High Performance Optimization*, Hans Frenk, Kees Roos, Tamás Terlaky, and Shuzhong Zhang (Eds.). Applied Optimization, Vol. 33. Springer, Berlin, 197–232. DOI: [http://dx.doi.org/10.1007/978-1-4757-3216-0\\_8](http://dx.doi.org/10.1007/978-1-4757-3216-0_8)
- Michael Armand, Germain Faure, Benjamin Grégoire, Chantal Keller, Laurent Théry, and Benjamin Werner. 2011. A modular integration of SAT/SMT solvers to coq through proof witnesses. In *Certified Programs and Proofs*, Jean-Pierre Jouannaud and Zhong Shao (Eds.). Lecture Notes in Computer Science, Vol. 7086. Springer, Berlin, 135–150. DOI: [http://dx.doi.org/10.1007/978-3-642-25379-9\\_12](http://dx.doi.org/10.1007/978-3-642-25379-9_12)
- Aharon Ben-Tal and Arkadi Semenovich Nemirovski. 2001. *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*. Society for Industrial and Applied Mathematics, Philadelphia, PA.
- Y. Bertot and P. Castéran. 2004. *Interactive Theorem Proving and Program Development: Coq'Art: The Calculus of Inductive Constructions*. Springer.
- Martin Berz and Kyoko Makino. 2009. Rigorous global search using Taylor models. In *Proceedings of the 2009 Conference on Symbolic Numeric Computation (SNC'09)*. ACM, New York, NY, 11–20.
- Jesse Bingham and Joe Leslie-Hurd. 2014. Verifying relative error bounds using symbolic simulation. In *Computer Aided Verification*, Armin Biere and Roderick Bloem (Eds.). Lecture Notes in Computer Science, Vol. 8559. Springer, 277–292. DOI: [http://dx.doi.org/10.1007/978-3-319-08867-9\\_18](http://dx.doi.org/10.1007/978-3-319-08867-9_18)
- B. Blanchet, P. Cousot, R. Cousot, J. Feret, L. Mauborgne, A. Miné, D. Monniaux, and X. Rival. 2003. A static analyzer for large safety-critical software. In *Proceedings of the ACM SIGPLAN 2003 Conference on Programming Language Design and Implementation (PLDI'03)*. ACM Press, San Diego, CA, 196–207.
- David Boland and George A. Constantinides. 2010. Automated precision analysis: A polynomial algebraic approach. In *FCCM'10*. 157–164.
- S. Boldo and G. Melquiond. 2011. Floccq: A unified library for proving floating-point algorithms in coq. In *2011 20th IEEE Symposium on Computer Arithmetic (ARITH)*. 243–252. DOI: <http://dx.doi.org/10.1109/ARITH.2011.40>
- Mateus Borges, Marcelo d'Amorim, Saswat Anand, David Bushnell, and Corina S. Pasareanu. 2012. Symbolic execution with interval solving and meta-heuristic search. In *Proceedings of the 2012 IEEE 5th International Conference on Software Testing, Verification and Validation (ICST'12)*. IEEE Computer Society, Washington, DC, 111–120.
- Samuel Boutin. 1997. Using reflection to build efficient and certified decision procedures. In *TACS'97. Springer-Verlag LNCS 1281*. Springer-Verlag, 515–529.
- Florian Bugarin, Didier Henrion, and JeanBernard Lasserre. 2015. Minimizing the sum of many rational functions. *Math. Program. Comput.* (2015), 1–29. <http://dx.doi.org/10.1007/s12532-015-0089-z>
- Liqian Chen, Antoine Miné, and Patrick Cousot. 2008. A sound floating-point polyhedra abstract domain. In *Programming Languages and Systems*, G. Ramalingam (Ed.). Lecture Notes in Computer Science, Vol. 5356. Springer, Berlin, 3–18. DOI: [http://dx.doi.org/10.1007/978-3-540-89330-1\\_2](http://dx.doi.org/10.1007/978-3-540-89330-1_2)
- Liqian Chen, Antoine Miné, Ji Wang, and Patrick Cousot. 2009. Interval polyhedra: An abstract domain to infer interval linear relationships. In *Static Analysis*, Jens Palsberg and Zhendong Su (Eds.). Lecture Notes in Computer Science, Vol. 5673. Springer, Berlin, 309–325. DOI: [http://dx.doi.org/10.1007/978-3-642-03237-0\\_21](http://dx.doi.org/10.1007/978-3-642-03237-0_21)
- Wei-Fan Chiang, Ganesh Gopalakrishnan, Zvonimir Rakamaric, and Alexey Solov'yev. 2014. Efficient search for inputs causing high floating-point errors. In *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP'14)*. ACM, New York, NY, 43–52.
- Coq 2016. The Coq Proof Assistant. (2016). Retrieved from <http://coq.inria.fr/>.
- P. Cousot and R. Cousot. 1977. Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Conference Record of the 4th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM Press, New York, NY, 238–252.
- Eva Darulova and Viktor Kuncak. 2014. Sound compilation of reals. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL'14)*. ACM, New York, NY, 235–248.
- Eva Darulova and Viktor Kuncak. 2016. *Towards a Compiler for Reals*. Technical Report. Ecole Polytechnique Federale de Lausanne.

- Marc Daumas and Guillaume Melquiond. 2010. Certification of bounds on expressions involving rounded operators. *ACM Trans. Math. Softw.* 37, 1, Article 2 (Jan. 2010), 20 pages.
- Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An efficient SMT solver. In *Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS'08 / ETAPS'08)*. Springer-Verlag, Berlin, 337–340.
- David Delmas, Eric Goubault, Sylvie Putot, Jean Souyris, Karim Tekkal, and Franck Védrine. 2009. Towards an industrial use of FLUCTUAT on safety-critical avionics software. In *Formal Methods for Industrial Critical Systems*, Maria Alpuente, Byron Cook, and Christophe Joubert (Eds.). Lecture Notes in Computer Science, Vol. 5825. Springer, Berlin, 53–69. DOI: [http://dx.doi.org/10.1007/978-3-642-04570-7\\_6](http://dx.doi.org/10.1007/978-3-642-04570-7_6)
- Christodoulos A. Floudas and Panos M. Pardalos. 1990. *A Collection of Test Problems for Constrained Global Optimization Algorithms*. Springer-Verlag, New York, NY.
- Sicun Gao, Soonho Kong, and Edmund M. Clarke. 2013. dReal: An SMT solver for nonlinear theories over the reals. In *Automated Deduction CADE-24*, Maria Paola Bonacina (Ed.). Lecture Notes in Computer Science, Vol. 7898. Springer, Berlin, 208–214. DOI: [http://dx.doi.org/10.1007/978-3-642-38574-2\\_14](http://dx.doi.org/10.1007/978-3-642-38574-2_14)
- Xitong Gao and George A. Constantinides. 2015. Numerical program optimization for high-level synthesis. In *Proceedings of the 2015 ACM / SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'15)*. ACM, New York, NY, 210–213.
- Khalil Ghorbal, Eric Goubault, and Sylvie Putot. 2010. A logical product approach to zonotope intersection. In *Computer Aided Verification*, Tayssir Touili, Byron Cook, and Paul Jackson (Eds.). Lecture Notes in Computer Science, Vol. 6174. Springer, Berlin, 212–226. DOI: [http://dx.doi.org/10.1007/978-3-642-14295-6\\_22](http://dx.doi.org/10.1007/978-3-642-14295-6_22)
- Benjamin Grégoire and Assia Mahboubi. 2005. Proving equalities in a commutative ring done right in coq. In *TPHOLs (Lecture Notes in Computer Science)*, Joe Hurd and Thomas F. Melham (Eds.), Vol. 3603. Springer, 98–113.
- Thomas C. Hales. 2006. Introduction to the flyspeck project. In *Mathematics, Algorithms, Proofs (Dagstuhl Seminar Proceedings)*, Thierry Coquand, Henri Lombardi, and Marie-Françoise Roy (Eds.). Dagstuhl, Germany.
- Leopold Haller, Alberto Griggio, Martin Brain, and Daniel Kroening. 2012. Deciding floating-point logic with systematic abstraction. In *Formal Methods in Computer-Aided Design (FMCAD)*. 131–140.
- David Handelman. 1988. Representing polynomials by positive linear functions on compact convex polyhedra. *Pacific J. Math.* 132, 1 (1988), 35–62. Retrieved from <http://projecteuclid.org/euclid.pjm/1102689794>
- John Harrison. 1996. HOL light: A tutorial introduction. In *FMCAD (Lecture Notes in Computer Science)*, Mandayam K. Srivas and Albert John Camilleri (Eds.), Vol. 1166. Springer, 265–269.
- John Harrison. 2000. Formal verification of floating point trigonometric functions. In *Formal Methods in Computer-Aided Design*, Jr. Hunt, Warren A. and Steven D. Johnson (Eds.). Lecture Notes in Computer Science, Vol. 1954. Springer, Berlin, 254–270. DOI: [http://dx.doi.org/10.1007/3-540-40922-X\\_14](http://dx.doi.org/10.1007/3-540-40922-X_14)
- N. J. Higham. 2002. *Accuracy and Stability of Numerical Algorithms: Second Edition*. Society for Industrial and Applied Mathematics.
2008. IEEE standard for floating-point arithmetic. *IEEE Std 754-2008* (Aug. 2008), 1–70. DOI: <http://dx.doi.org/10.1109/IEEESTD.2008.4610935>
- Bertrand Jeannot and Antoine Miné. 2009. Apron: A library of numerical abstract domains for static analysis. In *Computer Aided Verification*, Ahmed Bouajjani and Oded Maler (Eds.). Lecture Notes in Computer Science, Vol. 5643. Springer, Berlin, 661–667.
- Jean-Louis Krivine. 1964. Quelques propriétés des préordres dans les anneaux commutatifs unitaires. *Compt. Rend. Acad. Sci.* 258 (1964), 3417–3418.
- J. B. Lasserre. 2009. *Moments, Positive Polynomials and Their Applications*. Imperial College Press.
- Jean B. Lasserre. 2001. Global optimization with polynomials and the problem of moments. *SIAM J. Optim.* 11, 3 (2001), 796–817.
- Jean B. Lasserre. 2006. Convergent SDP-relaxations in polynomial optimization with sparsity. *SIAM J. Optim.* 17, 3 (2006), 822–843.
- Jean B. Lasserre. 2011. A new look at nonnegativity on closed sets and polynomial optimization. *SIAM J. Optim.* 21, 3 (2011), 864–885. DOI: <http://dx.doi.org/10.1137/100806990>
- Jean B. Lasserre and Mihai Putinar. 2010. Positivity and optimization for semi-algebraic functions. *SIAM J. Optim.* 20, 6 (2010), 3364–3383.
- Monique Laurent. 2009. Sums of squares, moment matrices and optimization over polynomials. In *Emerging Applications of Algebraic Geometry*. Springer, 157–270.
- Michael D. Linderman, Matthew Ho, David L. Dill, Teresa H. Meng, and Garry P. Nolan. 2010. Towards program optimization through automated analysis of numerical precision. In *Proceedings of the 8th Annual IEEE/ACM International Symposium on Code Generation and Optimization (CGO'10)*. ACM, New York, NY, 230–237.

- Victor Magron. 2013. *Formal Proofs for Global Optimization—Templates and Sums of Squares*. Ph.D. Dissertation. École Polytechnique.
- Victor Magron. 2014. NLCertify: A tool for formal nonlinear optimization. In *Mathematical Software ICMS 2014*, Hoon Hong and Chee Yap (Eds.). Lecture Notes in Computer Science, Vol. 8592. Springer, Berlin, 315–320.
- Victor Magron, Xavier Allamigeon, Stéphane Gaubert, and Benjamin Werner. 2015a. Certification of real inequalities: Templates and sums of squares. *Math. Program.* 151, 2 (2015), 477–506. DOI: <http://dx.doi.org/10.1007/s10107-014-0834-5>
- Victor Magron, Xavier Allamigeon, Stéphane Gaubert, and Benjamin Werner. 2015b. Formal proofs for nonlinear optimization. *J. Formal. Reason.* 8, 1 (2015), 1–24.
- Maxima. 2013. Maxima, a Computer Algebra System. Version 5.30.0. (2013). Retrieved from <http://maxima.sourceforge.net/>.
- W. M. McEneaney. 2006. *Max-plus Methods for Nonlinear Control and Estimation*. Birkhäuser Boston Inc., Boston, MA.
- Guillaume Melquiond. 2012. Floating-point arithmetic in the coq system. *Information and Computation* 216, 0 (2012), 14–23. DOI: <http://dx.doi.org/10.1016/j.ic.2011.09.005> Special Issue: 8th Conference on Real Numbers and Computers.
- Antoine Miné. 2006. The octagon abstract domain. *Higher-Order Symbol. Comput.* 19, 1 (2006), 31–100. DOI: <http://dx.doi.org/10.1007/s10990-006-8609-1>
- Ramon E. Moore. 1962. *Interval Arithmetic and Automatic Error Analysis in Digital Computing*. Ph.D. Dissertation. Department of Computer Science, Stanford University.
- César Muñoz and Anthony Narkawicz. 2013. Formalization of a representation of bernstein polynomials and applications to global optimization. *J. Automat. Reason.* 51, 2 (August 2013), 151–196. DOI: <http://dx.doi.org/10.1007/s10817-012-9256-3>
- Nasa 2011. NASA World Wind Java SDK. Retrieved from <http://worldwind.arc.nasa.gov/java/>.
- OCaml 2015. Objective Caml (OCaml) programming language website. Retrieved from <http://caml.inria.fr/>.
- Gabriele Paganelli and Wolfgang Ahrendt. 2013. Verifying (in-)stability in floating-point programs by increasing precision, using SMT solving. In *Proceedings of the 15th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing, (SYNASC'13)*, Nikolaj Bjørner, Viorel Negru, Tetsuo Ida, Tudor Jebelean, Dana Petcu, Stephen M. Watt, and Daniela Zaharie (Eds.). IEEE Computer Society, 209–216. DOI: <http://dx.doi.org/10.1109/SYNASC.2013.35>
- Cordian Riemer, Thorsten Theobald, Lina Jansson Andrén, and Jean B. Lasserre. 2013. Exploiting symmetries in SDP-relaxations for polynomial optimization. *Math. Operat. Res.* 38, 1 (2013), 122–141. DOI: <http://dx.doi.org/10.1287/moor.1120.0558>
- Pierre Roux. 2015. Formal proofs of rounding error bounds. *J. Automat. Reason.* (2015), 1–22. DOI: <http://dx.doi.org/10.1007/s10817-015-9339-z>
- Philipp Rümmer and Thomas Wahl. 2010. An SMT-LIB theory of binary floating-point arithmetic. In *Informal proceedings of 8th International Workshop on Satisfiability Modulo Theories (SMT)*.
- Alexey Solov'yev and Thomas C. Hales. 2013. Formal verification of nonlinear inequalities with taylor interval approximations. In *Proceedings of the 5th International Symposium on NASA Formal Methods (NFM'13). Proceedings (Lecture Notes in Computer Science)*, Guillaume Brat, Neha Rungta, and Arnaud Venet (Eds.), Vol. 7871. Springer, 383–397.
- Alexey Solov'yev, Charles Jacobsen, Zvonimir Rakamarić, and Ganesh Gopalakrishnan. 2015. Rigorous estimation of floating-point round-off errors with symbolic taylor expansions. In *Proceedings of the 20th International Symposium on Formal Methods (FM) (Lecture Notes in Computer Science)*, Nikolaj Bjørner and Frank de Boer (Eds.), Vol. 9109. Springer, 532–550.
- J. Stolfi and L. H. de Figueiredo. 2003. An introduction to affine arithmetic. *TEMA Tend. Mat. Apl. Comput.* 4, 3 (2003), 297–312.
- M. J. Todd. 2001. Semidefinite optimization. *Acta Numer.* 10 (2001), 515–560.
- Lieven Vandenbergh and Stephen Boyd. 1994. Semidefinite programming. *SIAM Rev.* 38 (1994), 49–95.
- Hayato Waki, Sunyoung Kim, Masakazu Kojima, and Masakazu Muramatsu. 2006. Sums of squares and semidefinite programming relaxations for polynomial optimization problems with structured sparsity. *SIAM J. Optim.* 17, 1 (2006), 218–242.
- Makoto Yamashita, Katsuki Fujisawa, Kazuhide Nakata, Maho Nakata, Mituhiro Fukuda, Kazuhiro Kobayashi, and Kazushige Goto. 2010. *A high-performance Software Package for Semidefinite Programs: SDPA7*. Technical Report. Dept. of Information Sciences, Tokyo Institute of Technology, Tokyo, Japan. Retrieved from [http://www.optimization-online.org/DB\\_FILE/2010/01/2531.pdf](http://www.optimization-online.org/DB_FILE/2010/01/2531.pdf).

Received March 2016; revised October 2016; accepted November 2016