

# Formal verification of a controller implementation in fixed-point arithmetic

Lars Flessing, Grigory Devadze, Stefan Streif\*

**Abstract**—For the implementations of controllers on digital processors, certain limitations, e.g. in the instruction set and register length, need to be taken into account, especially for safety-critical applications. This work aims to provide a computer-certified inductive definition for the control functions that are implemented on such processors accompanied with the fixed-point data type in a proof assistant. Using these inductive definitions we formally ensure correct realization of the controllers on a digital processor. Our results guarantee overflow-free computations of the implemented control algorithm. The method presented in this paper currently supports functions that are defined as polynomials within an arbitrary fixed-point structure. We demonstrate the verification process in the case study on an example with different scenarios of fixed-point type implementations.

## I. INTRODUCTION

Verification is an important part in the development of control systems. Since some control systems are used in safety-critical applications, such as autonomous driving [19] or avionics [7], [12], it is essential to ensure that they meet certain performance and safety requirements. One key challenge is to guarantee that the properties of the control algorithm which are usually proven mathematically, are retained after numerical implementation in software [14]. Furthermore, certain standards must be met by the final realization i.e. the software which runs on the digital processor.

Conventionally, it is often assumed that the numerical imperfections are captured by the input error or system uncertainty which are then dealt within a robust control setting [25]. However, in the case of an overflow within the calculations, the resulting values depend on the realization of the controller itself and may not be bounded [32]. In the so-called algorithmic control theory [28], the importance of dealing with these imperfections by suitable formal approaches has been highlighted. The present work contributes toward these goals and aims to guarantee the correctness of the implementation of the controllers by bridging the gap between mathematical notion and the actual software implementation. This approach differs from the classical numerical analysis, where it is often unclear whether the used concepts are connected to the actual software code. In contrast to that, we derive numerical and computational concepts within a special so-called type-theoretic realization. In particular, the present work relies on the use of the proof assistant software and the so-called *proofs-as-programs paradigm*.

\* The authors are with the Automatic Control and System Dynamics Lab, Technische Universität Chemnitz, Chemnitz, Germany {lars.flessing, grigory.devadze, stefan.streif}@etit.tu-chemnitz.de.

## II. RELATED WORKS

The main concern for certain safety-critical applications is the verification of performance characteristics for a given controller design. In particular, such verification can be done using closed-loop simulation [16]. This method provides the ability to test a predefined range of scenarios which are relevant for the behavior of the systems. A common method for the verification purposes is the so-called *fault injection* into the simulation system [21]. Yang et al. use a fault injection strategy for traction control systems in the simulation of a high-speed trains [30]. In [31], the authors extend this approach to the field tests where a hardware-in-the-loop solution mimics hardware faults.

Another solution to ensure safety-critical characteristics of a control system is to add fail-safe systems and fault modes, which is known as *fault-tolerant control* [15]. For example, the control function can be divided into different modes that handle different performance requirements. These additional modes can be used to specifically ensure that the desired requirements are met. As an example, in [29] the authors introduce an error mode for a linear time-invariant control system that suffers from a loss of control input. The error mode is defined as a separate control function which has to ensure system stability. Blanke et al. [4] presented a variety of approaches, e.g. fault tolerant  $\mathcal{H}_\infty$ -design or fault tolerant model fitting designs. Keel et al. discussed inaccuracy effects of linear robust controller. The authors introduced the notion of fragility that is the loss of the controller performance caused by the deviation of the control parameters and defined respective metrics to estimate the fragility of a robust controller [17].

One may aim to analyze correctness and error-proneness via statistical approaches. In [6], the authors use uncertainty quantification [26] to make a qualitative statement about the computational correctness. This approach has been used by Micheltmore et al. to provide statistical guarantees for autonomous vehicle control [22] or by Berning et al. for random constrained path planning of unmanned aircraft [2]. However, these approaches are more of numerical nature and do not have a formal connection to the actual software realization, i.e. all approaches assume that the underlying controllers are implemented correctly.

To address the software implementation and obtain specific certification there is no way around of the usage of formal methods. Formal methods can be used to assert a guarantee for a certain mathematical property within a logical system and prove that this property holds.

The core of several tools is a *satisfiability modulo theories*

(SMT) solver [1] that is able to check a set of first-order logical statement. Bessa et al. presents a framework that converts the source code of the digital controller with a processor definition into a set of logical instructions [3]. The tool *RoSa* uses a combination of SMT solving and interval arithmetic to derive finite precision error bounds via verification constraints [5].

Herencia-Zapana et al. use source code annotations over Lyapunov-stabilizing properties of the code lines described via the quadratic invariants [13]. With these annotations, the authors obtain a translation of the code into the representation of the associated logical statements, which are checked within the prototype verification system PVS [23]. Alternatively, Roozbehani et al. proposed a control-theoretic framework for the verification of numerical programs [24]. The authors rely on the framework of the so-called Lyapunov invariants and employ the convex optimization for their construction. Recently, Devadze et al. used computer-assisted proofs and program extraction techniques for the formalization of initial value problem for ODE and sum-of-squares certificates for stability of polynomial systems [8], [9]. An alternative direction of algorithmic verification is *FLUCTUAT*, which allows to propagate the errors of function on variables with an uncertainty margin using statistical analysis. Finally, we mention the framework *FPTaylor*, where roundoff errors are estimated by a symbolic Taylor expansion with the help of HOL [27] and the library *Real2Float*, where one is able to estimate an upper error bound for floating-point programs [20].

#### A. Contributions and outline of this work.

The focus of the present work is to provide a formal framework for verified controller implementations in which computations are guaranteed to be overflow-free (*OF*). This is obtained by modeling the behavior of the processor hardware as a logical type of the used mathematical interpretation. Thus we may derive a mathematical representation of the effects of the hardware definition, such as the overflow effect that is presented in this paper. We extensively use formal methods and the proof-as-programs paradigm for the sake of correctness. Section III introduces the problem setup and the motivating example, where the occurrence of the overflows have destabilizing effects on system behavior. For demonstration purposes the example is carried out with the simulation of different fixed-point type settings to show the essence of the verification process. Section IV presents the proof assistant *Minlog* and provides the formal definitions for the data types that model the processor arithmetic. The derivation of the correct controller setup for the overflow-free computation is given in Section V. In Section VI we demonstrate the verification process by revisiting the example from the Section III and show that the certain controller implementation, which led to the actual failure, did not satisfy the formal definition. Finally, we provide the conclusion and an outlook in Section VII.

### III. PROBLEM SETUP

We begin with a motivational example.

#### A. Controller implementation

We consider the implementation of a controller on a digital processor (e.g. micro-controller) using fixed-point arithmetic. Since the digital processor provides only a limited number of bits for the representation, the implementation is affected by the so-called *finite word length problem*. For a fixed-point specification with  $q \in \mathbb{N}$  bits for the integer part and  $p \in \mathbb{N}$  bits for the fraction part the respective value domain of the controller is defined by

$$D_V := \{x_1, x_2 \in \mathbb{R} : \\ - (2^q + 1 - 2^{-p}) \leq x_1, x_2 \leq 2^q + 1 - 2^{-p} \wedge \\ 2^{-p} \leq |x_1 - x_2|\}. \quad (1)$$

The sign for a value  $x \in \mathbb{R}$  is realized by an additional bit that is either 0 for  $0 \leq x$  or 1 for  $x < 0$ . Violation of the value domain is called *overflow* which may be caused during computations within the control algorithm. For this example the overflow will be handled by wrapping since this type of implementations works without any additional logical tests. The discretization is handled with a *closest value* method that rounds the real value to the closest fixed-point value. Not all digital processor provide overflow processing routines and some only support integer operations or even logical operation as well as the register shifts [10]. The aim of this work is to verify that a control algorithm does not produce overflow. The following example demonstrates that an overflow might lead to significant degradation of the controller performance.

#### B. Example

Consider the following unstable plant:

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}, b = \begin{bmatrix} 2 \\ 0 \end{bmatrix}, c = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}^\top, d = 0 \quad (2)$$

with the transfer function

$$G_P(s) = \frac{s+1}{s^2-s-1}.$$

Furthermore assume that this plant has an additive uncertainty that is compensated with the weighting function  $W_1(s) = \frac{2s+1.1}{-10s^2-10.3s-5}$  and a multiplicative uncertainty that can be compensated with the weighting function  $W_2(s) = -3$ . To obtain a robust controller with respect to these uncertainties, we use the  $\mathcal{H}_\infty$  synthesis [11] that results in the following transfer function:

$$G_C(s) = \frac{2290s^3 + 3775s^2 + 2603s + 707.8}{s^4 + 198.7s^3 - 1332s^2 - 1483s - 767.9}.$$

We convert the controller transfer function  $G_C(s)$  into a state space representation. With respect to the implementation context we restrict the range of values for the matrices  $A \in D_V^{4 \times 4}, b \in D_V^4, c \in D_V^{1 \times 4}$ . For the controller states  $z \in D_V^4$  we specify the discrete-time implementation of the

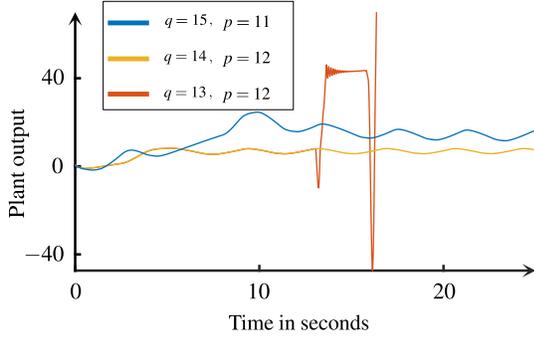


Fig. 1. Simulation result of the system output with different integer and fraction part lengths.

controller with system output  $y$  sampled at time  $h \cdot k$ ,  $k \in \mathbb{N}$  and time-step size  $h = 0.0001s$ :

$$z_{k+1} = z_k + h \begin{bmatrix} -1.03 & -0.5 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0.0511 & 0.0216 & -513.9303 & 510.7020 \\ 0.0314 & 0.0122 & -315.2434 & 316.2515 \end{bmatrix} z_k + h \begin{bmatrix} 5 \\ 0 \\ -1023.4 \\ -632.5 \end{bmatrix} y_k$$

$$u_{k+1} = [0.0002 \quad 0.001 \quad -1.6189 \quad -1.0018] z_k \quad (3)$$

We simulate the closed-loop system with different values for  $p$  and  $q$ . Figure 1 shows the result for an initial value for the system of  $x_0 = [0, 1]^T$  and  $z_0 = 0$  for the controller. The simulations only differ in the fixed-point specification. The blue and yellow lines show successful stabilization. They differ in the number of fraction bits but are similar in number of integer bits. The orange and yellow lines have the same number of fraction bits, but reducing the the integer bits form 14 to 13 the controller to stops stabilizing the system after  $t = 13s$ . This is caused by several overflow errors, which lead to the controller failure.

In the following sections we provide the approach to guarantee that a controller implementation will be overflow-free within the fixed-point arithmetic.

#### IV. PROOF ASSISTANT SYSTEM AND FORMAL DEFINITIONS

A partial goal of this work is to model the physical limitations of bit streams in *Minlog*. Due to the functional nature of the modelling, *Minlog* offers a solid set of tools for our purposes.

##### A. Introduction to *Minlog*

We introduce some preliminary concepts of the proof system *Minlog*. Roughly speaking, a proof assistant is a tool which is capable to address mathematical theory programmatically. Within *Minlog* we are able to define propositions  $P, Q$  to be arbitrary formulas and  $\rightarrow$  the logical implication. With these symbols we can state

- $P \rightarrow Q$ ,  $Q$  can be derived from  $P$
- $\forall \alpha. P(\alpha)$ , the proposition  $P$  holds for all  $\alpha$ .

That way, we can define base types by constructor functions, further on called *constructor*, that are either constant or recursive calls.

For example consider the natural numbers, that can be defined by `zero` and a successor function  $\text{Succ} := \mathbb{N} \Rightarrow \mathbb{N}$ , mapping from the natural numbers into the natural numbers. It is also possible to define types by other types. For instance the rational numbers, defined by numerator and denominator, may be represented by  $\text{RatConst} := \mathbb{Z} \Rightarrow \mathbb{P} \Rightarrow \mathbb{Q}$ .

With the type definitions and operators, proofs can be stated in the given logic and reasoned. This is the setting a proof assistant offers.

Functions that are defined in *Minlog* are called *program constant*. Since proofs and programs are equivalent, it is apparent that a program constant is proof of the correctness of the function [18]. This forces the program constants to guarantee an output in the final type in a finite time. This principle is called the *notion of totality* and it ensures that the program will always yield correct results effectively. In the present work, the totality of all program constants have been implemented.

To prove theorems with *Minlog*, assumptions need to be provided which are relevant for the problem statement, e.g. by introducing propositions. In *Minlog* it is also possible to introduce the so-called inductively defined predicate constant (IDPC). Informally speaking, these are assumptions specific to a variable. In following we separate the usual mathematical notation of numbers from the functional definitions within *Minlog* by using bold letters rather than the double stricken. For instance, the type realization for the rational numbers  $\mathbb{Q}$  is denoted by  $\mathbf{Q}$ .

##### B. Data representation types

To model the behavior of the data in a finite storage processing system, we need to encapsulate the essential properties of this type. Thus a representation for one and zero, a terminal symbol to limit the variable, and a terminal symbol marking an overflow. We introduce the type `fp` that we call a binary stream.

*Definition 1 (Binary Stream):* The type `fp` is defined by

- $\circ := \text{fp}$ , the end of a bit stream, denoted by  $\circ$ ,
- $\emptyset := \text{fp}$ , the end of a bit stream with an overflow, denoted by  $\emptyset$ ,
- $1 := \text{fp} \Rightarrow \text{fp}$ , represents 1, and
- $0 := \text{fp} \Rightarrow \text{fp}$ , represents 0

and is denoted by the variable  $v$ .

*Remark 2:* With these functions, we are capable to represent a generic bit-stream and mark whether it will return a correct value. A variable that has an overflow is called *dirty*. Checking for the overflow is realized via the program constant `FpNotOVL` which is defined such that `FpNotOVL` returns false if the terminal symbol of the variable is `OVL` or `T` when the variable is clean.

The main property of such bit-stream is the limitation of the length. Program constants need to take this into account. Thus, we define two versions for all program constants either

limit the variable length, denoted with a superscript ‘‘A’’, or not. Only limiting program constants can cause an overflow. With the function  $\text{FpLength} := \text{fp} \Rightarrow \mathbb{N}$ , denoted as  $\langle v \rangle$  the length of a variable can be obtained.

We now present an interpretation of the data type  $\text{fp}$  that implements the fixed-point numbers. The constructor functions themselves represent the sign of the typed variable while the arguments represent the fraction and integer part of the fixed-point type.

*Definition 3 (Signed Fixed Point Type):* The signed fixed point type  $\text{sfp}$  is defined by

$$\begin{aligned} \text{FpPos} &:= \text{fp} \Rightarrow \text{fp} \Rightarrow \text{sfp} \\ \text{FpNeg} &:= \text{fp} \Rightarrow \text{fp} \Rightarrow \text{sfp} \end{aligned}$$

and is denoted by the variable  $w$ . The constructor  $\text{FpPos}$  is represented with  $+(v_1, v_2)$  while  $\text{FpNeg}$  is read as  $\sim(v_1, v_2)$ , accordingly, where  $v_1$  is the fraction, and  $v_2$  integer.

The multiplication in the  $\text{sfp}$  type is defined by a composure of limiting and nonlimiting program constants of the  $\text{fp}$  type.

*Definition 4 (SfpTimes):* The multiplication of the type  $\text{sfp} + (v_5, v_6) = +(v_1, v_2) \cdot +(v_3, v_4)$  is defined by

$$\begin{aligned} v_5 &= \nu[(v_1 \cdot_{\text{fp}} v_4) +_{\text{fp}} (v_2 \cdot_{\text{fp}} v_3)] +_{\text{fp}} \\ &\quad \iota(v_1 \cdot_{\text{fp}} v_3, |v_2|, |v_2|), |v_1|] \\ v_6 &= (v_2 \cdot_{\text{fp}}^A v_4 +_{\text{fp}}^A \\ &\quad \iota((v_1 \cdot_{\text{fp}} v_4) +_{\text{fp}} (v_2 \cdot_{\text{fp}} v_3)), |v_1|, |v_2|)) \end{aligned}$$

for fraction and integer part. The function  $\iota := \text{sfp} \Rightarrow \mathbb{N} \Rightarrow \mathbb{N} \Rightarrow \text{fp}$  selects a part of a bit-stream defined by the starting bit and amount of selected bits.  $\nu := \text{fp} \Rightarrow \mathbb{N} \Rightarrow \text{fp}$  selects a given number of bits, but starting at the first bit.

To connect the type with normal number definitions conversion functions are necessary. Using  $\text{Minlog}$  machinery we were able to prove the following technical theorems.

*Theorem 5 (SfpToRatToSfpId):* The conversion of a variable of the type  $\text{sfp}$  into the type of the rationals is without loss of accuracy. That is

$$\forall w. \eta_{\text{sfp}}^{\mathbf{Q}}(\eta_{\mathbf{Q}}^{\text{sfp}}(w)) = w.$$

*Theorem 6 (RatToSfpToRatError):* The conversion of a variable  $a$  of the type  $\mathbf{Q}$  into a variable of the type  $\text{sfp}$  has, for the function  $\eta_{\mathbf{Q}}^{\text{sfp}}(a, p, q)$ , an error bound of

$$\forall a. |a| \leq 2^q \rightarrow |\eta_{\text{sfp}}^{\mathbf{Q}}(\eta_{\mathbf{Q}}^{\text{sfp}}(a)) - a| \leq \frac{1}{2^{p+1}}$$

with  $p$  the amount of fraction bits and  $q$  the amount of integer bits.

The errors of the conversion in the operations are estimated as follows.

*Theorem 7 (SfpPlusError):* For the summation of two rational numbers  $a_1, a_2$  in the type  $\text{sfp}$  with  $p$  fraction bits and  $q$  integer bits the following property holds:

$$\begin{aligned} \forall a_1, a_2. |a_1 + a_2| \leq 2^q \rightarrow \\ |(\eta_{\text{sfp}}^{\mathbf{Q}}(\eta_{\mathbf{Q}}^{\text{sfp}}(a_1) +_{\text{sfp}} \eta_{\mathbf{Q}}^{\text{sfp}}(a_2)))| \leq 2\delta_{\text{sfp}} \end{aligned}$$

with  $\delta_{\text{sfp}} = 2^{-p}$ .

*Proof:* We use the the error estimates (5) and (6).

$$\begin{aligned} \forall a_1, a_2. |a_1 + a_2| \leq 2^q \rightarrow \\ |(\eta_{\text{sfp}}^{\mathbf{Q}}(\eta_{\mathbf{Q}}^{\text{sfp}}(a_1) +_{\text{sfp}} \eta_{\mathbf{Q}}^{\text{sfp}}(a_2)) - (a_1 + a_2))| \leq \\ |a_1 + \delta_{\text{sfp}} + a_2 + \delta_{\text{sfp}} - (a_1 + a_2)| = |2\delta_{\text{sfp}}| \end{aligned}$$

*Theorem 8 (SfpTimesError):* For the multiplication of two rational numbers  $a_1, a_2$  in the type  $\text{sfp}$  with  $p$  fraction bits and  $q$  integer bits the following property holds:

$$\begin{aligned} \forall a_1, a_2. |a_1 \cdot a_2| \leq 2^q \rightarrow \\ |\eta_{\text{sfp}}^{\mathbf{Q}}(\eta_{\mathbf{Q}}^{\text{sfp}}(a_1) \cdot \eta_{\mathbf{Q}}^{\text{sfp}}(a_2)) - (a_1 \cdot a_2)| \leq \\ |(a_1 + a_2)|\delta_{\text{sfp}} + \delta_{\text{sfp}}^2 \end{aligned} \quad (4)$$

with  $\delta_{\text{sfp}} = 2^{-p}$ .

*Proof:* By Definitions 4, Theorem 6 and under assumption  $|a_1 \cdot a_2| \leq 2^q$  we can write

$$\begin{aligned} |\eta_{\text{sfp}}^{\mathbf{Q}}(\eta_{\mathbf{Q}}^{\text{sfp}}(a_1) \cdot \eta_{\mathbf{Q}}^{\text{sfp}}(a_2)) - (a_1 \cdot a_2)| \\ \leq |(a_1 + \delta_{\text{sfp}}) \cdot (a_2 + \delta_{\text{sfp}}) - a_1 \cdot a_2| \\ = |(a_1 + a_2)|\delta_{\text{sfp}} + \delta_{\text{sfp}}^2 \end{aligned}$$

In the next section we derive the domain of the implementation that guarantees the overflow free computation of controller implementations.

## V. MAIN RESULTS

Using the preliminary definitions from Section IV the main result is provided in Theorem 15. We formally derive the property and proof that the computations within the given so-called reliable domain can be bounded.

We define a domain for the arguments of the implementation where we can assure that no overflow occurs, calling the results *overflow free*, OF. This domain is dependent on the sequence of operations of the implementation.

*Definition 9 (Sequenced Algorithm):* A sequenced algorithm  $\Gamma(w)$  is a function where each binary base operation  $+, -, \cdot$ , further depicted by  $\circ$ , is assigned a value  $i \in \mathbb{N}$  that represents the place in the order of execution. We further on assume that  $\Gamma(w)$  is a function in fixed-point numbers  $D_V$  that are defined by the number of integer bits and the number of fraction bits and  $w \in D_V$ .

We now define the domain of the arguments where the results of the algorithm are correct.

*Definition 10 (Reliable Domain):* For a fixed-point type the implementation of the domain

$$D_R = \{w : \Gamma_{\min} \leq \Gamma(w) \leq \Gamma_{\max}\}.$$

is called *reliable domain*.

For a specific fixed-point type definition, with  $q$  integer bits and  $p$  fraction bits, we can find, with Definition 10 a minimum and a maximum value for  $\Gamma(w)$ .

*Theorem 11 (Bounded domain configuration):* For a sequenced algorithm  $\Gamma(w)$  with  $n$  operations  $\circ_i$  with the

second argument given by  $\gamma_i \in D_V$  and for a given  $w \in D_V$  and  $p, q \in \mathbb{N}$  it holds that

$$\forall w \in D_V, -2^q \leq \Gamma_{\min} \leq \Gamma(w) \leq \Gamma_{\max} \leq 2^q.$$

Furthermore, the domain  $D_V$  is dependent on  $\Gamma(w)$ ,  $\Gamma_{\min}$  and  $\Gamma_{\max}$ .

*Proof:* We proof by induction over  $n$ . We proceed by case distinction on arithmetical operations. The first case is for (+):

$$\begin{aligned} -2^q &\leq \gamma_1 + w \leq 2^q \\ D_{V,\min} = 2^q - \gamma_1 &\leq w \leq 2^q - \gamma_1 = D_{V,\max}. \end{aligned}$$

The second case for the operation ( $\cdot$ ):

$$\begin{aligned} -2^q &\leq \gamma_1 \cdot w \leq 2^q \\ D_{V,\min} = \frac{-2^q}{\gamma_1} &\leq w \leq \frac{2^q}{\gamma_1} = D_{V,\max}. \end{aligned}$$

With the limitation of  $D_V$  we obtain the maximum and minimum value for the implementation  $\Gamma_{\min} = \Gamma(D_{V,\min})$  and  $\Gamma_{\max} = \Gamma(D_{V,\max})$ .

In the induction step we again need to separate the two cases. For (+) the lower bound is

$$D_{V,\min} \leq w + \Gamma(w) \leq w + \Gamma_{\min} \rightarrow D_{V,\min} - \Gamma_{\min} \leq w$$

and for the upper bound

$$w + \Gamma(w) \leq w + \Gamma_{\max} \leq D_{V,\max} \rightarrow w \leq D_{V,\max} - \Gamma_{\max}.$$

For ( $\cdot$ ) and a negative  $\Gamma_{\min}$  we obtain for the lower bound

$$D_{V,\min} \leq w \cdot \Gamma(w) \leq w \cdot \Gamma_{\min} \rightarrow w \leq \frac{D_{V,\min}}{\Gamma_{\min}}.$$

and for the upper bound

$$w \cdot \Gamma(w) \leq w \cdot \Gamma_{\max} \leq D_{V,\max} \rightarrow \frac{D_{V,\max}}{\Gamma_{\max}} \leq w.$$

Additionally the values  $\Gamma_{\min}, \Gamma_{\max}$  can be influenced by the parameters. We follow the induction again, this time we assume operations without the argument  $w$ . We have one case for (+) since the case for ( $\cdot$ ) is trivial.

$$\Gamma_{\min} = -(2^q + 1 - 2^{-p}) \leq \gamma_1 \leq (2^q + 1 - 2^{-p}) = \Gamma_{\max}$$

For the induction step we have again two cases. For (+) we have for the upper bound

$$\Gamma_{\min,i-1} \leq \gamma_i + \Gamma(w) \leq \gamma_i + \Gamma_{\min,i} \rightarrow \Gamma_{\min,i-1} - \gamma_i \leq \Gamma_{\min,i}.$$

For the lower bound we have

$$\gamma_i + \Gamma(w) \leq \gamma_i + \Gamma_{\max,i} \leq \Gamma_{\max,i-1} \rightarrow \Gamma_{\max,i} \leq \Gamma_{\max,i-1} - \gamma_i.$$

For ( $\cdot$ ) we have and under assumption that  $0 \leq |\gamma_i|$  we have for the upper bound

$$\Gamma_{\min,i-1} \leq \gamma_i \cdot \Gamma(w) \leq \gamma_i \cdot \Gamma_{\min,i} \rightarrow \frac{\Gamma_{\min,i-1}}{\gamma_i} \leq \Gamma_{\min,i}.$$

For the lower bound we have

$$\gamma_i \cdot \Gamma(w) \leq \gamma_i \cdot \Gamma_{\max,i} \leq \Gamma_{\max,i-1} \rightarrow \Gamma_{\max,i} \leq \frac{\Gamma_{\max,i-1}}{\gamma_i}$$

With the Definition 10 and the Theorem 11 we can now formalize the type `imp`.

*Definition 12 (Implementation):* The type `imp` is defined by the constructor

$$\text{imp} := \text{sfp} \Rightarrow \text{sfp} \Rightarrow (\text{sfp} \Rightarrow \text{sfp})$$

The first two variables are the bounds of the reliable set  $D_R$ , followed by the specification of the sequenced algorithm  $\Gamma := \text{sfp} \Rightarrow \text{sfp}$ .

We now further define the type with an IDPC that ensures correct processing of the given algorithm.

*Definition 13 (Reliable Implementation):* An implementation  $\text{imp} := w_1 \Rightarrow w_2 \Rightarrow \Gamma(w)$  is called reliable if

$$\forall w. w_1 \leq w \leq w_2 \rightarrow \Gamma_{\min} \leq \Gamma(w) \leq \Gamma_{\max} \rightarrow \emptyset(\Gamma(w)).$$

*Remark 14:* This ensures that the reliable domain is stated within the constructor and not an arbitrary set.

We will now show the correctness of this type of implementation.

*Theorem 15 (ImpRelCorr):* For a reliable implementation  $\Gamma(w)$ , with the operation parameters  $\gamma_i$  and  $\gamma_m = \max |\gamma_i| \in D_V$ , the error to the rational valued polynomial  $g(a)$  is a bounded function of the argument  $a$  if  $\eta_{fp}^Q(a) \in D_R$ .

*Proof:* The error of the function is defined by

$$\delta = \left| g(a) - \eta_{fp}^Q \left( \eta_{fp}^{fp} (g(a)) \right) \right|.$$

We can estimate the error with Theorem 7 and 8 for a  $w \in D_R$  with  $\delta_c$  the conversion error and  $\delta_t$  the summation error.

$$\begin{aligned} \delta &= \left| \sum_{i=0}^n \gamma_i a^i - \eta_{fp}^Q \left( \eta_{fp}^{fp} \left( \sum_{i=0}^n \gamma_i a^i \right) \right) \right| \\ &\leq \left| \sum_{i=0}^n \gamma_m a^i - ((\gamma_m + \delta_c)(x + \delta_c)^i + \delta_t) \right| \\ &\leq \left| \sum_{i=0}^{n+1} \gamma_m \left( \sum_{k=1}^n \binom{n}{k} x^{n-k} \delta_c^k \right) + \delta_c (a + \delta_c)^n + \delta_t \right| \\ &= \left| \frac{n^2 + n}{2} \left( \delta_c a^n + (\gamma_m + \delta_c) \left( \sum_{k=1}^n \binom{n}{k} a^{n-k} \delta_c^k \right) + \delta_t \right) \right|. \end{aligned}$$

This is a function of  $a$  and  $\Gamma$ . ■

*Remark 16:* If the  $w \notin D_R$  the implementation  $\Gamma(w)$  will have an overflow somewhere within the computation steps. This reduces the result of  $\Gamma(w)$  to a arbitrary value that might not stabilize the system anymore.

## VI. CASE STUDY

To illustrate the formalization, we will show that the failing example from Section III is not fulfilling Definition 13. When we use the maximum controller state  $z_{max}$  during the stabilization process we can estimate the Equation (3) with an upper bound:

$$\begin{aligned} u_{k+1} &\leq c_C^\top (A_C z_{max} + b_C y(t)) \\ &= c_C^\top A_C z_{max} + (c_C^\top b_C) y(t). \end{aligned}$$

We estimate these values for the different fixed-point type specifications and use them to obtain the bounds for the reliable domain  $D_R$ . These are also dependent on the value

TABLE I

THE SIMULATION PARAMETERS WITH THE RESULTING BOUND OF THE RELIABLE DOMAIN.

	$p$	$q$	$D_{R,\min}$	$D_{R,\max}$
1	11	15	-1.0352	27.5781
2	12	14	-2.5435	11.7632
3	12	13	1.0332	8.1863

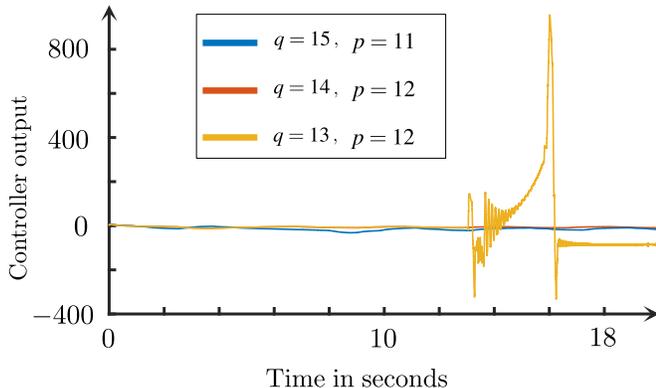


Fig. 2. The output of the controller from Figure 2 with their different fixed-point realizations.

domain and shown in Table I for the different specifications in Section III. We see that the first implementation contains the control target  $y = 0$  and is capable of yielding correct results. This is also shown in Figure 2. For the yellow and red line the difference is that the control target is not within the reliable domain of the yellow simulation. This lead to the destabilization of the feedback loop due to overflows. What is interesting is that the controller performed reasonably well before the degradation occurred  $t = 13s$ . It is believed that such effects may be overlooked when computations become more complex.

Thus, with this formalization of an implementation it becomes possible to check a priori whether the used specification for the fixed-point type is sufficient to achieve the control goals. This way one may guarantee that the computations are overflow-free within the given domain and under the given assumptions, for all values of the domain. In the case of the examples one may see that one additional bit for the integer part of the fixed-point type is sufficient to prevent overflows in the controller computation. With the formalization this can be detected and the implementation can be adjusted to ensure that it is reliable.

## VII. CONCLUSION AND OUTLOOK

Even in case of robustly stabilizing controllers, the final software implementation of the control algorithms on digital processors can become unstable due to the overflow of the computations using fixed-point arithmetic. In this work we have presented a formal method to verify that fixed-point representations on a processor are overflow free. This was achieved by a formalization of the problem and by creating a formal type that incorporates the limitations of the digital processor. For the derivations, the proof assistant system

Minlog was used. The performance of the formally verified controllers was demonstrated in a case study.

One advantage of the presented approach is that it allows guaranteeing correctness of the mathematical derivation and of the software implementation simultaneously. While there are methods for overflow handling it is more important to ensure that they can be avoided. This then leads to smaller source code implementations and faster execution times. Future work will consider more complex function evaluations as well as the predictive controllers.

## REFERENCES

- [1] C. Barrett and C. Tinelli. Satisfiability modulo theories. In *Handbook of Model Checking*, pages 305–343. Springer, 2018.
- [2] A. W. Berning, A. Girard, I. Kolmanovsky, and S. N. D’Souza. Rapid uncertainty propagation and chance-constrained path planning for small unmanned aerial vehicles. *Advanced Control for Applications: Engineering and Industrial Systems*, 2(1):e23, 2020.
- [3] I. Bessa, H. Ismail, L. Cordeiro, and J. E. Filho. Verification of fixed-point digital controllers using direct and delta forms realizations. *Design Automation for Embedded Systems*, 20, 06 2016.
- [4] M. Blanke, M. Kinnaert, J. Lunze, M. Staroswiecki, and J. Schröder. *Diagnosis and fault-tolerant control*, volume 2. Springer, 2006.
- [5] E. Darulova and V. Kuncak. Towards a compiler for reals. *ACM Trans. Program. Lang. Syst.*, 39(2), Mar. 2017.
- [6] M. L. Daza-Torres, J. C. Montesinos-López, M. A. Capistrán, J. A. Christen, and H. Haario. Error control in the numerical posterior distribution in the bayesian UQ analysis of a semilinear evolution PDE. *International Journal for Uncertainty Quantification*, 2020.
- [7] D. Delmas, E. Goubault, S. Putot, J. Souyris, K. Tekkal, and F. Védrine. Towards an industrial use of fluctuat on safety-critical avionics software. In *International Workshop on Formal Methods for Industrial Critical Systems*, pages 53–69. Springer, 2009.
- [8] G. Devadze, L. Flessing, and S. Streif. Extraction of a computer-certified ODE solver. *Accepted for publication on European control conference 2021. arXiv preprint*, 2021.
- [9] G. Devadze, V. Magron, and S. Streif. Computer-assisted proofs for Lyapunov stability via Sums of Squares certificates and Constructive Analysis. *arXiv preprint arXiv:2006.09884*, 2020.
- [10] ELAN Microelectronics Corp. EM78P156K Datasheet. <http://www.emc.com.tw/emc/tw/Product/Product/detail/158>. Accessed : 18.2.2021.
- [11] K. Glover and J. C. Doyle. State-space formulae for all stabilizing controllers that satisfy an  $H_\infty$ -norm bound and relations to relations to risk sensitivity. *Systems & Control Letters*, 11(3):167–172, 1988.
- [12] A. E. Goodloe, C. Muñoz, F. Kirchner, and L. Correnson. Verification of numerical programs: From real numbers to floating point numbers. In *NASA Formal Methods Symposium*, pages 441–446. Springer, 2013.
- [13] H. Herencia-Zapana, R. Jobredeaux, S. Owre, P.-L. Garoche, E. Feron, G. Perez, and P. Ascariz. PVS linear algebra libraries for verification of control software algorithms in c/ACSL. In A. E. Goodloe and S. Person, editors, *NASA Formal Methods*, pages 147–161, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg.
- [14] G. J. Holzmann and M. H. Smith. Software model checking: Extracting verification models from source code. *Software Testing, Verification and Reliability*, 11(2):65–79, 2001.
- [15] J. Jiang and X. Yu. Fault-tolerant control systems: A comparative study between active and passive approaches. *Annual Reviews in Control*, 36(1):60–72, 2012.
- [16] J. Kapinski, J. V. Deshmukh, X. Jin, H. Ito, and K. Butts. Simulation-based approaches for verification of embedded control systems: An overview of traditional and advanced modeling, testing, and verification techniques. *IEEE Control Systems Magazine*, 36(6):45–64, 2016.
- [17] L. Keel and S. Bhattacharyya. Robust, Fragile, or Optimal? *Automatic Control, IEEE Transactions on*, 42:1098 – 1105, 09 1997.
- [18] J. Kennedy and R. Kossak. *Set Theory, Arithmetic, and Foundations of Mathematics: Theorems, Philosophies*, volume 36. Cambridge University Press, 2011.
- [19] J. Kong, M. Pfeiffer, G. Schildbach, and F. Borrelli. Kinematic and dynamic vehicle models for autonomous driving control design. In *2015 IEEE Intelligent Vehicles Symposium (IV)*, pages 1094–1099, 2015.

- [20] V. Magron, G. Constantinides, and A. Donaldson. Certified roundoff error bounds using semidefinite programming. *ACM Trans. Math. Softw.*, 43(4), Jan. 2017.
- [21] Mei-Chen Hsueh, T. K. Tsai, and R. K. Iyer. Fault injection techniques and tools. *Computer*, 30(4):75–82, 1997.
- [22] R. Michelmore, M. Wicker, L. Laurenti, L. Cardelli, Y. Gal, and M. Kwiatkowska. Uncertainty quantification with statistical guarantees in end-to-end autonomous driving control. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7344–7350, 2020.
- [23] S. Owre, J. M. Rushby, and N. Shankar. PVS: A prototype verification system. In *International Conference on Automated Deduction*, pages 748–752. Springer, 1992.
- [24] M. Roozbehani, A. Megretski, and E. Feron. Optimization of Lyapunov Invariants in verification of software systems. *IEEE Transactions on Automatic Control*, 58(3):696–711, 2013.
- [25] P. Schülting, C. H. van der Broeck, and R. W. De Doncker. Analysis and design of repetitive controllers for applications in distorted distribution grids. *IEEE Transactions on Power Electronics*, 34(1):996–1004, 2019.
- [26] R. C. Smith. *Uncertainty quantification: theory, implementation, and applications*, volume 12. SIAM, 2013.
- [27] A. Solovyev, M. S. Baranowski, I. Briggs, C. Jacobsen, Z. Rakamarić, and G. Gopalakrishnan. Rigorous estimation of floating-point round-off errors with symbolic Taylor expansions. *ACM Trans. Program. Lang. Syst.*, 41(1), Dec. 2018.
- [28] P. Tsiotras and M. Mesbahi. Toward an algorithmic control theory. *Journal of Guidance, Control, and Dynamics*, 40(2):194–196, 2017.
- [29] W. Xiang, G. Zhai, and C. Briat. Stability analysis for LTI control systems with controller failures and its application in failure tolerant control. *IEEE Transactions on Automatic Control*, 61(3):811–816, 2016.
- [30] C. Yang, C. Yang, T. Peng, X. Yang, and W. Gui. A fault-injection strategy for traction drive control systems. *IEEE Transactions on Industrial Electronics*, 64(7):5719–5727, 2017.
- [31] X. Yang, C. Yang, T. Peng, Z. Chen, B. Liu, and W. Gui. Hardware-in-the-loop fault injection for traction control system. *IEEE Journal of Emerging and Selected Topics in Power Electronics*, 6(2):696–706, 2018.
- [32] R. Yates. Fixed-point arithmetic: An introduction. *Digital Signal Labs*, 81(83):198, 2009.