

# Reliable evaluation of the Worst-Case Peak Gain matrix in multiple precision

Anastasia Volkova, Thibault Hilaire, Christoph Lauter

► **To cite this version:**

Anastasia Volkova, Thibault Hilaire, Christoph Lauter. Reliable evaluation of the Worst-Case Peak Gain matrix in multiple precision. ARITH 22 - 22nd IEEE Symposium on Computer Arithmetic, Jun 2015, Lyon, France. pp.96-103, 10.1109/ARITH.2015.14 . hal-01083879v2

**HAL Id: hal-01083879**

**<https://hal.sorbonne-universite.fr/hal-01083879v2>**

Submitted on 10 Dec 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Reliable evaluation of the Worst-Case Peak Gain matrix in multiple precision

Anastasia Volkova

Thibault Hilaire

Christoph Lauter

**Abstract**—The worst-case peak gain (WCPG) of an LTI filter is an important measure for the implementation of signal processing algorithms. It is used in the error propagation analysis for filters, thus a reliable evaluation with controlled precision is required. The WCPG is computed as an infinite sum and has matrix powers in each summand. We propose a direct formula for the lower bound on truncation order of the infinite sum in dependency of desired truncation error. Several multiprecision methods for complex matrix operations are developed and their error analysis performed. We present a multiprecision complex matrix inversion algorithm using Newton-type iteration, along with its error analysis and proof of convergence. A multiprecision matrix powering method is presented. All methods yield a rigorous solution with an absolute error bounded by an a priori given value. The results are illustrated with numerical examples.

## INTRODUCTION

The majority of control and digital signal processing algorithms are dedicated to linear time-invariant systems with finite or infinite impulse response. Most of them are implemented for application in embedded systems, which use finite-precision arithmetic. Unfortunately, the quantification of coefficients and further roundoff errors lead to degradation of the algorithms. Therefore, an accurate error analysis of implementation of such algorithms is required.

However, this analysis is complicated by the non-linear propagation of errors through the filter as they are amplified on each step by internal state of the system. A solution is proposed in [8], based on a property of bounded-input bounded output systems [1], [2] where the largest possible peak value of the output is determined by the use of the Worst-Case Peak Gain (WCPG) matrix. Therefore, the analysis of error propagation in LTI systems is directly dependent on the reliable evaluation of the WCPG.

This measure is computed with an infinite sum and has matrix powers in each summand. These problems are both known to be non-trivial. In this article we propose a detailed algorithm for the reliable evaluation of the WCPG matrix with multiple precision. This algorithm ensures that the WCPG is computed with an absolute error rigorously bounded by an a priori given value  $\varepsilon$ . For these purposes several multiprecision algorithms for complex entries were developed. Our methods ensure that their absolute error of computations is rigorously bounded by an a priori given value. This is achieved by adapting the precision of intermediate computations and correct rounding. Therefore, we present not only the error analysis of the approximations made on each step of the WCPG computation, but we also deduce the required accuracy for our kernel multiprecision algorithms such that the overall error bound is satisfied.

We propose an analysis for the error induced by truncating the infinite sum and a direct formula for the computation of a lower bound on truncation order in dependency with the desired absolute error. The truncation order algorithm involves Interval Arithmetic computations and uses Theory of Verified Inclusions.

Some preliminary definitions about LTI systems are recalled in Section I. Section II describes the global algorithm used to reliably evaluate the WCPG matrix  $\mathbf{W}$ . The truncation order and the truncation error are analyzed in Section III. Section IV is focused on the different steps used for the summation and the associated error analysis, whereas Section V details some basic bricks in multiple precision. Finally, numerical examples are presented in Section VI before conclusion.

**Notation:** Throughout the article matrices are in uppercase boldface (for example  $\mathbf{A}$ ), vectors are in lowercase boldface (for example  $\mathbf{v}$ ), scalars are in lowercase (for example  $\alpha$ ). Operators  $\otimes$ , and  $\oplus$  denote floating-point (FP) multiplication and addition respectively,  $\mathbb{F}$  the set of FP numbers.  $[x]$  corresponds to an interval.  $\mathbf{A}^*$  denote the conjugate transpose of the matrix  $\mathbf{A}$ . All absolute values and inequalities with matrices are considered to be element-by-element, for example  $|\mathbf{A}| < |\mathbf{B}|$  denotes  $|A_{ij}| < |B_{ij}| \forall i, j$ . In addition,  $\mathbf{A} < \varepsilon$  denotes  $A_{ij} < \varepsilon \forall i, j$ .  $\mathbf{I}_n$  denotes the identity matrix of size  $n \times n$  and  $\rho(\mathbf{A})$  the spectral radius of  $\mathbf{A}$ .

## I. LTI FILTERS AND WORST-CASE PEAK GAIN

A Linear Time Invariant (LTI) filter is a system used in signal processing, image processing, control theory, etc. It is defined by an input-output relationship in time-domain or equivalently in frequency-domain. Linear controllers, Finite Impulse Response (FIR) filters, Infinite Impulse Response (IIR) are classical examples of LTI system. We focus here only on discrete-time systems: a discrete-time LTI system (filter) is a numerical application that transforms an input signal  $\{\mathbf{u}(k)\}_{k \geq 0}$  into an output signal  $\{\mathbf{y}(k)\}_{k \geq 0}$  ( $\mathbf{u}(k)$  and  $\mathbf{y}(k)$  may be vectors or scalars), where  $k \in \mathbb{N}$  is the step time.

Unlike a mathematical function, the output at time  $k$  depends not only on the input at time  $k$  but also on the internal state of the filter (generally determined from the previous inputs and outputs). A common input-output relationship is the state-space representation [9]. It describes the evolution of the state vector  $\mathbf{x}(k)$  from the previous step and the input:

$$\begin{cases} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \mathbf{y}(k) &= \mathbf{C}\mathbf{x}(k) + \mathbf{D}\mathbf{u}(k) \end{cases} \quad (1)$$

where  $\mathbf{u}(k) \in \mathbb{R}^{q \times 1}$  is the input vector,  $\mathbf{y}(k) \in \mathbb{R}^{p \times 1}$  the output vector,  $\mathbf{x}(k) \in \mathbb{R}^{n \times 1}$  the state vector and  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times q}$ ,  $\mathbf{C} \in \mathbb{R}^{p \times n}$  and  $\mathbf{D} \in \mathbb{R}^{p \times q}$  are the state-space matrices of the system.

**Proposition 1 (Bounded Input Bounded Output systems)**

Let  $\mathcal{H}$  be such a state-space system. If an input  $\{\mathbf{u}(k)\}_{k \geq 0}$  is known to be bounded by  $\bar{\mathbf{u}}$  ( $\forall k \geq 0, |\mathbf{u}_i(k)| \leq \bar{u}_i, 1 \leq i \leq q$ ), then the output  $\{\mathbf{y}(k)\}_{k \geq 0}$  will be bounded iff the spectral radius  $\rho(\mathbf{A})$  is strictly less than 1. This property is known as the Bounded Input Bounded Output (BIBO) stability [9].

Moreover, in that case, the output is (component-wise) bounded by  $\bar{\mathbf{y}}$  with  $\bar{\mathbf{y}} = \mathbf{W}\bar{\mathbf{u}}$  where  $\mathbf{W} \in \mathbb{R}^{p \times q}$  is the Worst-Case Peak Gain (WCPG) matrix [1] of the system  $\mathcal{H}$ , defined by

$$\mathbf{W} := |\mathbf{D}| + \sum_{k=0}^{\infty} |\mathbf{C}\mathbf{A}^k\mathbf{B}| \quad (2)$$

*Proof:* Let  $\{\mathbf{J}(k)\}_{k \geq 0}$  be the impulse response matrix of the system, i.e.  $\mathbf{J}_{ij}(k)$  is the response on the  $i^{\text{th}}$  output to the Dirac impulse at time  $k = 0$  (i.e.  $\delta(0) = 1$  and  $\delta(k) = 0, \forall k \neq 0$ ) on the  $j^{\text{th}}$  input.

With (1), we have

$$\mathbf{J}(k) = \begin{cases} \mathbf{D} & \text{if } k = 0 \\ \mathbf{C}\mathbf{A}^{k-1}\mathbf{B} & \text{if } k > 0. \end{cases} \quad (3)$$

Since the input  $\{\mathbf{u}(k)\}_{k \geq 0}$  can be seen as a weighted sum of Dirac impulses (shifted in time), and thanks to the linearity and time invariance property of LTI systems [9], we get

$$\mathbf{y}(k) = \sum_{l=0}^k \mathbf{J}(l)\mathbf{u}(k-l) \quad (4)$$

( $\{\mathbf{y}\}_{k \geq 0}$  is the result of the convolution of  $\{\mathbf{J}\}_{k \geq 0}$  by  $\{\mathbf{u}\}_{k \geq 0}$ ).

Then the output is (component-wise) bounded by

$$\mathbf{y}(k) \leq \left( \sum_{l=0}^k |\mathbf{J}(l)| \right) \bar{\mathbf{u}}, \quad \forall k \geq 0. \quad (5)$$

We have equality for the  $i^{\text{th}}$  output if the input is such that  $\mathbf{u}_j(l) = \bar{u}_j \cdot \text{sign}(\mathbf{J}_{ij}(k-l)), \forall 0 \leq l \leq k$ , where  $\text{sign}(x)$  returns  $\pm 1$  or 0 depending on the sign of  $x$ . Finally

$$\forall k \geq 0, \quad \mathbf{y}(k) < \left( \sum_{l=0}^{\infty} |\mathbf{J}(l)| \right) \bar{\mathbf{u}}. \quad (6)$$

■

**Remark 1**  $\mathbf{W}\bar{\mathbf{u}}$  is the supremum of the output  $\{\mathbf{y}\}_{k \geq 0}$ , since it is possible to build a finite input  $\{\mathbf{u}(k)\}_{0 \leq k \leq K}$  to approach it on any given output at any given distance.

**Remark 2** This proposition can be completed when considering intervals for the input, instead of bounds (corresponding to symmetric intervals). In that case, the Worst-Case Peak Gain matrix indicates by how much the radius of the input interval is amplified on the output [8] (although this is not valid for the transient phase, i.e. for the few first steps). However, even in that case,  $\mathbf{W}\bar{\mathbf{u}}$  is a supremum we need to compute.

This proposition can be used to bound outputs, states or intermediate variables in the context of finite precision implementation of algorithms, and more specially in Fixed-Point

arithmetic. In [7], an extension of the state-space has been presented, in order to represent and encompass all the possible algorithms for linear filters (i.e. all the input-to-output data flows based on additions, multiplications by constant and delay, such as state-space, direct forms,  $\rho$ DFIIT [18], etc.), and the same approach was applied.

First, it is used to bound all the variables involved in the algorithm, and then to determine their fixed-point representation (position of the Most Significant Bit and scaling) while preserving by construction from overflow.

Second, it is used to determine the impact on the output of the computational errors. Classical error analysis cannot be used in that context due to the feedback scheme of the computation (Interval Arithmetic or Affine Arithmetic do not provide tight bounds [12]).

Since the filter is linear, the implemented filter  $\mathcal{H}^*$  can be seen as the exact filter  $\mathcal{H}$  where the output is corrupted by the vector of errors  $e(k)$  occurring at each sum of product through a given linear filter  $\mathcal{H}_e$  (see Figure 1).

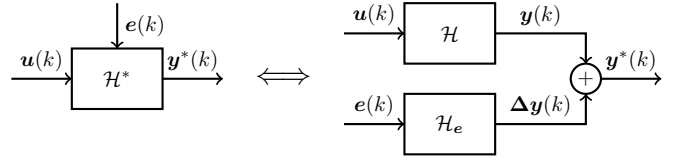


Fig. 1. The implemented filter is equivalent to the exact filter where the output is corrupted by the computational errors passing themselves through a filter.

State-space matrices of  $\mathcal{H}_e$  can be obtained analytically [8] and Proposition 1 can be used to determine the output error  $\Delta \mathbf{y}$  due to finite-precision arithmetic.

For all these reasons, the reliable computation of the Worst-Case Peak Gain matrix is a required step for the accurate error analysis of LTI systems in finite precision.

## II. ALGORITHM FOR WCPG EVALUATION

Given an LTI filter in state-space realization (1) and  $\varepsilon$ , a desired absolute approximation error, we want to determine the Worst-Case Peak Gain matrix  $\mathbf{W}$  of this filter, defined in (2). While computing such an approximation, various errors, such as truncation and summation errors, are made.

Instead of directly computing the infinite sum  $|\mathbf{C}\mathbf{A}^k\mathbf{B}|$  for any  $k \geq 0$ , we will use an approximate eigenvalue decomposition of  $\mathbf{A}$  (i.e.  $\mathbf{A} \approx \mathbf{V}\mathbf{T}\mathbf{V}^{-1}$ ) and compute the FP sum  $|\mathbf{C}\mathbf{V}\mathbf{T}^k\mathbf{V}^{-1}\mathbf{B}|$  for  $0 \leq k \leq N$ .

Our approach to compute the approximation  $\mathbf{S}_N$  of  $\mathbf{W}$  is summarized in algorithm 1 where all the operations ( $\otimes$ ,  $\oplus$ , inv, abs, etc.) are FP multiple precision operations done at various precisions, to be determined but set such that the overall error be less than  $\varepsilon$ :

$$|\mathbf{W} - \mathbf{S}_N| \leq \varepsilon. \quad (7)$$

The overall error analysis decomposes into 6 steps, where each one expresses the impact of a particular approximation (or truncation), and provides the accuracy requirements for the associated operations such that the result is rigorously bounded

---

**Algorithm 1:** Floating-point evaluation of the WCPG.

---

**Input:**  $A \in \mathbb{F}^{n \times n}, B \in \mathbb{F}^{n \times q}, C \in \mathbb{F}^{p \times n}, D \in \mathbb{F}^{p \times q}, \varepsilon > 0$   
**Output:**  $S_N \in \mathbb{F}^{p \times q}$

Step 1: Compute  $N$   
Step 2: Compute  $V$  from an eigendecomposition of  $A$   
 $T \leftarrow \text{inv}(V) \otimes A \otimes V$   
Check that  $\|T\|_2 \leq 1$   
Step 3:  $B' \leftarrow \text{inv}(V) \otimes B$   
 $C' \leftarrow C \otimes V$   
 $S_{-1} \leftarrow |D|, P_{-1} \leftarrow I_n$   
**for**  $k$  **from** 0 **to**  $N$  **do**  
Step 4:  $P_k \leftarrow T \otimes P_{k-1}$   
Step 5:  $L_k \leftarrow C' \otimes P_k \otimes B'$   
Step 6:  $S_k \leftarrow S_{k-1} \oplus \text{abs}(L_k)$   
**end**  
**return**  $S_N$

---

by  $\varepsilon$ . These steps are discussed in details in Sections III and IV:

Step 1: Let  $W_N$  be the truncated sum

$$W_N := \sum_{k=0}^N |CA^k B| + |D|. \quad (8)$$

We compute a truncation order  $N$  of the infinite sum  $W$  such that the truncation error is less than  $\varepsilon_1 > 0$ :

$$|W - W_N| \leq \varepsilon_1. \quad (9)$$

See Section III for more details.

Step 2: Error analysis for computing the powers  $A^k$  of a full matrix  $A$ , when the  $k$  reaches several hundreds, is a significant problem, especially when the norm of  $A$  is larger than 1 and its eigenvalues are close to 1. However, if  $A$  is diagonalizable, i.e. may be represented as  $A = XEX^{-1}$  with  $E \in \mathbb{C}^{n \times n}$  strictly diagonal and  $X \in \mathbb{C}^{n \times n}$  unitary, then powering of  $A$  reduces to powering the diagonal matrix  $E$ , which is more convenient.

Suppose we have an *almost* unitary matrix  $V$  approximating  $X$ . We require this approximation to be just quite accurate so that we are able to discern the different associated eigenvalues and be sure they are less than 1.

We may then consider the matrix  $V$  to be exact and compute an approximation  $T$  to  $V^{-1}AV$  with sufficient accuracy such that the error of computing  $VT^kV^{-1}$  instead of matrix  $A^k$  is less than  $\varepsilon_2 > 0$ :

$$\left| W_N - \sum_{k=0}^N |CVT^kV^{-1}B| \right| \leq \varepsilon_2. \quad (10)$$

See Section IV-A.

Step 3: We compute approximations  $B'$  and  $C'$  of  $V^{-1}B$  and  $CV$ , respectively. We require that the propagated error committed in using  $B'$  instead of  $V^{-1}B$  and  $C'$  instead of  $CV$  be less than  $\varepsilon_3 > 0$ :

$$\left| \sum_{k=0}^N |CVT^kV^{-1}B| - \sum_{k=0}^N |C'T^kB'| \right| \leq \varepsilon_3. \quad (11)$$

See Section IV-B.

Step 4: We compute in  $P_k$  the powers  $T^k$  of  $T$  with a certain accuracy. It is required that the error of computations be less than  $\varepsilon_4 > 0$ :

$$\left| \sum_{k=0}^N |C'T^kB'| - \sum_{k=0}^N |C'P_kB'| \right| \leq \varepsilon_4. \quad (12)$$

See Section IV-C.

Step 5: We compute in  $L_k$  each summand  $C'P_kB'$  with a error small enough such that the overall approximation error induced by this step is less than  $\varepsilon_5 > 0$ :

$$\left| \sum_{k=0}^N |C'P_kB'| - \sum_{k=0}^N |L_k| \right| \leq \varepsilon_5. \quad (13)$$

See Section IV-D.

Step 6: Finally, we sum  $L_k$  in  $S_N$  with enough precision so that the absolute error bound for summation is bounded by  $\varepsilon_6 > 0$ :

$$\left| \sum_{k=0}^N |L_k| - S_N \right| \leq \varepsilon_6. \quad (14)$$

See Section IV-E.

By ensuring that each step verifies its bound  $\varepsilon_i$ , and taking  $\varepsilon_i = \frac{1}{6}\varepsilon$ , we get  $\varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \varepsilon_4 + \varepsilon_5 + \varepsilon_6 \leq \varepsilon$ , hence (7) will be satisfied if inequalities (9) to (14) are.

Our approach hence determines first a truncation order  $N$  and then performs summation up to that truncation error, whilst adjusting accuracy in the different summation steps. A competing approach would be not to start with truncation order determination but to immediately go for summation and to stop when adding more terms does not improve accuracy. However, such an approach would not allow the final error to be bounded in an a priori way. As we shall see, the multiple precision FP summation needs to know a bound on the number of terms to be summed, beforehand.

### III. TRUNCATION ORDER AND TRUNCATION ERROR

The goal of the first step is to determine a lower bound on the truncation order  $N$  of the infinite sum (2) such that its 'tail' is smaller than the given  $\varepsilon_1$ .

In [1] Balakrishnan and Boyd propose "simple" lower and upper bounds on  $N$  but their algorithm has a few drawbacks that make it unusable in applications: they use an iterative scheme to find the exact minimal truncation order  $N_{min}$  and, on each step, compute a power and two Hankel singular values of matrix  $A$ . Additionally, they describe their algorithm in terms of exact arithmetic, i.e. do not propose any error analysis. Finally, although they propose a method for an exact  $N_{min}$  determination, the complexity of error analysis and implementation of their algorithm is too large and even unnecessary for our purpose.

Obviously,  $W_N$  is a lower bound on  $W$  and increases monotonically to  $W$  with increasing  $N$ . Hence the truncation error is

$$|W - W_N| = \sum_{k>N} |CA^k B|. \quad (15)$$

### A. A bound on the truncation error

Many simple bounds on (15) are possible. For instance, if matrix  $\mathbf{A}$  is diagonalizable and the eigendecomposition of  $\mathbf{A}$  is computed

$$\mathbf{A} = \mathbf{X}\mathbf{E}\mathbf{X}^{-1} \quad (16)$$

where  $\mathbf{X} \in \mathbb{C}^{n \times n}$  is the right hand eigenvector matrix, and  $\mathbf{E} \in \mathbb{C}^{n \times n}$  is a diagonal matrix holding the eigenvalues  $\lambda_i$ , the terms  $\mathbf{C}\mathbf{A}^k\mathbf{B}$  can be written

$$\mathbf{C}\mathbf{A}^k\mathbf{B} = \mathbf{\Phi}\mathbf{E}^k\mathbf{\Psi} \quad (17)$$

$$= \sum_{l=1}^n \mathbf{R}_l \lambda_l^k \quad (18)$$

where  $\mathbf{\Phi} \in \mathbb{C}^{p \times n}$ ,  $\mathbf{\Psi} \in \mathbb{C}^{n \times q}$  and  $\mathbf{R}_i \in \mathbb{C}^{p \times q}$  are defined by

$$\mathbf{\Phi} := \mathbf{C}\mathbf{X} \quad (19)$$

$$\mathbf{\Psi} := \mathbf{X}^{-1}\mathbf{B} \quad (20)$$

$$(\mathbf{R}_l)_{ij} := \mathbf{\Phi}_{il}\mathbf{\Psi}_{lj}. \quad (21)$$

In this setting, we obtain

$$|\mathbf{W} - \mathbf{W}_N| = \sum_{k>N} \sum_{l=1}^n |\mathbf{R}_l \lambda_l^k|. \quad (22)$$

As required by Proposition 1, all eigenvalues  $\lambda_l$  of matrix  $\mathbf{A}$  must be strictly smaller than one in magnitude. We may therefore notice that the outer sum is in geometric progression with a common ratio  $|\lambda_l| < 1$ . So the following bound is possible (we remind the reader that inequalities and absolute values are considered to be element by element):

$$\begin{aligned} |\mathbf{W} - \mathbf{W}_N| &\leq \sum_{k=N+1}^{\infty} \sum_{l=1}^n |\mathbf{R}_l| |\lambda_l^k| \\ &\leq \sum_{l=1}^n |\mathbf{R}_l| \frac{|\lambda_l^{N+1}|}{1 - |\lambda_l|} \\ &= \rho(\mathbf{A})^{N+1} \sum_{l=1}^n \frac{|\mathbf{R}_l|}{1 - |\lambda_l|} \left( \frac{|\lambda_l|}{\rho(\mathbf{A})} \right)^{N+1}. \end{aligned} \quad (23)$$

Since  $\frac{|\lambda_l|}{\rho(\mathbf{A})} \leq 1$  holds for all terms, we may leave out the powers:

$$|\mathbf{W} - \mathbf{W}_N| \leq \rho(\mathbf{A})^{N+1} \sum_{l=1}^n \frac{|\mathbf{R}_l|}{1 - |\lambda_l|} \frac{|\lambda_l|}{\rho(\mathbf{A})}. \quad (24)$$

Notate

$$\mathbf{M} := \sum_{l=1}^n \frac{|\mathbf{R}_l|}{1 - |\lambda_l|} \frac{|\lambda_l|}{\rho(\mathbf{A})} \in \mathbb{R}^{p \times q}. \quad (25)$$

The tail of the infinite sum left after truncation is hence bounded by

$$|\mathbf{W} - \mathbf{W}_N| \leq \rho(\mathbf{A})^{N+1} \mathbf{M}. \quad (26)$$

**Remark 3** Another tighter bound is possible

$$|\mathbf{W} - \mathbf{W}_N| \leq \rho(\mathbf{A})^{N+1-K} \sum_{l=1}^n \frac{|\mathbf{R}_l|}{1 - |\lambda_l|} \left( \frac{|\lambda_l|}{\rho(\mathbf{A})} \right)^K, \quad \forall N > K. \quad (27)$$

### B. Deducing a lower bound on the truncation order

In order to get (26) bounded by  $\varepsilon_1$ , it is required that

$$\rho(\mathbf{A})^{N+1} \mathbf{M} \leq \varepsilon_1.$$

Solving this inequality for  $N$  leads us to the following bound:

$$N \geq \left\lceil \frac{\log \frac{\varepsilon_1}{m}}{\log \rho(\mathbf{A})} \right\rceil \quad (28)$$

where  $m$  is defined as  $m := \min_{i,j} |\mathbf{M}_{i,j}|$ .

However we cannot compute exact values for all quantities occurring in (27) when using finite-precision arithmetic. We only have approximations for them. Thus, in order to reliably determine a lower bound on  $N$ , we must compute lower bounds on  $m$  and  $\rho(\mathbf{A})$ , from which we can deduce an upper bound on  $\log \frac{\varepsilon_1}{m}$  and a lower bound on  $\log \rho(\mathbf{A})$  to eventually obtain a lower bound on  $N$ .

### C. A rigorous algorithm to determine truncation order

Due to the implementation of eq. (16) to (21) with the finite-precision arithmetic, only approximations on  $\lambda, \mathbf{X}, \mathbf{\Phi}, \mathbf{\Psi}, \mathbf{R}_i$  can be obtained. There exist many FP libraries, such as LAPACK<sup>1</sup>, providing functions for an eigendecomposition as needed for (16) and to solve linear systems of equations (20), but there is a drawback: they usually deliver good and fast approximations to the solution of a given numerical problem but there is neither verification nor guarantee about the accuracy of that approximation.

For these reasons we propose to combine LAPACK FP arithmetic with Interval Arithmetic [3] enhanced with the Theory of Verified Inclusions [14], [15], [16], [17] in order to obtain trusted intervals on the eigensystem and, eventually, a rigorous bound on  $N$ .

In Interval Arithmetic real numbers are represented as sets of reals with addition, subtraction, multiplication and division defined [3]. The Theory of Verified Inclusions is a set of algorithms computing guaranteed bounds on solutions of various numerical problems, developed by S. Rump [14]. The verification process is performed by means of checking an interval fixed point and yields to a trusted interval for the solution, i.e. it is verified that the result interval contains an exact solution of given numerical problem.

It permits us to quickly obtain trusted error bounds on the truncation order without significant impact on algorithm performance, since this computation is done only once. In addition, if the spectral radius of  $\mathbf{A}$  cannot be shown less than 1, we can stop the algorithm.

Using the ideas proposed by Rump in [17], we obtain trusted intervals for the eigensystem with the following steps:

- 1) Using the LAPACK eigensolver, we compute FP approximations  $\mathbf{V}$  for the eigenvectors  $\mathbf{X}$  and  $\alpha$  for the eigenvalues  $\lambda$ , along with error estimates  $\varepsilon_X$  and  $\varepsilon_\lambda$ . These error estimates are such that  $|\lambda - \alpha| \leq \varepsilon_\lambda$  and  $|\mathbf{X} - \mathbf{V}| \leq \varepsilon_X$  should be not far from the truth.
- 2) We construct, verify and possibly adjust intervals for  $[\lambda] = [\alpha - \varepsilon_\lambda, \alpha + \varepsilon_\lambda]$  and  $[\mathbf{X}] = [\mathbf{V} - \varepsilon_X, \mathbf{V} + \varepsilon_X]$

<sup>1</sup><http://www.netlib.org/lapack/>

such that for all vectors  $\lambda' \in [\lambda]$  there exists a matrix  $X' \in [X]$  satisfying  $AX' = X' \text{diag}(\lambda')$  and such that for all matrices  $X' \in [X]$  there exists a vector  $\lambda' \in [\lambda]$  satisfying  $AX' = X' \text{diag}(\lambda')$ .

In this process, first intervals for the eigensystem are constructed from the error estimates  $\varepsilon_\alpha$  and  $\varepsilon_V$  as radii and the approximate solutions  $V$  and  $\alpha$  as mid-points. Further, these intervals are verified with inclusion algorithms [17]. If the verification does not succeed, the intervals are extended by some small factor and process is repeated until it succeeds or until there exists an eigenvalue interval which contains 1.

For the solution of the linear system of equations (LSE) appearing in (20), the algorithm for interval verification is based on [15] and consists in two steps:

- 1) Using LAPACK, compute a FP approximation  $\Omega$  on the solution of  $V\Psi = B$  along with an error estimate  $\varepsilon_\Psi$  such that  $|\Psi - \Omega| \leq \varepsilon_\Psi$  should be true.
- 2) Construct, verify and adjust intervals  $[\Psi] = [\Omega - \varepsilon_\Psi, \Omega + \varepsilon_\Psi]$  such that for all matrices  $X' \in [X]$  there exists  $\Psi' \in [\Psi]$  such that  $X'\Psi' = B$  holds. The intervals for verification are constructed in the same way as for the eigensystem solution. We require the existence of the exact solution of the linear system of equations not for the system  $V\Psi = B$  but for  $[X]\Psi = B$ , i.e.  $[\Psi]$  must contain the exact solution for each element of the already verified interval  $[X]$ .

Finally, the intervals for (19), (21), (25) and (27) are computed with Interval Arithmetic. Our complete algorithm to determine a reliable lower bound on  $N$  is given with Algorithm 2.

---

**Algorithm 2:** Lower bound of truncation order

---

**Input:**  $A \in \mathbb{F}^{n \times n}, B \in \mathbb{F}^{n \times q}, C \in \mathbb{F}^{p \times n}, \varepsilon_1 > 0$

**Output:**  $N \in \mathbb{N}$

- 1  $\alpha, V, \varepsilon_\alpha, \varepsilon_V \leftarrow$  LAPACK eigendecomposition for  $A$ ;
  - 2  $\Omega, \varepsilon_\Psi \leftarrow$  LAPACK solver for  $V\Psi = B$ ;
  - 3  $[\lambda], [X] \leftarrow$  Eigensystem verification algorithm;
  - 4  $[\Psi] \leftarrow$  LSE solution verification algorithm;
  - 5  $[\Phi] \leftarrow C[X]$ ;
  - 6  $[R_l]_{i,j} \leftarrow [\Phi_{i,l}][\Psi_{l,j}]$ ;
  - 7  $[\rho] \leftarrow \max_i |\lambda_i|$ ;
  - 8  $[M] \leftarrow \sum_{i=1}^n \frac{|[R_i]|}{1 - |\lambda_i|} \frac{|\lambda_i|}{[\rho]}$ ;
  - 9  $[m] \leftarrow \min_{i,j} |[M]_{i,j}|$ ;
  - 10  $N \leftarrow \sup \left( \left[ \frac{\log \frac{\varepsilon_1}{[m]}}{\log [\rho]} \right] \right)$ ;
  - 11 **return**  $N$
- 

#### IV. SUMMATION

Once the truncation order determined, we need to provide a summation scheme reliable in FP arithmetic, i.e. such that the error of computations is bounded by an a priori given value. To do so we propose to perform all operations in multiple precision arithmetic whilst adapting precision dynamically where needed. For this purpose several multiple precision algorithms are developed:

- $\text{multiplyAndAdd}(A, B, C, \delta)$ : given matrices  $A \in \mathbb{C}^{p \times n}, B \in \mathbb{C}^{n \times q}, C \in \mathbb{C}^{p \times q}$ , computes a matrix  $D \in \mathbb{C}^{p \times n}$  such that  $D = A \cdot B + C + \Delta$ , where the error matrix  $\Delta$  is bounded by  $|\Delta| < \delta$ , for a certain scalar absolute error bound  $\delta$ , given in argument to the algorithm. The algorithm can be reused to compute plain matrix products  $A \cdot B$ , by setting  $C$  to the zero matrix; in this case we notate  $A \otimes B$  for the output of  $\text{multiplyAndAdd}$ .
- $\text{sumAbs}(A, B, \delta)$ : given  $A \in \mathbb{R}^{p \times n}, B \in \mathbb{C}^{p \times n}$ , computes a matrix  $C \in \mathbb{R}^{p \times n}$  such that  $C = A + |B| + \Delta$ , where the error matrix  $\Delta$  is bounded by  $|\Delta| < \delta$ , for a certain scalar absolute error bound  $\delta$ , given in argument to the algorithm. With a slight notational abuse, we shall also notate  $A \oplus \text{abs}(B)$  for  $\text{sumAbs}$ , even though we condense both the error in computing the absolute value of a complex matrix and the addition error into one, using  $\text{sumAbs}$ .
- $\text{inv}(V, \delta)$ : given a complex square matrix  $V \in \mathbb{C}^{n \times n}$ , computes a matrix  $U \in \mathbb{C}^{n \times n}$  such that  $U = V^{-1} + \Delta$ , where the error matrix  $\Delta$  is bounded by  $|\Delta| < \delta$ , for a certain scalar absolute error bound  $\delta$ , given in argument to the algorithm. The algorithm we use is based on Newton-Raphson matrix iteration, requires a seed matrix inverse in argument and works only under certain conditions, which we can easily verify in our case. See Section (V).

These computation kernels adapt the precision of their intermediate computations where needed. The algorithms we use for these basic bricks will be discussed in Section (V).

##### A. Step 2: using the Eigendecomposition

1) *Propagation of  $\Delta_2$* : As seen, in each step of the summation, a matrix power,  $A^k$ , must be computed. In [6] Higham devotes an entire chapter to error analysis of matrix powers but this theory is in most cases inapplicable for state matrices  $A$  of linear filters, as the requirement  $\rho(|A|) < 1$  does not necessarily hold here. Therefore, despite taking  $A$  to just a finite power  $k$ , the sequence of computed matrices may explode in norm since  $k$  may take an order of several hundreds or thousands. Thus, even extending the precision is not a solution, as an enormous number of bits would be required.

However, the state matrices  $A$  usually have a good structure. Suppose  $A$  is diagonalizable, i.e. there exists an unitary matrix  $X \in \mathbb{C}^{n \times n}$  and diagonal  $E \in \mathbb{C}^{n \times n}$  such that  $A = XEX^{-1}$ . Then  $A^k = XE^kX^{-1}$ . A good choice of  $X$  and  $E$  are the eigenvector and eigenvalue matrices obtained with eigendecomposition (16). However, with LAPACK we can compute only approximations on them and we cannot control their accuracy. Therefore, we propose following method to *almost* diagonalize matrix  $A$ . The method does not make any assumptions on matrix  $V$  except for it being an *almost* unitary. Therefore, for simplicity of further reasoning we treat  $V$  as an exact matrix.

Let  $K$  be some matrix,  $K \in \mathbb{C}^{n \times m}$ , then its Frobenius

norm  $\|\mathbf{K}\|_F$  is defined by:

$$\|\mathbf{K}\|_F := \sqrt{\sum_{i=1}^n \sum_{j=1}^m |\mathbf{K}_{ij}|^2}. \quad (29)$$

The Frobenius norm is sub-multiplicative and its following properties are used in the discussions below:

$$|\mathbf{K}_{ij}| \leq \|\mathbf{K}\|_F \quad \forall i, j \quad (30)$$

$$\|\mathbf{K}\|_2 \leq \|\mathbf{K}\|_F \leq \sqrt{\min(m, n)} \|\mathbf{K}\|_2, \quad (31)$$

where  $\|\mathbf{K}\|_2$  is the spectral-norm, i.e. equal to the largest singular value of  $\mathbf{K}$ .

Moreover, if  $\mathbf{K}$  is a square  $n \times n$  matrix such that  $\|\mathbf{K}\|_2 \leq 1$ , then for all  $k$ ,  $\|\mathbf{K}^k\|_2 \leq 1$  and

$$\|\mathbf{K}^k\|_F \leq \sqrt{n}. \quad (32)$$

Using our multiprecision algorithms for matrix inverse and multiplication we may compute a complex  $n \times n$  matrix  $\mathbf{T}$ :

$$\mathbf{T} := \mathbf{V}^{-1} \mathbf{A} \mathbf{V} - \Delta_2, \quad (33)$$

where  $\mathbf{V} \in \mathbb{C}^{n \times n}$  is an approximation on  $\mathbf{X}$ , i.e. an *almost* unitary matrix,  $\Delta_2 \in \mathbb{C}^{n \times n}$  is a matrix representing the element-by-element errors due to the two matrix multiplications and the inversion of matrix  $\mathbf{V}$ .

Although the matrix  $\mathbf{E}$  is strictly diagonal, here  $\mathbf{V}$  is not exactly unitary and consequently  $\mathbf{T}$  is a full matrix. However it has prevailing elements on the main diagonal. Thus  $\mathbf{T}$  is an approximation on  $\mathbf{E}$ .

We require for matrix  $\mathbf{T}$  to satisfy  $\|\mathbf{T}\|_2 \leq 1$ . This condition is stronger than  $\rho(\mathbf{A}) < 1$ , and Section IV-A2 provides a way to test this condition.

Notate  $\Xi_k := (\mathbf{T} + \Delta_2)^k - \mathbf{T}^k$ . Hence  $\Xi_k \in \mathbb{C}^{n \times n}$  represents an error matrix which captures the propagation of error  $\Delta_2$  when powering  $\mathbf{T}$ . Since

$$\mathbf{A}^k = \mathbf{V}(\mathbf{T} + \Delta_2)^k \mathbf{V}^{-1}, \quad (34)$$

therefore

$$\mathbf{C} \mathbf{A}^k \mathbf{B} = \mathbf{C} \mathbf{V} \mathbf{T}^k \mathbf{V}^{-1} \mathbf{B} + \mathbf{C} \mathbf{V} \Xi_k \mathbf{V}^{-1} \mathbf{B}. \quad (35)$$

Thus the error of computing  $\mathbf{V} \mathbf{T}^k \mathbf{V}^{-1}$  instead of  $\mathbf{A}^k$  in (8) is bounded by

$$\left| \sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B}| - \sum_{k=0}^N |\mathbf{C} \mathbf{V} \mathbf{T}^k \mathbf{V}^{-1} \mathbf{B}| \right| \leq \quad (36)$$

$$\sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B} - \mathbf{C} \mathbf{V} \mathbf{T}^k \mathbf{V}^{-1} \mathbf{B}| \leq \sum_{k=0}^N |\mathbf{C} \mathbf{V} \Xi_k \mathbf{V}^{-1} \mathbf{B}| \quad (37)$$

Here and further on each step of the algorithm we use rather inequalities with left side in form (36) rather than (35), i.e. we will instantly use the triangulation property  $||a| - |b|| \leq |a - b| \quad \forall a, b$  applied element-by-element to matrices.

In order to determine the accuracy of the computations on this step such that (36) is bounded by  $\varepsilon_2$ , we need to perform detailed analysis of  $\Xi_k$ , with spectral-norm.

$$\begin{aligned} \|\Xi_k\|_2 &= \|(\mathbf{T} + \Delta_2)^k - \mathbf{T}^k\|_2 \\ &= \|\mathbf{T}(\mathbf{T} + \Delta_2)^{k-1} + \Delta_2(\mathbf{T} + \Delta_2)^{k-1} - \mathbf{T} \mathbf{T}^{k-1}\|_2 \\ &= \|\mathbf{T}((\mathbf{T} + \Delta_2)^{k-1} - \mathbf{T}^{k-1}) + \Delta_2(\mathbf{T} + \Delta_2)^{k-1}\|_2 \\ &= \|\mathbf{T} \Xi_{k-1} + \Delta_2(\Xi_{k-1} + \mathbf{T}^{k-1})\|_2 \\ &\leq \|\mathbf{T}\|_2 \|\Xi_{k-1}\|_2 + \|\Delta_2\|_2 (\|\Xi_{k-1}\|_2 + \|\mathbf{T}\|_2^{k-1}) \\ &\leq \|\Xi_{k-1}\|_2 + \|\Delta_2\|_2 (\|\Xi_{k-1}\|_2 + 1) \end{aligned} \quad (38)$$

If  $\|\Xi_{k-1}\|_2 \leq 1$ , which must hold in our case since  $\Xi_k$  represent an error-matrix, then

$$\|\Xi_k\|_2 \leq \|\Xi_{k-1}\|_2 + 2 \|\Delta_2\|_2 \quad (39)$$

As  $\|\Xi_1\|_2 = \|\Delta_2\|_2$  we get the desired bound capturing the propagation of  $\Delta_2$

$$\|\Xi_k\|_2 \leq 2(k+1) \|\Delta_2\|_2, \quad (40)$$

or, with Frobenius norm,

$$\|\Xi_k\|_F \leq 2\sqrt{n}(k+1) \|\Delta_2\|_F. \quad (41)$$

Substituting this bound to (36) and folding the sum, we obtain

$$\sum_{i=0}^N |\mathbf{C} \mathbf{V} \Xi_k \mathbf{V}^{-1} \mathbf{B}| \leq \beta \|\Delta_2\|_F \|\mathbf{C} \mathbf{V}\|_F \|\mathbf{V}^{-1} \mathbf{B}\|_F, \quad (42)$$

with  $\beta = \sqrt{n}(N+1)(N+2)$ . Thus, we get a bound on the error of approximation of  $\mathbf{A}$  by  $\mathbf{V} \mathbf{T} \mathbf{V}^{-1}$ . Since we require it to be less than  $\varepsilon_2$  we obtain a condition for the error on the inversion and two matrix multiplications:

$$\|\Delta_2\|_F \leq \frac{1}{\beta} \frac{\varepsilon_2}{\|\mathbf{C} \mathbf{V}\|_F \|\mathbf{V}^{-1} \mathbf{B}\|_F}. \quad (43)$$

Using this bound we can deduce the desired accuracy of our multiprecision algorithms for complex matrix multiplication and inverse as a function of  $\varepsilon_2$ .

2) *Checking  $\|\mathbf{T}\|_2 \leq 1$ :* Since  $\|\mathbf{T}\|_2^2 = \rho(\mathbf{T}^* \mathbf{T})$ , we study the eigenvalues of  $\mathbf{T}^* \mathbf{T}$ . According to Gershgorin's circle theorem [5], each eigenvalue  $\mu_i$  of  $\mathbf{T}^* \mathbf{T}$  is in the disk centered in  $(\mathbf{T}^* \mathbf{T})_{ii}$  with radius  $\sum_{j \neq i} |(\mathbf{T}^* \mathbf{T})_{ij}|$ .

Let us decompose  $\mathbf{T}$  into  $\mathbf{T} = \mathbf{F} + \mathbf{G}$ , where  $\mathbf{F}$  is diagonal and  $\mathbf{G}$  contains all the other terms ( $\mathbf{F}$  contains the approximate eigenvalues,  $\mathbf{G}$  contains small terms and is zero on its diagonal). Denote  $\mathbf{Y} := \mathbf{T}^* \mathbf{T} - \mathbf{F}^* \mathbf{F} = \mathbf{F}^* \mathbf{G} + \mathbf{G}^* \mathbf{F} + \mathbf{G}^* \mathbf{G}$ . Then

$$\begin{aligned} \sum_{j \neq i} |(\mathbf{T}^* \mathbf{T})_{ij}| &= \sum_{j \neq i} |\mathbf{Y}_{ij}| \\ &\leq (n-1) \|\mathbf{Y}\|_F \\ &\leq (n-1) \left( 2 \|\mathbf{F}\|_F \|\mathbf{G}\|_F + \|\mathbf{G}\|_F^2 \right) \\ &\leq (n-1) (2\sqrt{n} + \|\mathbf{G}\|_F) \|\mathbf{G}\|_F. \end{aligned} \quad (44)$$

Each eigenvalue of  $\mathbf{T}^* \mathbf{T}$  is in the disk centered in  $(\mathbf{F}^* \mathbf{F})_{ii} + (\mathbf{Y})_{ii}$  with radius  $\gamma$ , where  $\gamma$  is equal to  $(n -$

1)  $(2\sqrt{n} + \|\mathbf{G}\|_F) \|\mathbf{G}\|_F$  computed in a rounding mode that makes the result become an upper bound (round-up).

As  $\mathbf{G}$  is zero on its diagonal, the diagonal elements  $(\mathbf{Y})_{ii}$  of  $\mathbf{Y}$  are equal to the diagonal elements  $(\mathbf{G}^* \mathbf{G})_{ii}$  of  $\mathbf{G}^* \mathbf{G}$ . They can hence be bounded as follows:

$$\begin{aligned} |(\mathbf{Y})_{ii}| &= |(\mathbf{G}^* \mathbf{G})_{ii}| \\ &\leq \|\mathbf{G}\|_F^2. \end{aligned} \quad (45)$$

Then, it is easy to see that the Gershgorin circles enclosing the eigenvalues of  $\mathbf{F}^* \mathbf{F}$  can be increased, meaning that if  $(\mathbf{F}^* \mathbf{F})_{ii}$  is such that

$$\forall i, \quad |(\mathbf{F}^* \mathbf{F})_{ii}| \leq 1 - \|\mathbf{G}\|_F^2 - \gamma, \quad (46)$$

it holds that  $\rho(\mathbf{T}^* \mathbf{T}) \leq 1$  and  $\|\mathbf{T}\|_2 \leq 1$ .

This condition can be tested by using FP arithmetic with directed rounding modes (round-up for instance).

After computing  $\mathbf{T}$  out of  $\mathbf{V}$  and  $\mathbf{A}$  according to (32), the condition on  $\mathbf{T}$  should be tested in order to determine if  $\|\mathbf{T}\|_2 \leq 1$ . This test failing means that  $\mathbf{V}$  is not a sufficient approximate of  $\mathbf{X}$  or that the error  $\Delta_2$  done computing (32) is too large, i.e. the accuracy of our multiprecision algorithm for complex matrix multiplication and inverse should be increased. The test is required for rigor only. We do perform the test in the implementation of our WCPG method, and, on the examples we tested, never saw it fail.

### B. Step 3: computing $\mathbf{CV}$ and $\mathbf{V}^{-1}\mathbf{B}$

We compute approximations on matrices  $\mathbf{CV}$  and  $\mathbf{V}^{-1}\mathbf{B}$  with a certain precision and need to determine the required accuracy of these multiplications such that the impact of these approximations is less than  $\varepsilon_3$ .

Notate  $\mathbf{C}' := \mathbf{CV} + \Delta_{3C}$  and  $\mathbf{B}' := \mathbf{V}^{-1}\mathbf{B} + \Delta_{3B}$ , where  $\Delta_{3C} \in \mathbb{C}^{p \times n}$  and  $\Delta_{3B} \in \mathbb{C}^{n \times q}$  are error-matrices containing the errors of the two matrix multiplications and the inversion.

Using Frobenius norm, we can bound the error in the approximation of  $\mathbf{CV}$  and  $\mathbf{V}^{-1}\mathbf{B}$  by  $\mathbf{C}'$  and  $\mathbf{B}'$  as follows:

$$\begin{aligned} \sum_{k=0}^N |\mathbf{CVT}^k \mathbf{V}^{-1}\mathbf{B} - \mathbf{C}'\mathbf{T}^k \mathbf{B}'| &\leq \quad (47) \\ \sum_{k=0}^N \|\Delta_{3C} \mathbf{T}^k \mathbf{B}' + \mathbf{C}' \mathbf{T}^k \Delta_{3B} + \Delta_{3C} \mathbf{T}^k \Delta_{3B}\|_F. \end{aligned}$$

Since  $\|\mathbf{T}\|_2 < 1$  holds we have (thanks to eq (31))

$$\|\Delta_{3C} \mathbf{T}^k \mathbf{B}' + \mathbf{C}' \mathbf{T}^k \Delta_{3B} + \Delta_{3C} \mathbf{T}^k \Delta_{3B}\|_F \leq \quad (48)$$

$$\sqrt{n} (\|\Delta_{3C}\|_F (\|\mathbf{B}'\|_F + \|\Delta_{3B}\|_F) + \|\mathbf{C}'\|_F \|\Delta_{3B}\|_F).$$

This bound represents the impact of our approximations for each  $k = 0 \dots N$ . If it is less than  $\frac{1}{N+1} \cdot \varepsilon_3$ , then the overall error is less than  $\varepsilon_3$ :

$$\begin{aligned} \|\Delta_{3C}\|_F (\|\mathbf{B}'\|_F + \|\Delta_{3B}\|_F) + \|\mathbf{C}'\|_F \|\Delta_{3B}\|_F &\leq \\ \frac{1}{\sqrt{n}(N+1)} \varepsilon_3. \end{aligned} \quad (49)$$

And then, it is evident that the two following conditions imply (48):

$$\|\Delta_{3C}\|_F \leq \frac{1}{3\sqrt{n}} \cdot \frac{1}{N+1} \frac{\varepsilon_3}{\|\mathbf{C}'\|_F} \quad (50)$$

$$\|\Delta_{3B}\|_F \leq \frac{1}{3\sqrt{n}} \cdot \frac{1}{N+1} \frac{\varepsilon_3}{\|\mathbf{B}'\|_F}. \quad (51)$$

Therefore, using bounds on  $\|\Delta_{3C}\|_F$  and  $\|\Delta_{3B}\|_F$ , we can deduce the required accuracy of our multiprecision matrix multiplication and inversion according to  $\varepsilon_3$ .

### C. Step 4: powering $\mathbf{T}$

Given a square complex matrix  $\mathbf{T}$  with prevailing main diagonal we need to compute its  $k^{\text{th}}$  power. Notate

$$\mathbf{P}_k := \mathbf{T}^k - \mathbf{\Pi}_k, \quad (52)$$

where  $\mathbf{\Pi}_k \in \mathbb{C}^{n \times n}$  represents element-by-element the error on the matrix powers, including error propagation from the first to the last power. Using the same simplification as in (35) and (36) we get the error of computing the approximations  $\mathbf{P}_k$  rather than the exact powers bounded by

$$\sum_{k=0}^N |\mathbf{C}' \mathbf{T}^k \mathbf{B}' - \mathbf{C}' \mathbf{P}_k \mathbf{B}'| \leq \sum_{k=0}^N |\mathbf{C}' \mathbf{\Pi}_k \mathbf{B}'|. \quad (53)$$

Thus a bound on  $|\mathbf{\Pi}_k|$  is required.

Since we need all the powers of  $\mathbf{T}$  from 1 to  $N$ , we use an iterative scheme to compute them. It is then evident, that we may write the recurrence

$$\mathbf{P}_k = \mathbf{T} \mathbf{P}_{k-1} + \mathbf{\Gamma}_k, \quad (54)$$

where  $\mathbf{\Gamma}_k \in \mathbb{C}^{n \times n}$  is the error matrix representing the error of the matrix multiplication at step  $k$ .

With  $\mathbf{P}_0 = \mathbf{I}$ ,  $\mathbf{P}_1 = \mathbf{T}$  and using the recurrence (53) we obtain

$$\mathbf{P}_k = \mathbf{T}^k + \sum_{l=2}^k \mathbf{T}^{k-l} \mathbf{\Gamma}_l. \quad (55)$$

Therefore,  $\mathbf{\Pi}_k = - \sum_{l=2}^k \mathbf{T}^{k-l} \mathbf{\Gamma}_l$ . Using the condition  $\|\mathbf{T}\|_2 < 1$  and properties of the Frobenius norm we get

$$\|\mathbf{\Pi}_k\|_F \leq \left\| \sum_{l=2}^k \mathbf{T}^{k-l} \mathbf{\Gamma}_l \right\|_F \leq \sqrt{n} \sum_{l=2}^k \|\mathbf{\Gamma}_l\|_F. \quad (56)$$

Therefore the impact of approximation of the matrix powers is bounded by

$$\begin{aligned} \sum_{k=0}^N |\mathbf{C}' \mathbf{\Pi}_k \mathbf{B}'| &\leq \sqrt{n} \sum_{k=0}^N \sum_{l=2}^k \|\mathbf{C}'\|_F \|\mathbf{\Gamma}_l\|_F \|\mathbf{B}'\|_F \\ &\leq \sqrt{n} (N+1) \sum_{l=2}^N \|\mathbf{C}'\|_F \|\mathbf{\Gamma}_l\|_F \|\mathbf{B}'\|_F. \end{aligned} \quad (57)$$



Obviously, if the error of matrix multiplication  $\Gamma_l$  satisfies

$$\|\Gamma_l\|_F \leq \frac{1}{\sqrt{n}} \cdot \frac{1}{N-1} \cdot \frac{1}{N+1} \cdot \frac{\varepsilon_4}{\|\mathbf{C}'\|_F \|\mathbf{B}'\|_F}. \quad (58)$$

for  $l = 2 \dots N$ , then we have (56) to be less than  $\varepsilon_4$ . Hence using (57) we may deduce the required accuracy of matrix multiplications on each step in dependency of  $\varepsilon_4$ .

#### D. Step 5: computing $L_k$

Once the matrices  $\mathbf{C}'$ ,  $\mathbf{B}'$  and  $\mathbf{P}_k$  are pre-computed and the error of their computation is bounded, we must evaluate their product. Let  $L_k$  be the approximate product of these three matrices at step  $k$ :

$$L_k := \mathbf{C}' \mathbf{P}_k \mathbf{B}' + \Upsilon_k, \quad (59)$$

where  $\Upsilon_k \in \mathbb{C}^{p \times q}$  is the matrix of element-by-element errors for the two matrix multiplications.

Then it may be shown, that the error of computations induced by this step is bounded by

$$\sum_{k=0}^N |\mathbf{C}' \mathbf{P}_k \mathbf{B}' - L_k| \leq \sum_{k=0}^N |\Upsilon_k|. \quad (60)$$

If we want the overall error of approximation on this step to be less than  $\varepsilon_5$  then we can choose  $\forall k = 0 \dots N$ :

$$|\Upsilon_k| \leq \frac{1}{N+1} \cdot \varepsilon_5. \quad (61)$$

As  $\Upsilon_k$  represents the matrix of the error due to two multiplications, with bound (60) we may deduce the required accuracy of each of those multiplications on every iteration of summation algorithm in dependency with  $\varepsilon_5$ .

#### E. Step 6: final summation

Finally the absolute value of the  $L_k$  must be taken and the result accumulated in the sum. We remind the reader that if all previous computations were exact, the matrix  $L_k$  would be a real matrix and the absolute-value-operation would have been an exact sign manipulation. However, as the computations were in finite-precision arithmetic,  $L_k$  is complex with a small imaginary part, which is naturally caused by the errors of computations and must not be neglected. Therefore the element-by-element absolute value of complex matrix must be computed.

Since we perform  $N+1$  accumulations of absolute values in the result sum  $\mathbf{S}_N$ , it is evident that bounding the error of each such computation by  $\frac{1}{N+1} \varepsilon_6$  is sufficient.

Therefore, using this bound for each invocation of our basic brick algorithm sumAbs we guarantee bound (14).

## V. BASIC BRICKS

In Section IV, we postulated the existence of three basic FP algorithms, multiplyAndAdd, sumAbs and inv, computing, respectively, the product-sum, the sum in absolute value and the inverse of matrices. Each of these operators was required to satisfy an absolute error bound  $|\Delta| < \delta$  to be ensured by the

matrix of errors  $\Delta$  with respect to scalar  $\delta$ , given in argument to the algorithm.

Ensuring such an *absolute* error bound is not possible in general when fixed-precision FP arithmetic is used. Any such algorithm, when returning its result, must round into that fixed-precision FP format. Hence, when the output grows sufficiently large, the unit-in-the-last-place of that format and hence that final rounding error in fixed-precision FP arithmetic will grow larger than a set absolute error bound.

In multiple precision FP arithmetic, such as offered by software packages like MPFR<sup>2</sup> [4], it is however possible to have algorithms determine themselves the output precision of the FP variables they return their results in. Hence an absolute error bound as the one we require can be guaranteed. In contrast to classical FP arithmetic, such as Higham analyzes, there is no longer any clear, overall *computing precision*, though. Variables just bear the precision that had been determined for them by the previous compute step.

This preliminary clarification made, description of our three basic bricks multiplyAndAdd, sumAbs and inv is easy.

For sumAbs( $\mathbf{A}, \mathbf{B}, \delta$ ) =  $\mathbf{A} + |\mathbf{B}| + \Delta$ , we can reason element by element. We need to approximate  $A_{ij} + \sqrt{\Re B_{ij}^2 + \Im B_{ij}^2}$  with absolute error no larger than  $\delta$ , where  $\Re z$  and  $\Im z$  are the real and imaginary parts of the complex  $z$ . This can be ensured by considering the FP exponents of each of  $A_{ij}$ ,  $\Re B_{ij}$  and  $\Im B_{ij}$  with respect to the FP exponent of  $\delta$ .

For multiplyAndAdd( $\mathbf{A}, \mathbf{B}, \mathbf{C}, \delta$ ) =  $\mathbf{A} \cdot \mathbf{B} + \mathbf{C} + \Delta$ , we can reason in terms of scalar products between  $\mathbf{A}$  and  $\mathbf{B}$ . The scalar products boil down to summation of products which, in turn, can be done exactly, as we can determine the precision of the  $A_{ik}$  and  $B_{kj}$ . As a matter of course the very same summation can capture the matrix elements  $C_{ij}$ . Finally, multiple precision FP summation with an absolute error bound can be performed with a modified, software-simulated Kulisch accumulator [10], which does not need to be exact but bear just enough precision to satisfy the absolute accuracy bound  $\delta$ .

Finally, once the multiplyAndAdd operator is available, it is possible to implement the matrix inversion algorithm inv using a Newton-Raphson-like iteration [13]:

$$\begin{aligned} \mathbf{U}_0 &\leftarrow \text{some seed inverse matrix for } \mathbf{V}^{-1} \\ \mathbf{R}_k &\leftarrow \mathbf{V} \mathbf{U}_k - \mathbf{I} \\ \mathbf{U}_{k+1} &\leftarrow \mathbf{U}_k - \mathbf{U}_k \mathbf{R} \end{aligned} \quad (62)$$

where the iterated matrices  $\mathbf{U}_k$  converge to  $\mathbf{V}^{-1}$  provided the multiplyAndAdd operations computing  $\mathbf{R}_k$  and  $\mathbf{U}_{k+1}$  are performed with enough accuracy, i.e. small enough  $\delta$  and  $\mathbf{V}$  satisfies some additional properties (in particular  $\|\mathbf{V}\|_2 < 1.5$  and  $\|\mathbf{V}^{-1}\|_2 < 1.5$ , which our matrices  $\mathbf{V}$  do satisfy).

## VI. NUMERICAL EXAMPLES

The algorithms discussed above were implemented in C, using GNU MPFR version 3.1.12, GNU MPFI<sup>3</sup> version 1.5.1

<sup>2</sup><http://www.mpfr.org/>

<sup>3</sup><https://gforge.inria.fr/projects/mpfi/>

	Example 1			Example 2			Example 3			Example 4		
sizes $n$ , $p$ and $q$	$n = 10$ ,	$p = 11$ ,	$q = 1$	$n = 12$ ,	$p = 1$ ,	$q = 25$	$n = 60$ ,	$p = 28$ ,	$q = 14$	$n = 3$ ,	$p = 1$ ,	$q = 4$
$1 - \rho(\mathbf{A})$	$1.39 \times 10^{-2}$			$8.65 \times 10^{-3}$			$1.46 \times 10^{-2}$			$2^{-60}$		
$\max(\mathbf{S}_N)$	$3.88 \times 10^1$			$5.50 \times 10^9$			$2.64 \times 10^2$			-		
$\min(\mathbf{S}_N)$	$1.29 \times 10^0$			$1.0 \times 10^0$			$1.82 \times 10^1$			-		
$\varepsilon$	$2^{-5}$	$2^{-53}$	$2^{-600}$	$2^{-5}$	$2^{-53}$	$2^{-600}$	$2^{-5}$	$2^{-53}$	$2^{-600}$	$2^{-5}$		
$N$	220	2153	29182	308	4141	47811	510	1749	27485	-		
Inversion iterations	0	2	4	2	3	5	1	2	4	-		
overall max precision (bits)	212	293	1401	254	355	1459	232	306	1416	-		
$\mathbf{V}^{-1}$ max precision (bits)	106	173	727	148	204	756	126	177	732	-		
$\mathbf{P}_N$ max precision (bits)	64	84	639	74	86	640	64	87	642	-		
$\mathbf{S}_N$ max precision (bits)	64	79	630	74	107	658	64	85	636	-		
Overall execution time (sec)	0.11	1.53	60.06	0.85	11.54	473.20	45.62	177.90	9376.86	0.00...		
$N$ algo execution time (sec)	0.05	0.61	0.43	0.13	1.29	1.29	26.87	28.26	33.91	0.00...		
Summation algorithm (sec)	0.04	0.92	59.63	0.72	10.25	471.91	18.75	149.64	9410.77	0.00...		

TABLE I. NUMERICAL RESULTS FOR 3 REAL-WORLD AND 1 CONSTRUCTED EXAMPLE

and CLAPACK<sup>4</sup> version 3.2.1. Our implementation was tested on several real-life and random examples:

- The first example comes from Control Theory: the LTI system is extracted from an active controller of vehicle longitudinal oscillation [11], and WCPG matrix is used to determine fixed-point arithmetic scaling of the states and the output.
- The second is a 12<sup>th</sup>-order Butterworth filter, described in  $\rho$ -Direct Form II transposed [18] (a particular algorithm, with low complexity and high robustness to quantization and computational errors), where the errors-to-output LTI system  $\mathcal{H}_e$  is considered (see Figure 1).
- The third one is a large random BIBO stable filter (obtained from the `drss` command of Matlab), with 60 states, 14 inputs and 28 outputs.
- The last one is built with a companion matrix  $\mathbf{A}$  with spectral radius equal to  $1 - 2^{-60}$ .

Experiments were done on a laptop computer with an Intel Core i5 processor running at 2.8 GHz and 16 GB of DDR3 RAM.

The numerical results detailed in Table I show that our algorithm for Worst-Case Peak Gain matrix evaluation with a priori error bound exhibits reasonable performance on typical examples. Even when the a priori error bound is pushed to compute WCPG results with an accuracy way beyond double precision, the algorithm succeeds in computing a result, even though execution time grows pretty high.

Our algorithm includes checks testing that certain properties of matrices are verified, in particular that  $\rho(\mathbf{A}) < 1$  and  $\|\mathbf{T}\|_2 \leq 1$ . Our Example 4, not taken from a real-world application but constructed on purpose, shows that the algorithm correctly detects that the conditions are not fulfilled for that example and refuses to compute any result.

## VII. CONCLUSIONS

With this paper, a reliable, rigorous multiprecision method to compute the Worst-Case Peak Gain matrix has been developed. It relies on Theory of Verified Inclusion, eigenvalue

decomposition to perform matrix powering, some multiple-precision arithmetic basic bricks developed to satisfy absolute error bounds and a detailed step-by-step error analysis.

A C program has been developed and now can be used as a tool for the implementation error analysis of LTI systems, and then the design of reliable finite precision digital algorithms for signal processing and control.

However, some efforts are still required to overcome machine (IEEE754 double) precision eigenvalue decomposition in LAPACK (specially for close-to-instability LTI systems) by using a multiple precision eigensolver. Additionally, as the proofs on the error bounds are pretty complicated, they should be formalized in a Formal Proof Checker, such as Coq or HolLight.

## REFERENCES

- [1] V. Balakrishnan and S. Boyd. On computing the worst-case peak gain of linear systems. *Systems & Control Letters*, 19:265–269, 1992.
- [2] J. Carletta, R. Veillette, F. Krach, and Zhengwei F. Determining appropriate precisions for signals in fixed-point iir filters. In *Design Automation Conference, 2003. Proceedings*, pages 656–661, June 2003.
- [3] H. Dawood. *Theories of Interval Arithmetic: Mathematical Foundations and Applications*. LAP Lambert Academic Publishing, Saarbrücken, 2011.
- [4] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier, and P. Zimmermann. MPFR: A multiple-precision binary floating-point library with correct rounding. *ACM Transactions on Mathematical Software*, 33(2):13:1–13:15, June 2007.
- [5] S. Gershgorin. Über die Abgrenzung der Eigenwerte einer Matrix. *Bull. Acad. Sci. URSS*, 1931(6):749–754, 1931.
- [6] N. J. Higham. *Accuracy and stability of numerical algorithms (2. ed.)*. SIAM, 2002.
- [7] T. Hilaire, P. Chevrel, and J.F. Whidborne. A unifying framework for finite wordlength realizations. *IEEE Trans. on Circuits and Systems*, 8(54):1765–1774, August 2007.
- [8] T. Hilaire and B. Lopez. Reliable implementation of linear filters with fixed-point arithmetic. In *Proc. IEEE Workshop on Signal Processing Systems (SiPS)*, 2013.
- [9] T. Kailath. *Linear Systems*. Prentice-Hall, 1980.
- [10] U. Kulisch and V. Snyder. The exact dot product as basic tool for long interval arithmetic. *Computing*, 91(3):307–313, March 2011.
- [11] D. Lefebvre, P. Chevrel, and S. Richard. An  $H_\infty$  based control design methodology dedicated to the active control of longitudinal oscillations. *IEEE Trans. on Control Systems Technology*, 11(6):948–956, November 2003.

<sup>4</sup><http://www.netlib.org/clapack/>

- [12] J.A. Lopez, C. Carreras, and O. Nieto-Taladriz. Improved interval-based characterization of fixed-point LTI systems with feedback loops. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 26(11):1923–1933, November 2007.
- [13] V. Pan and J. Reif. Efficient parallel solution of linear systems. In *Proceedings of the Seventeenth Annual ACM Symposium on Theory of Computing, STOC '85*, pages 143–152, New York, NY, USA, 1985. ACM.
- [14] S. M. Rump. New results on verified inclusions. In *Accurate Scientific Computations, Symposium, Bad Neuenahr, FRG, March 12-14, 1985, Proceedings*, pages 31–69, 1985.
- [15] S. M. Rump. Solution of linear systems with verified accuracy. *Applied numerical mathematics*, 3(3):233–241, June 1987.
- [16] S. M. Rump. Reliability in computing: The role of interval methods in scientific computing. chapter Algorithms for Verified Inclusions — Theory and Practice, pages 109–126. Academic Press Professional, Inc., San Diego, CA, USA, 1988.
- [17] S. M. Rump. Guaranteed inclusions for the complex generalized eigenproblem. *Computing*, 42(2-3):225–238, September 1989.
- [18] Z. Zhao and G. Li. Roundoff noise analysis of two efficient digital filter structures. *IEEE Transactions on Signal Processing*, 54(2):790–795, February 2006.