

Towards reliable implementation of digital filters

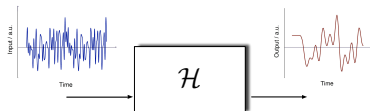
Anastasia Volkova

Sorbonne Universités, UPMC, LIP6

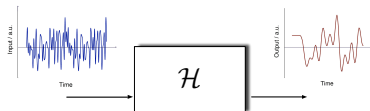
Seminar at ENS Lyon, LIP, AriC team
January 5, 2017



Context: digital filters



Context: digital filters

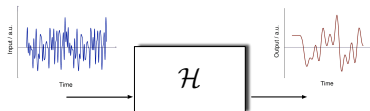


On the one hand

- LTI filter with Infinite Impulse Response
- Its transfer function:

$$H(z) = \frac{\sum_{i=0}^n b_i z^{-i}}{1 + \sum_{i=1}^n a_i z^{-i}}$$

Context: digital filters



On the one hand

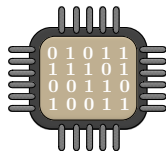
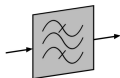
- LTI filter with Infinite Impulse Response
- Its transfer function:

$$H(z) = \frac{\sum_{i=0}^n b_i z^{-i}}{1 + \sum_{i=1}^n a_i z^{-i}}$$

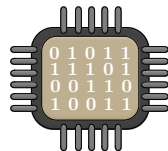
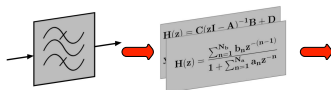
On the other hand

- Hardware or Software target
- Implementation in Fixed- or Floating-Point Arithmetic

Context: implementation of LTI filters

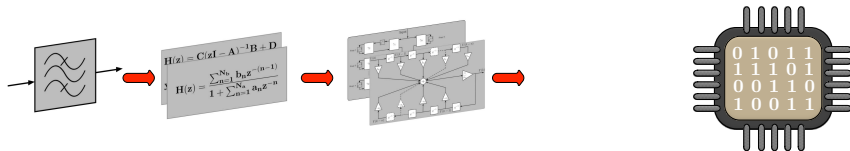


Context: implementation of LTI filters



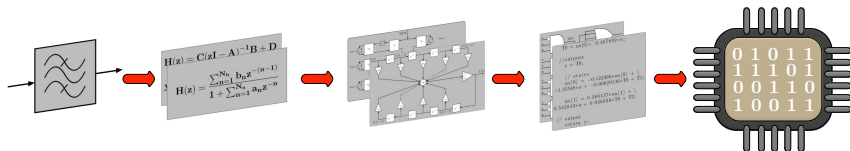
- Transfer function generation

Context: implementation of LTI filters



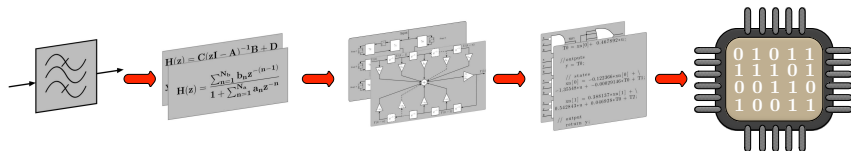
- Transfer function generation
- Algorithm choice: State-space, Direct Form I, Direct Form II, ...

Context: implementation of LTI filters



- Transfer function generation
- Algorithm choice: State-space, Direct Form I, Direct Form II, ...
- Software or Hardware implementation

Context: implementation of LTI filters



- Transfer function generation
 - ☞ Coefficient quantization
- Algorithm choice: State-space, Direct Form I, Direct Form II, ...
 - ☞ Large variety of structures with no common quality criteria
- Software or Hardware implementation
 - ☞ Constraints: power consumption, area, error, speed, etc.
 - ☞ Computational errors due to finite-precision implementation

Filter-to-code generator

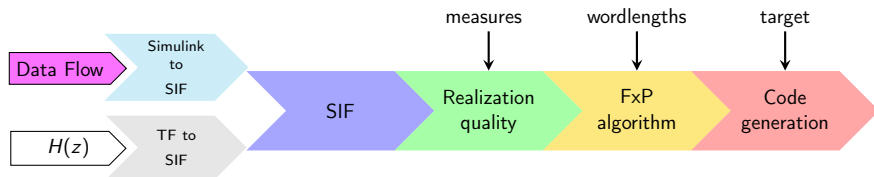


Figure: Automatic Filter Generator Flow

- Use unified *analytical* representation of linear data-flows
- Can describe any analytical and graphical representation
- Adopted numerous classical and developed new quality measures
- Provide **fully rigorous** and **reliable** implementation in Fixed-Point arithmetic
- Generate C (for μ Cs and DSPs) and VHDL code (for ASICs and FPGAs)

Filter-to-code generator

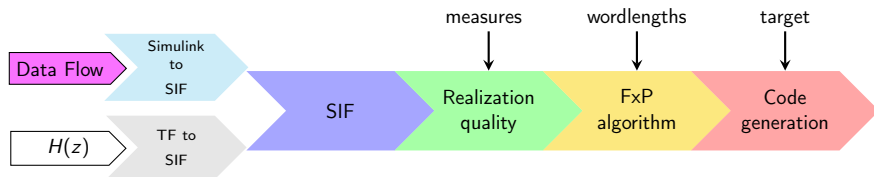


Figure: Automatic Filter Generator Flow

- Use unified *analytical* representation of linear data-flows
- Can describe any analytical and graphical representation
- Adopted numerous classical and developed new quality measures
- Provide **fully rigorous and reliable** implementation in Fixed-Point arithmetic
- Generate C (for μ Cs and DSPs) and VHDL code (for ASICs and FPGAs)

Filter-to-code generator

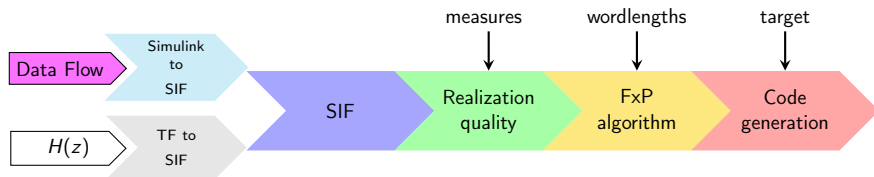


Figure: Automatic Filter Generator Flow

- Use unified *analytical* representation of linear data-flows
- Can describe any analytical and graphical representation
- Adopted numerous classical and developed new quality measures
- Provide **fully rigorous and reliable** implementation in Fixed-Point arithmetic
- Generate C (for μ Cs and DSPs) and **VHDL** code (for ASICs and **FPGAs**)

Reliable implementation of digital filters in Fixed-Point Arithmetic

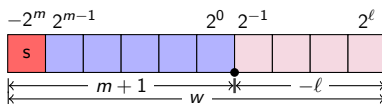
Let $\mathcal{H} := (\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ be a LTI filter:

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

The filter \mathcal{H} is considered Bounded Input Bounded Output stable if

$$\rho(\mathbf{A}) < 1$$

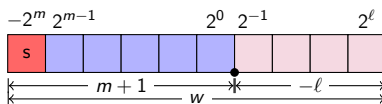
Two's complement Fixed-Point arithmetic



$$y = -2^m y_m + \sum_{i=\ell}^{m-1} 2^i y_i$$

- Wordlength: w
- Most Significant Bit position: m
- Least Significant Bit position: $\ell := m - w + 1$

Two's complement Fixed-Point arithmetic



$$y = -2^m y_m + \sum_{i=\ell}^{m-1} 2^i y_i$$

- $y(k) \in \mathbb{R}$
- wordlength w bits
- minimal Fixed-Point Format (FPF) is the least m :

$$\forall k, \quad y(k) \in [-2^m; 2^m - 2^{m-w+1}]$$

Reliable Fixed-Point implementation

Input:

- $\mathcal{H} = (\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$
- bound on the input interval
- wordlength constraints

Determine: the Fixed-Point Formats s.t.

- the least MSBs
 - no overflows occur
- ↪ must take into account computational errors

Reliable Fixed-Point implementation

Input:

- $\mathcal{H} = (\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$
- bound on the input interval
- wordlength constraints

Determine: the Fixed-Point Formats s.t.

- the least MSBs
 - no overflows occur
- ↪ must take into account computational errors

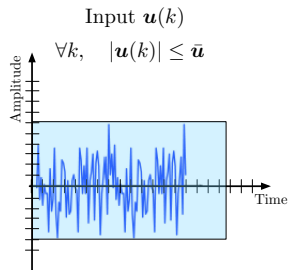
How to proceed:

1. determine the output interval of all variables
2. analyze propagation of the error in filter implementation and determine the FxPF

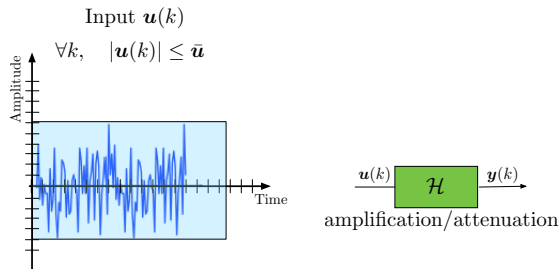
Deducing the output interval¹

¹A.V. et al., "Reliable Evaluation of the Worst-Case Peak Gain Matrix in Multiple Precision", ARITH22, 2015

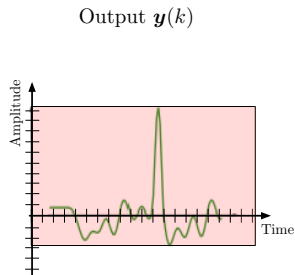
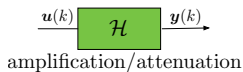
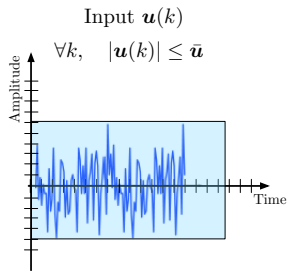
Basic brick: the Worst-Case Peak Gain theorem



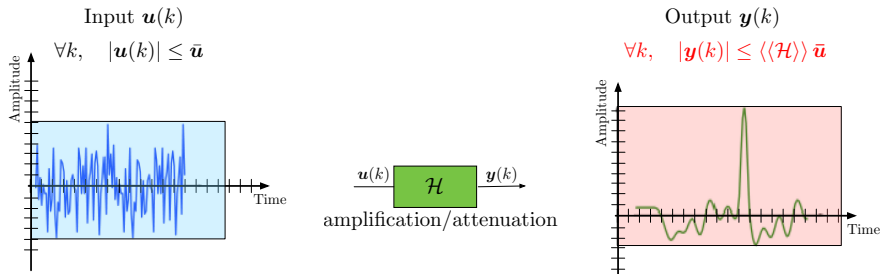
Basic brick: the Worst-Case Peak Gain theorem



Basic brick: the Worst-Case Peak Gain theorem



Basic brick: the Worst-Case Peak Gain theorem



Worst-Case Peak Gain

$$\langle\langle \mathcal{H} \rangle\rangle = |\mathbf{D}| + \sum_{k=0}^{\infty} |\mathbf{C}\mathbf{A}^k\mathbf{B}|$$

Computing the Worst-Case Peak Gain

Problem: compute the Worst-Case Peak Gain with arbitrary precision.

$$\langle\langle\mathcal{H}\rangle\rangle = |\mathbf{D}| + \sum_{k=0}^{\infty} \left| \mathbf{C} \mathbf{A}^k \mathbf{B} \right|$$

- Cannot sum infinitely \Rightarrow need to truncate the sum
- Once the sum is truncated, evaluate it in multiple precision

Truncation

$$\sum_{k=0}^{\infty} |\mathbf{C}\mathbf{A}^k\mathbf{B}| \longrightarrow \sum_{k=0}^N |\mathbf{C}\mathbf{A}^k\mathbf{B}|$$

Truncation

$$\left| \sum_{k=0}^{\infty} |\mathbf{C}\mathbf{A}^k\mathbf{B}| - \sum_{k=0}^N |\mathbf{C}\mathbf{A}^k\mathbf{B}| \right| \leq \varepsilon_1$$

Compute an approximate lower bound on truncation order N such that the truncation error is smaller than ε_1 .

Lower bound on truncation order N

$$N \geq \left\lceil \frac{\log \frac{\varepsilon_1}{\|\mathbf{M}\|_{\min}}}{\log \rho(\mathbf{A})} \right\rceil, \quad \text{with } \mathbf{M} := \sum_{l=1}^n \frac{|\mathbf{R}_l|}{1 - |\lambda_l|} \frac{|\lambda_l|}{\rho(\mathbf{A})}$$

where

λ – eigenvalues of matrix \mathbf{A}

\mathbf{R}_l – l^{th} residue matrix computed out of $\mathbf{C}, \mathbf{B}, \lambda$

$$\sum_{k=0}^N |C \mathbf{A}^k B|$$

Powering

$$\sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B}|$$



×



=

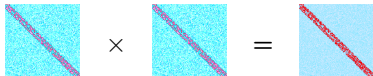


cancellation

$$\sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B}|$$



cancellation

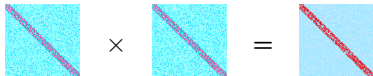


less cancellation

$$\sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B}|$$



cancellation



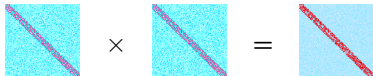
less cancellation

$$\mathbf{A} = \mathbf{X} \mathbf{E} \mathbf{X}^{-1}$$

$$\sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B}|$$



cancellation



less cancellation

$$\mathbf{A} = \mathbf{X} \mathbf{E} \mathbf{X}^{-1}$$

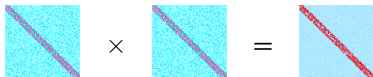
$$\mathbf{V} \approx \mathbf{X} \text{ and } \mathbf{T} \approx \mathbf{E}$$

Powering

$$\sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B}|$$



cancellation



less cancellation

$$\mathbf{A} = \mathbf{X} \mathbf{E} \mathbf{X}^{-1}$$

$$\mathbf{V} \approx \mathbf{X} \text{ and } \mathbf{T} \approx \mathbf{E}$$

$$\mathbf{T} \approx \mathbf{V}^{-1} \times \mathbf{A} \times \mathbf{V}$$

$$\mathbf{A}^k \approx \mathbf{V} \times \mathbf{T}^k \times \mathbf{V}^{-1}$$

$$\left| \sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B}| - \sum_{k=0}^N |\mathbf{C} \mathbf{V} \mathbf{T}^k \mathbf{V}^{-1} \mathbf{B}| \right| \leq \varepsilon_2$$

Given matrix \mathbf{V} compute \mathbf{T} such that the error of substitution of the product $\mathbf{V} \mathbf{T}^k \mathbf{V}^{-1}$ instead of \mathbf{A}^k is less than ε_2 .

Further steps

$$\left| \sum_{k=0}^N |\mathbf{C} \mathbf{A}^k \mathbf{B}| - \sum_{k=0}^N |\mathbf{C} \mathbf{V} \mathbf{T}^k \mathbf{V}^{-1} \mathbf{B}| \right| \leq \varepsilon_2$$

Apply the same approach for the other steps:

$$\left| \sum_{k=0}^N |\mathbf{C} \mathbf{V} \mathbf{T}^k \mathbf{V}^{-1} \mathbf{B}| - \sum_{k=0}^N |\mathbf{C}' \mathbf{T}^k \mathbf{B}'| \right| \leq \varepsilon_3$$

$$\left| \sum_{k=0}^N |\mathbf{C}' \mathbf{T}^k \mathbf{B}'| - \sum_{k=0}^N |\mathbf{C}' \mathbf{P}_k \mathbf{B}'| \right| \leq \varepsilon_4$$

$$\left| \sum_{k=0}^N |\mathbf{C}' \mathbf{P}_k \mathbf{B}'| - \sum_{k=0}^N |\mathbf{L}_k| \right| \leq \varepsilon_5$$

$$\left| \sum_{k=0}^N |\mathbf{L}_k| - \mathbf{S}_N \right| \leq \varepsilon_6$$

Requirement:

Provide matrix operations which satisfy an element-by-element absolute error bound δ given in the argument.

Requirement:

Provide matrix operations which satisfy an element-by-element absolute error bound δ given in the argument.

Problem:

In fixed-precision FP arithmetic such absolute bound is not generally possible.

Basic bricks

Requirement:

Provide matrix operations which satisfy an element-by-element absolute error bound δ given in the argument.

Problem:

In fixed-precision FP arithmetic such absolute bound is not generally possible.

Solution:

Use multiple-precision FP arithmetic and dynamically adapt precision of the result variables.

What if coefficients are not exact?

Cases when \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} are not exact:

- coefficients are results of finite-precision computations (e.g. quantization, SIF \leftrightarrow State-Space transformation etc.)

What if coefficients are not exact?

Cases when \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} are not exact:

- coefficients are results of finite-precision computations (e.g. quantization, SIF \leftrightarrow State-Space transformation etc.)

To take these properties into account we use Interval Arithmetic.
—→ Need to compute the WCPG in interval arithmetic.

What if coefficients are not exact?

Cases when \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} are not exact:

- coefficients are results of finite-precision computations (e.g. quantization, SIF \leftrightarrow State-Space transformation etc.)

To take these properties into account we use Interval Arithmetic.
—→ Need to compute the WCPG in interval arithmetic.

Notation: interval matrix $\mathbf{M}^{\mathcal{I}}$ is centered at $mid(\mathbf{M}^{\mathcal{I}})$ and has radius $rad(\mathbf{M}^{\mathcal{I}})$.

We suppose all interval arithmetic to be in multiple-precision.

Interval WCPG computation

Problem: compute the interval Worst-Case Peak Gain matrix

$$\langle\langle\mathcal{H}^{\mathcal{I}}\rangle\rangle = |\mathbf{D}^{\mathcal{I}}| + \sum_{k=0}^{\infty} \left| \mathbf{C}^{\mathcal{I}} \mathbf{A}^{\mathcal{I}^k} \mathbf{B}^{\mathcal{I}} \right|.$$

Ensure:

- enclosure property: $\forall \mathcal{H} \in \mathcal{H}^{\mathcal{I}} \implies \langle\langle\mathcal{H}\rangle\rangle \in \langle\langle\mathcal{H}_N^{\mathcal{I}}\rangle\rangle$
- if coefficients' radii $\rightarrow 0$ and the precision $\rightarrow \infty$, then $\langle\langle\mathcal{H}_N^{\mathcal{I}}\rangle\rangle$ is a ε -neighbourhood of the exact WCPG matrix for arbitrary $\varepsilon > 0$

Computing the eigensystem of interval matrix

Eigenvalues of interval matrix

Compute enclosures $\lambda^{\mathcal{I}}$ such that $\forall \mathbf{A} \in \mathbf{A}^{\mathcal{I}}, \lambda(\mathbf{A}) \in \lambda^{\mathcal{I}}$

Approach

Following the works of Xu and Rachid (1996) and Rohn(1998), use the Generalized Gershgorin's Circles theorem.

Computing the eigensystem of interval matrix

Eigenvalues of interval matrix

Compute enclosures $\lambda^{\mathcal{I}}$ such that $\forall \mathbf{A} \in \mathbf{A}^{\mathcal{I}}, \lambda(\mathbf{A}) \in \lambda^{\mathcal{I}}$

Approach

Following the works of Xu and Rachid (1996) and Rohn(1998), use the Generalized Gershgorin's Circles theorem.

Eigenvectors of interval matrix

Given the enclosures on eigenvalues $\lambda^{\mathcal{I}}$, compute enclosures $\mathbf{V}^{\mathcal{I}}$ such that $\forall \lambda \in \lambda^{\mathcal{I}}, \forall \mathbf{A} \in \mathbf{A}^{\mathcal{I}}$ if $\mathbf{A}\mathbf{V} = \lambda\mathbf{V}$, then $\mathbf{V} \in \mathbf{V}^{\mathcal{I}}$.

Approach

Use Rump's theory of Verified Inclusions.

Determining the Fixed-Point Formats²

²A.V. et al., "Determining Fixed-Point Formats for a Digital Filter Implementation using the Worst-Case Peak Gain Measure", Asilomar 49, 2015

Determining the Fixed-Point Formats

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

We know that if $\forall k, |\mathbf{u}_i(k)| \leq \bar{\mathbf{u}}_i$, then

$$\forall k, \quad |\mathbf{y}_i(k)| \leq (\langle\langle \mathcal{H} \rangle\rangle \bar{\mathbf{u}})_i.$$

Determining the Fixed-Point Formats

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

We know that if $\forall k, |\mathbf{u}_i(k)| \leq \bar{\mathbf{u}}_i$, then

$$\forall k, \quad |\mathbf{y}_i(k)| \leq (\langle\langle \mathcal{H} \rangle\rangle \bar{\mathbf{u}})_i.$$

We need to find the least integer \mathbf{m}_y such that

$$\forall k, \quad |\mathbf{y}_i(k)| \leq 2^{\mathbf{m}_{y_i}} - 2^{\mathbf{m}_{y_i} - \mathbf{w}_{y_i} + 1}.$$

Determining the Fixed-Point Formats

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

We know that if $\forall k, |\mathbf{u}_i(k)| \leq \bar{\mathbf{u}}_i$, then

$$\forall k, \quad |\mathbf{y}_i(k)| \leq (\langle\langle \mathcal{H} \rangle\rangle \bar{\mathbf{u}})_i.$$

We need to find the least integer \mathbf{m}_y such that

$$\forall k, \quad |\mathbf{y}_i(k)| \leq 2^{\mathbf{m}_{y_i}} - 2^{\mathbf{m}_{y_i} - \mathbf{w}_{y_i} + 1}.$$

It easy to show that \mathbf{m}_y can be computed with

$$\mathbf{m}_{y_i} = \lceil \log_2 (\langle\langle \mathcal{H} \rangle\rangle \bar{\mathbf{u}})_i - \log_2 (1 - 2^{1 - \mathbf{w}_{y_i}}) \rceil.$$

Determining the Fixed-Point Formats

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

We know that if $\forall k, |\mathbf{u}_i(k)| \leq \bar{\mathbf{u}}_i$, then

$$\forall k, \quad |\mathbf{y}_i(k)| \leq (\langle\langle\mathcal{H}\rangle\rangle \bar{\mathbf{u}})_i.$$

We need to find the least integer \mathbf{m}_y such that

$$\forall k, \quad |\mathbf{y}_i(k)| \leq 2^{\mathbf{m}_{y_i}} - 2^{\mathbf{m}_{y_i} - \mathbf{w}_{y_i} + 1}.$$

It easy to show that \mathbf{m}_y can be computed with

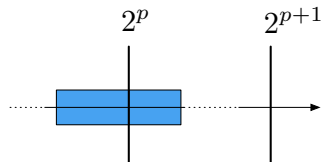
$$\mathbf{m}_{y_i} = \lceil \log_2 (\langle\langle\mathcal{H}\rangle\rangle \bar{\mathbf{u}})_i - \log_2 (1 - 2^{1-\mathbf{w}_{y_i}}) \rceil.$$

Control the accuracy of the WCPG such that $0 \leq \hat{\mathbf{m}}_{y_i} - \mathbf{m}_{y_i} \leq 1$

Open question: the off-by-one problem

$$\widehat{m}_{y_i} = \lceil \mathfrak{m} \rceil$$

Problem: interval \mathfrak{m} contains a power of 2.

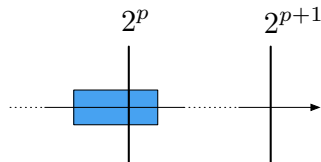


Open question: the off-by-one problem

$$\widehat{m}_{y_i} = \lceil \mathfrak{m} \rceil$$

Problem: interval \mathfrak{m} contains a power of 2.

Technique: Ziv's strategy to reduce interval



Open question: the off-by-one problem

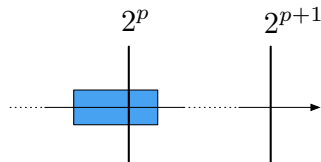
$$\widehat{m}_{y_i} = \lceil \mathfrak{m} \rceil$$

Problem: interval \mathfrak{m} contains a power of 2.

Technique: Ziv's strategy to reduce interval

Dilemma:

- propagation of computational errors
- or overestimation in linear filter decomposition?

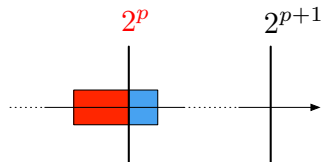


Open question: the off-by-one problem

$$\widehat{\mathbf{m}}_{y_i} = \lceil \mathfrak{m} \rceil$$

Problem: interval \mathfrak{m} contains a power of 2.

Technique: Ziv's strategy to reduce interval



Dilemma:

- propagation of computational errors
- or overestimation in linear filter decomposition?

Possible approach:

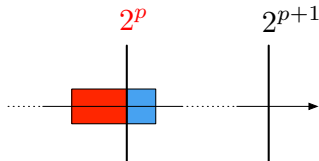
- Assume the format $\widehat{\mathbf{m}} = p$
- Does there exist a *reachable* $\mathbf{x}^\diamond(k)$ s.t. $\mathbf{y}^\diamond(k)$ overflows ?

Open question: the off-by-one problem

$$\widehat{\mathbf{m}}_{y_i} = \lceil \mathfrak{m} \rceil$$

Problem: interval \mathfrak{m} contains a power of 2.

Technique: Ziv's strategy to reduce interval



Dilemma:

- propagation of computational errors
- or overestimation in linear filter decomposition?

Possible approach:

- Assume the format $\widehat{\mathbf{m}} = p$
- Does there exist a *reachable* $\mathbf{x}^\diamond(k)$ s.t. $\mathbf{y}^\diamond(k)$ overflows ?

Technique: SMT? integer linear programming? lattice algorithms?

Taking the quantization errors into account

The exact filter \mathcal{H} is:

$$\mathcal{H} \quad \left\{ \begin{array}{lcl} \mathbf{x}(k+1) & = & \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \mathbf{y}(k) & = & \mathbf{C}\mathbf{x}(k) + \mathbf{D}\mathbf{u}(k) \end{array} \right.$$

Taking the quantization errors into account

The actually implemented filter \mathcal{H}^\diamond is:

$$\mathcal{H}^\diamond \begin{cases} \mathbf{x}^\diamond(k+1) &= \diamond_{m_x}(\mathbf{A}\mathbf{x}^\diamond(k) + \mathbf{B}\mathbf{u}(k)) \\ \mathbf{y}^\diamond(k) &= \diamond_{m_y}(\mathbf{C}\mathbf{x}^\diamond(k) + \mathbf{D}\mathbf{u}(k)) \end{cases}$$

where \diamond_m is some operator ensuring faithful rounding:

$$|\diamond_m(x) - x| \leq 2^{m-w+1}.$$

Taking the quantization errors into account

The actually implemented filter \mathcal{H}^\diamond is:

$$\mathcal{H}^\diamond \begin{cases} \mathbf{x}^\diamond(k+1) = & \mathbf{A}\mathbf{x}^\diamond(k) + \mathbf{B}\mathbf{u}(k) + \boldsymbol{\varepsilon}_x(k) \\ \mathbf{y}^\diamond(k) = & \mathbf{C}\mathbf{x}^\diamond(k) + \mathbf{D}\mathbf{u}(k) + \boldsymbol{\varepsilon}_y(k) \end{cases}$$

with

$$|\varepsilon_x(k)| \leq 2^{m_x - w_x + 1} \quad \text{and} \quad |\varepsilon_y(k)| \leq 2^{m_y - w_y + 1}.$$

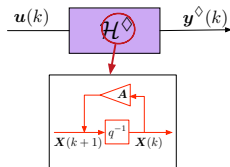
Taking the quantization errors into account

The actually implemented filter \mathcal{H}^\diamond is:

$$\mathcal{H}^\diamond \begin{cases} \mathbf{x}^\diamond(k+1) = & \mathbf{A}\mathbf{x}^\diamond(k) + \mathbf{B}\mathbf{u}(k) + \epsilon_x(k) \\ \mathbf{y}^\diamond(k) = & \mathbf{C}\mathbf{x}^\diamond(k) + \mathbf{D}\mathbf{u}(k) + \epsilon_y(k) \end{cases}$$

with

$$|\epsilon_x(k)| \leq 2^{m_x - w_x + 1} \quad \text{and} \quad |\epsilon_y(k)| \leq 2^{m_y - w_y + 1}.$$



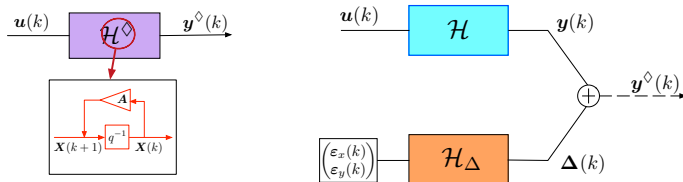
Taking the quantization errors into account

The actually implemented filter \mathcal{H}^\diamond is:

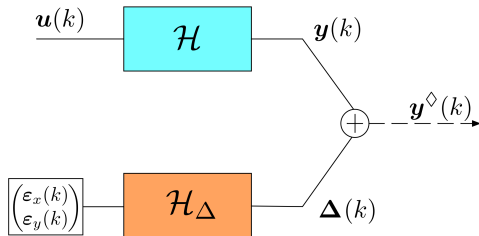
$$\mathcal{H}^\diamond \begin{cases} \mathbf{x}^\diamond(k+1) = & \mathbf{A}\mathbf{x}^\diamond(k) + \mathbf{B}\mathbf{u}(k) + \varepsilon_x(k) \\ \mathbf{y}^\diamond(k) = & \mathbf{C}\mathbf{x}^\diamond(k) + \mathbf{D}\mathbf{u}(k) + \varepsilon_y(k) \end{cases}$$

with

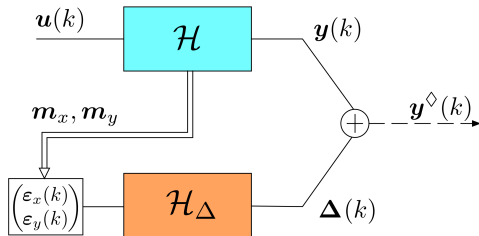
$$|\varepsilon_x(k)| \leq 2^{m_x - w_x + 1} \quad \text{and} \quad |\varepsilon_y(k)| \leq 2^{m_y - w_y + 1}.$$



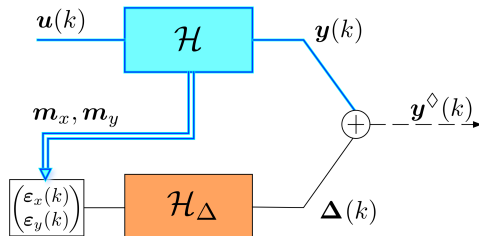
Algorithm



Algorithm

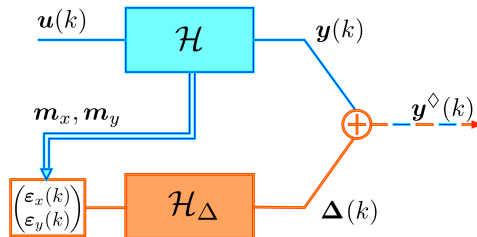


Algorithm



Step 1: Determine the initial guess MSBs \mathbf{m}_y for the exact filter \mathcal{H}

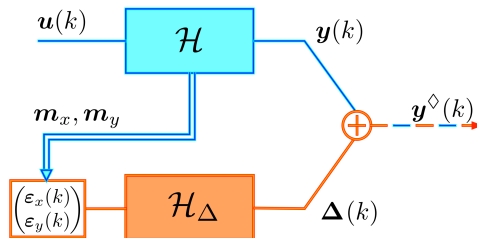
Algorithm



Step 1: Determine the initial guess MSBs \mathbf{m}_y for the exact filter \mathcal{H}

Step 2: Compute the error-filter \mathcal{H}_Δ , induced by the format \mathbf{m}_y and deduce the MSBs $\mathbf{m}_\zeta^\diamond$

Algorithm



- Step 1: Determine the initial guess MSBs \mathbf{m}_y for the exact filter \mathcal{H}
- Step 2: Compute the error-filter \mathcal{H}_Δ , induced by the format \mathbf{m}_y and deduce the MSBs $\mathbf{m}_\zeta^\diamond$
- Step 3: If $\mathbf{m}_{y_i}^\diamond == \mathbf{m}_{y_i}$ then return $\mathbf{m}_{y_i}^\diamond$
otherwise $\mathbf{m}_{y_i} \leftarrow \mathbf{m}_{y_i} + 1$ and go to Step 2.

Numerical example

Example:

- Random filter with 3 states, 1 input, 1 output
- $\bar{u} = 5.125$, wordlengths set to 7 bits

	states			output
	$x_1(k)$	$x_2(k)$	$x_3(k)$	$y(k)$
Step 1	6	7	5	6
Step 2	6	7	6	6
Step 3	6	7	6	6

Table: Evolution of the MSB positions

Numerical examples

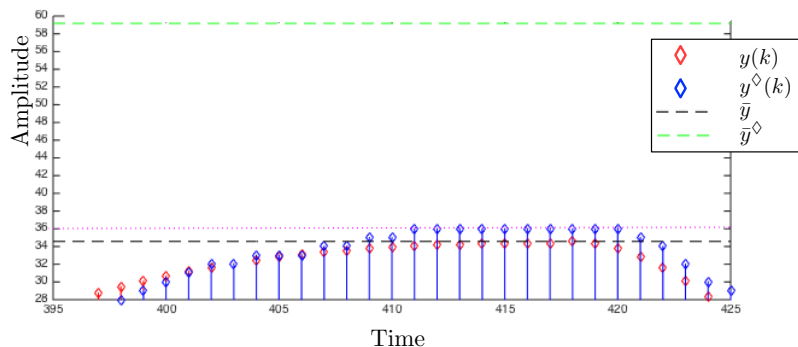


Figure: The exact and quantized outputs of the example.
Quantized output does not pass over to the next binade.

Numerical examples

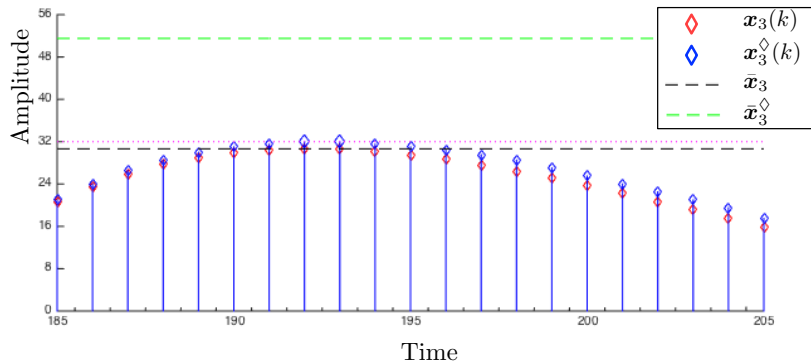


Figure: The exact and quantized third state of the example.
Quantized state passes over to the next binade.

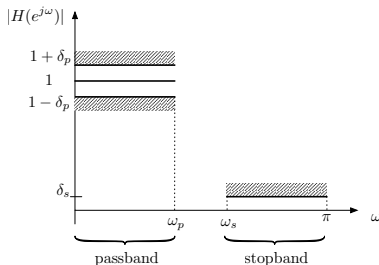
Work in progress:
"Verification of digital *implemented* filters against
specification."

Filter specifications

The filters are designed from transfer function specifications in frequency domain:

$$\underline{\beta} \leq |H(e^{j\omega})| \leq \overline{\beta}, \quad \forall \omega \in [\omega_1, \omega_2]$$

Example: lowpass filter



$$\begin{cases} 1 - \delta_p \leq |H(e^{j\omega})| \leq 1 + \delta_p & \forall \omega \in [0, \omega_p] \quad (\text{passband}) \\ |H(e^{j\omega})| \leq \delta_s & \forall \omega \in [\omega_s, \pi] \quad (\text{stopband}) \end{cases}$$

Reformulating the problem

Need to show that $\forall z \in \{e^{j\omega} | \omega \in \Omega \subset [0, \pi]\}$

$$\underline{\beta} \leq |H(z)| \leq \overline{\beta}$$

Reformulating the problem

Need to show that $\forall z \in \{e^{j\omega} | \omega \in \Omega \subset [0, \pi]\}$

$$\underline{\beta}^2 \leq |H(z)|^2 \leq \overline{\beta}^2$$

Reformulating the problem

Need to show that $\forall z \in \{e^{j\omega} | \omega \in \Omega \subset [0, \pi]\}$

$$\underline{\beta}^2 \leq |H(z)|^2 \leq \overline{\beta}^2$$

We have that

$$|H(z)|^2 = \frac{|b(z)|^2}{|a(z)|^2} = \frac{b(z)\overline{b(\overline{z})}}{a(z)\overline{a(\overline{z})}} = \frac{b(z)b(\frac{1}{\overline{z}})}{a(z)a(\frac{1}{\overline{z}})} =: \frac{P(z)}{Q(z)},$$

where $P(z)$ and $Q(z)$ polynomials with real coefficients.

Reformulating the problem

Need to show that $\forall z \in \{e^{j\omega} | \omega \in \Omega \subset [0, \pi]\}$

$$\underline{\beta}^2 \leq |H(z)|^2 \leq \overline{\beta}^2$$

We have that

$$|H(z)|^2 = \frac{|b(z)|^2}{|a(z)|^2} = \frac{b(z)\overline{b(\overline{z})}}{a(z)\overline{a(\overline{z})}} = \frac{b(z)b(\frac{1}{\overline{z}})}{a(z)a(\frac{1}{\overline{z}})} =: \frac{P(z)}{Q(z)},$$

where $P(z)$ and $Q(z)$ polynomials with real coefficients.

Therefore, we need to show:

$$\underline{\beta}' \leq \frac{P(z)}{Q(z)} \leq \overline{\beta}'$$

We use the Sollya tool to prove it.

Basic brick: verifying the bounds of a rational function

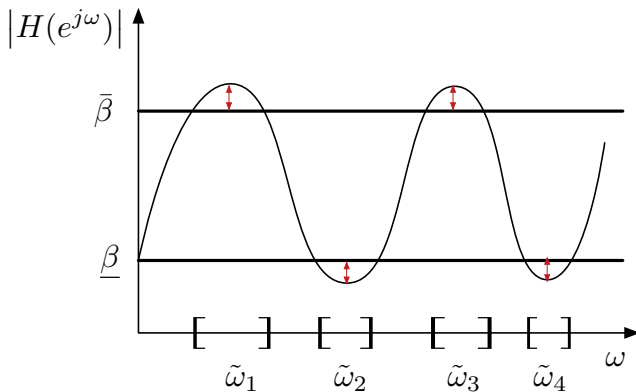


Figure: If a specification is not satisfied, our algorithm returns the problematic frequencies as small intervals $\tilde{\omega}_i$ and maximum overflows within those frequencies.

Verifying a LTI filter implementation

Given a filter implementation compute the corresponding transfer function.

Transfer function of a single-input single-output state-space system:

$$H(z) = \mathbf{c}(z\mathbf{I} - \mathbf{A})^{-1}\mathbf{b} + d.$$

Verifying a LTI filter implementation

Given a filter implementation compute the corresponding transfer function.

Transfer function of a single-input single-output state-space system:

$$H(z) = \mathbf{c}(z\mathbf{I} - \mathbf{A})^{-1}\mathbf{b} + d.$$

Using the eigendecomposition $\mathbf{A} = \mathbf{V}\mathbf{E}\mathbf{V}^{-1}$:

$$H(z) = \frac{P(z)}{Q(z)} + d$$

$$P(z) = \sum_{i=1}^n (\mathbf{c}\mathbf{V})_i (\mathbf{V}^{-1}\mathbf{b})_i \prod_{j \neq i} (z - \lambda_j)$$

$$Q(z) = \prod_{j=1}^n (z - \lambda_j)$$

Verifying a LTI filter implementation

Given a filter implementation compute the corresponding transfer function.

Transfer function of a single-input single-output state-space system:

$$H(z) = \mathbf{c}(z\mathbf{I} - \mathbf{A})^{-1}\mathbf{b} + d.$$

Using the eigendecomposition $\mathbf{A} = \mathbf{V}\mathbf{E}\mathbf{V}^{-1}$:

$$H(z) = \frac{P(z)}{Q(z)} + d$$

$$P(z) = \sum_{i=1}^n (\mathbf{c}\mathbf{V})_i (\mathbf{V}^{-1}\mathbf{b})_i \prod_{j \neq i} (z - \lambda_j)$$

$$Q(z) = \prod_{j=1}^n (z - \lambda_j)$$

We can compute an approximation $\hat{H}(z)$ in Multiple Precision arithmetic.

Verifying a LTI filter implementation

Bounding the approximation error:

$$\begin{array}{ccccc} \text{dSS} & - & \widehat{\text{dSS}} & = & \Delta \text{dSS} \\ & \nearrow \text{red dashed} & \uparrow \text{green solid} & & \searrow \text{red dashed} \\ H & - & \widehat{H} & = & \Delta H \end{array}$$

Figure: In red: inexact transformation. In green: exact transformation.

Verifying a LTI filter implementation

Bounding the approximation error:

$$\begin{array}{ccccc} \text{dSS} & \xrightarrow{\text{green}} & \widehat{\text{dSS}} & = & \Delta\text{dSS} \\ & \searrow \text{red} & \uparrow \text{green} & & \nearrow \text{red} \\ H & \xrightarrow{\text{red}} & \widehat{H} & = & \Delta H \end{array}$$

Figure: In red: inexact transformation. In green: exact transformation.

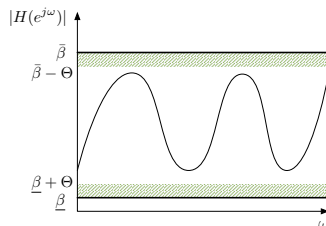
Lemma

Given a discrete state-space system, the error of a multiple precision approximation $\widehat{H}(z)$ on its transfer function is bounded by:

$$\left| H(z) - \widehat{H}(z) \right| \leq \Theta,$$

where $\Theta = \langle\langle \Delta\text{dSS} \rangle\rangle + \varepsilon$ is the Worst-Case Peak Gain of the system ΔdSS computed with arbitrary small absolute error bounded by the $\varepsilon > 0$.

Verifying a LTI filter implementation



Lemma

If the approximation $\hat{H}(z)$ satisfies

$$\underline{\beta} + \Theta \leq \left| \hat{H}(e^{i\omega}) \right| \leq \bar{\beta} - \Theta, \quad \forall \omega \in \Omega \quad (1)$$

Then the exact transfer function $H(z)$ verifies the initial specifications

$$\underline{\beta} \leq |H(e^{i\omega})| \leq \bar{\beta}, \quad \forall \omega \in \Omega \quad (2)$$

Numerical results: comparing filter design tools

		Butterworth	Chebyshev	Elliptic
lowpass	MATLAB	1.29e-17	7.93e-17	✓
	SciPy	2.14e-15	4.48e-2	4.48e-2
highpass	MATLAB	2.77e-16	6.94e-17	4.48e-2
	SciPy	3.02e-15	2.29e-16	4.48e-2
bandpass	MATLAB	3.04e-17	✓	✓
	SciPy	✓	4.48e-2	4.48e-2
bandstop	MATLAB	4.59e-16	3.09e-15	✓
	SciPy	✓	6.36e-15	7.02e-6

Table: Quality of the filters designed with MATLAB vs. SciPy for some simple filters. For example, a lowpass filter here has $\delta_p = 10^{-\frac{1}{20}}$, $\delta_s = 10^{-1}$, $\omega_p = 0.4\pi$, $\omega_s = 0.5\pi$

Numerical results: comparing filter structures

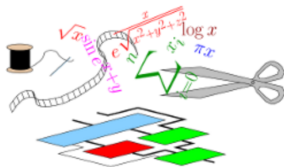
filter		Butterworth			Chebyshev		
	wordlength	32	16	8	32	16	8
DFilt	margin	✓	✓	1.17e-2	1.97e-10	7.03e-6	7.42e-3
	time	12s	15s	1min17s	1min22s	35s	18s
ρ DFilt	margin	✓	6.78e-2	✗	1.21e-9	4.72e-4	8.35e-2
	time	13s	3min17s	unstable	1min26s	1min02s	17s
SS	margin	✓	9.11e-7	9.64e-3	2.06e-10	7.64e-6	2.23e-3
	time	15s	2min4s	1min01s	23s	49s	18s

Table: Filter structures for a lowpass MATLAB-designed filter.

Work in progress:
"Code generation for FPGAs"

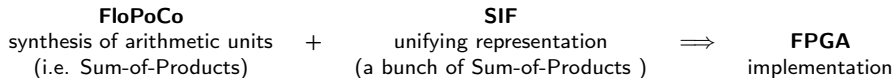
with Matei Istoan and Florent de Dinechin

FloPoCo

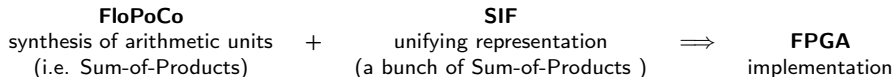


Circuits computing just right

Binding with FloPoCo



Binding with FloPoCo

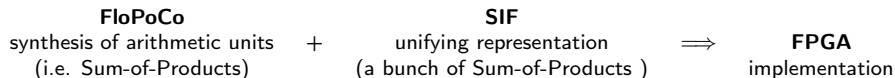


FloPoCo: requires the the desired Fixed-Point Format of the output of a Sum-of-Product.

SIF: we deduce a lower bound on the error of computation of each Sum-of-Product s.t. the error-bound *on the filter's output* is respected.

Naive approach: the error-budget is distrubuted equally

Binding with FloPoCo



FloPoCo: requires the the desired Fixed-Point Format of the output of a Sum-of-Product.

SIF: we deduce a lower bound on the error of computation of each Sum-of-Product s.t. the error-bound *on the filter's output* is respected.

Naive approach: the error-budget is distrubuted equally

We can do better

Take into account different impact of different Sums-of-Products.
BUT need to solve non-linear optimization problem

Conclusion

- Proposed a new completely rigorous approach for the Fixed-Point implementation of linear digital filters
- Provided reliable evaluation of the WCPG measure
- Applied the WCPG measure to determine the FxPF that guarantee that overflow occurs
- Proposed a new approach for the accurate computation of the transfer function of an implemented filter
- Proposed a rigorous approach on the verification of a LTI filter implementation against its specification
- Binding with FloPoCo gives possibility to fairly compare different filter realizations

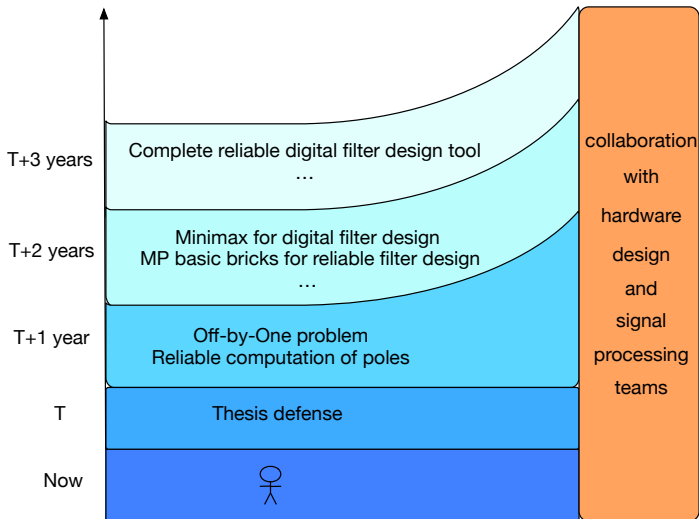
Filter design and implementation:

- Plug all the stages of the generator into optimization routines
- Fairly compare finite precision implementation under various constraints
- Analysis of (short) floating-point implementations

From computer arithmetic point of view:

- Eigendecomposition with a priori absolute error bound
- Off-by-one problem
- Optimize the quantization of filter coefficients using our specification verification algorithm
- Basic bricks multiple precision algorithms (transfer function, Lyapunov equations,...)

Future plans...



Thank you!
Questions?

L_2 -norm evaluation

Another related problem is the reliable evaluation of the L_2 -norm. If \mathbf{H} is a transfer function, then its L_2 -norm is defined by

$$\|\mathbf{H}\|_2 \triangleq \sqrt{\frac{1}{2\pi} \int_0^{2\pi} \|\mathbf{H}(e^{j\omega})\|_F^2 d\omega}$$

Parseval's theorem gives another expression when \mathbf{H} is described with state-space matrices \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{D} :

$$\begin{aligned} \|\mathbf{H}\|_2 &= \sqrt{\text{tr}(\mathbf{C}\mathbf{W}_c\mathbf{C}^\top + \mathbf{D}\mathbf{D}^\top)} \\ &= \sqrt{\text{tr}(\mathbf{B}^\top\mathbf{W}_o\mathbf{B} + \mathbf{D}^\top\mathbf{D})} \end{aligned}$$

where \mathbf{W}_c and \mathbf{W}_o are the controllability and observability Gramians of the system.

- \mathbf{W}_c is the controllability Gramian of the system.

$$\mathbf{W}_c \triangleq \sum_{k=0}^{\infty} (\mathbf{A}^k \mathbf{B})(\mathbf{A}^k \mathbf{B})^\top$$

\mathbf{W}_c is the solution of the discrete-time Lyapunov equation

$$\mathbf{W}_c = \mathbf{A} \mathbf{W}_c \mathbf{A}^\top + \mathbf{B} \mathbf{B}^\top$$

- \mathbf{W}_o is the observability Gramian of the system.

$$\mathbf{W}_o \triangleq \sum_{k=0}^{\infty} (\mathbf{C} \mathbf{A}^k)^\top (\mathbf{C} \mathbf{A}^k)$$

\mathbf{W}_o is the solution of the discrete-time Lyapunov equation

$$\mathbf{W}_o = \mathbf{A}^\top \mathbf{W}_o \mathbf{A} + \mathbf{C}^\top \mathbf{C}$$

Computation of the Gramians

The Gramians are usually computed by solving the discrete-time Lyapunov equation $\mathbf{X} = \mathbf{A}\mathbf{X}\mathbf{A}^\top + \mathbf{Q}$

The following methods can be used:

- solve $(\mathbf{I} - \mathbf{A} \otimes \mathbf{A})\mathbf{x} = \mathbf{q}$
where $\mathbf{x} = \text{Vec}(\mathbf{X})$ and $\mathbf{q} = \text{Vec}(\mathbf{Q})$
→ numerically inefficient
- use infinite sum $\sum_{k=0}^{\infty} \mathbf{A}^k \mathbf{Q} \mathbf{A}^{k\top}$
→ may required a lot of computation
- use Hammarling's method, based on Schur decomposition of matrix \mathbf{A}
→ efficient, but required a deep analysis of the computational errors of the algorithm

see "Computational methods for linear matrix equations", V. Simoncini

Questions

- How to have a reliable evaluation of the L_2 -norm in multiple precision
- How to proceed when \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} are interval matrices (small radii, containing previously computed errors)

Step 1. Bound on truncation error

Truncation error is the tail of the infinite sum:

$$\sum_{k>N} \left| \mathbf{C} \mathbf{A}^k \mathbf{B} \right|$$

Step 1. Bound on truncation error

Truncation error is the tail of the infinite sum:

$$\sum_{k>N} |\mathbf{C}\mathbf{A}^k\mathbf{B}|$$

Suppose $\mathbf{A} = \mathbf{X}\mathbf{E}\mathbf{X}^{-1}$, where $\mathbf{E} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is the eigenvalue matrix and \mathbf{X} is the eigenvector matrix. Then,

$$\mathbf{C}\mathbf{A}^k\mathbf{B} = \mathbf{C}\mathbf{X}\mathbf{E}^k\mathbf{X}^{-1}\mathbf{B} = \sum_{l=1}^n R_l \lambda_l^k$$

Step 1. Bound on truncation error

Truncation error is the tail of the infinite sum:

$$\sum_{k>N} |\mathbf{C}\mathbf{A}^k\mathbf{B}|$$

Suppose $\mathbf{A} = \mathbf{X}\mathbf{E}\mathbf{X}^{-1}$, where $\mathbf{E} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is the eigenvalue matrix and \mathbf{X} is the eigenvector matrix. Then,

$$\mathbf{C}\mathbf{A}^k\mathbf{B} = \mathbf{C}\mathbf{X}\mathbf{E}^k\mathbf{X}^{-1}\mathbf{B} = \sum_{l=1}^n R_l \lambda_l^k$$

Bound on truncation error

$$\sum_{k>N} |\mathbf{C}\mathbf{A}^k\mathbf{B}| \leq \rho(\mathbf{A})^{N+1} M$$

$$M := \sum_{l=1}^n \frac{|R_l|}{1 - |\lambda_l|} \frac{|\lambda_l|}{\rho(\mathbf{A})}$$

Step 1. Bound on truncation error

Truncation error is the tail of the infinite sum:

$$\sum_{k>N} |\mathbf{C}\mathbf{A}^k\mathbf{B}|$$

Suppose $\mathbf{A} = \mathbf{X}\mathbf{E}\mathbf{X}^{-1}$, where $\mathbf{E} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is the eigenvalue matrix and \mathbf{X} is the eigenvector matrix. Then,

$$\mathbf{C}\mathbf{A}^k\mathbf{B} = \mathbf{C}\mathbf{X}\mathbf{E}^k\mathbf{X}^{-1}\mathbf{B} = \sum_{l=1}^n R_l \lambda_l^k$$

Bound on truncation error

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

$$\mathbf{M} := \sum_{l=1}^n \frac{|R_l|}{1 - |\lambda_l|} \frac{|\lambda_l|}{\rho(\mathbf{A})}$$

Step 1. Bound on truncation order

Lower bound on truncation order

$$N \geq \left\lceil \frac{\log \frac{\varepsilon_1}{m}}{\log \rho(\mathbf{A})} \right\rceil$$
$$\mathbf{M} := \sum_{l=1}^n \frac{|\mathbf{R}_l|}{1 - |\lambda_l|} \frac{|\lambda_l|}{\rho(\mathbf{A})}$$

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

Step 1. Bound on truncation order

Lower bound on truncation order

$$N \geq \left\lceil \frac{\log \frac{\varepsilon_1}{m}}{\log \rho(\mathbf{A})} \right\rceil$$
$$\mathbf{M} := \sum_{l=1}^n \frac{|\mathbf{R}_l|}{1 - |\lambda_l|} \frac{|\lambda_l|}{\rho(\mathbf{A})}$$

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

Reliable evaluation

Interval Arithmetic and Rump's Theory of Verified Inclusions are used to determine a rigorous bound of N .

Step 2. "Diagonalization" of matrix **A**

$$\mathbf{T} := \mathbf{V}^{-1}\mathbf{A}\mathbf{V} - \mathbf{\Delta}_2$$

Step 2. "Diagonalization" of matrix **A**

$$\mathbf{T} := \mathbf{V}^{-1}\mathbf{AV} - \mathbf{\Delta}_2$$

- **V** is some approximation on **X**
- **Δ₂** represents the element-by-element errors due to the two matrix multiplications and the inversion of matrix **V**

Step 2. "Diagonalization" of matrix **A**

$$\mathbf{T} := \mathbf{V}^{-1}\mathbf{A}\mathbf{V} - \mathbf{\Delta}_2$$

- **V** is some approximation on **X**
- $\mathbf{\Delta}_2$ represents the element-by-element errors due to the two matrix multiplications and the inversion of matrix **V**
- **T** diagonal in dominant with very small other elements
- $\|\mathbf{T}\|_2 \leq 1$

Step 2. "Diagonalization" of matrix \mathbf{A}

$$\begin{aligned}\mathbf{T} &:= \mathbf{V}^{-1}\mathbf{A}\mathbf{V} - \mathbf{\Delta}_2 \\ \mathbf{A}^k &= \mathbf{V}(\mathbf{T} + \mathbf{\Delta}_2)^k \mathbf{V}^{-1}\end{aligned}$$

The error of substitution of \mathbf{A} by $\mathbf{V}\mathbf{T}\mathbf{V}^{-1}$:

$$\sqrt{n}(N+1)(N+2) \|\mathbf{\Delta}_2\|_F \|\mathbf{C}\mathbf{V}\|_F \|\mathbf{V}^{-1}\mathbf{B}\|_F$$

Step 2. "Diagonalization" of matrix \mathbf{A}

$$\begin{aligned}\mathbf{T} &:= \mathbf{V}^{-1} \mathbf{A} \mathbf{V} - \mathbf{\Delta}_2 \\ \mathbf{A}^k &= \mathbf{V} (\mathbf{T} + \mathbf{\Delta}_2)^k \mathbf{V}^{-1}\end{aligned}$$

The error of substitution of \mathbf{A} by $\mathbf{V} \mathbf{T} \mathbf{V}^{-1}$:

$$\sqrt{n}(N+1)(N+2) \|\mathbf{\Delta}_2\|_F \|\mathbf{C} \mathbf{V}\|_F \|\mathbf{V}^{-1} \mathbf{B}\|_F \stackrel{!}{\leq} \varepsilon_2$$

Step 2. "Diagonalization" of matrix \mathbf{A}

$$\begin{aligned}\mathbf{T} &:= \mathbf{V}^{-1} \mathbf{A} \mathbf{V} - \mathbf{\Delta}_2 \\ \mathbf{A}^k &= \mathbf{V}(\mathbf{T} + \mathbf{\Delta}_2)^k \mathbf{V}^{-1}\end{aligned}$$

The error of substitution of \mathbf{A} by $\mathbf{V} \mathbf{T} \mathbf{V}^{-1}$:

$$\sqrt{n}(N+1)(N+2) \|\mathbf{\Delta}_2\|_F \|\mathbf{C} \mathbf{V}\|_F \|\mathbf{V}^{-1} \mathbf{B}\|_F \stackrel{!}{\leq} \varepsilon_2$$

A condition on the error-matrix $\mathbf{\Delta}_2$:

$$\|\mathbf{\Delta}_2\|_F \leq \frac{1}{\sqrt{n}(N+1)(N+2)} \frac{\varepsilon_2}{\|\mathbf{C} \mathbf{V}\|_F \|\mathbf{V}^{-1} \mathbf{B}\|_F}$$

Step 3. Computing products \mathbf{C}' and \mathbf{B}'

$$\mathbf{C}' := \mathbf{C}\mathbf{V} + \mathbf{\Delta}_{3_C}$$

$$\mathbf{B}' := \mathbf{V}^{-1}\mathbf{B} + \mathbf{\Delta}_{3_B}$$

where $\mathbf{\Delta}_{3_C} \in \mathbb{C}^{p \times n}$ and $\mathbf{\Delta}_{3_B} \in \mathbb{C}^{n \times q}$ are error-matrices.

Bound on the multiplication errors $\mathbf{\Delta}_{3_C}$ and $\mathbf{\Delta}_{3_B}$:

$$\|\mathbf{\Delta}_{3_C}\|_F \leq \frac{1}{3\sqrt{n}} \cdot \frac{1}{N+1} \frac{\varepsilon_3}{\|\mathbf{C}'\|_F}$$

$$\|\mathbf{\Delta}_{3_B}\|_F \leq \frac{1}{3\sqrt{n}} \cdot \frac{1}{N+1} \frac{\varepsilon_3}{\|\mathbf{B}'\|_F}.$$

Step 4. Powering \mathbf{T}

$$\mathbf{P}_k := \mathbf{T}^k - \mathbf{\Delta}_{4_k}$$

$\mathbf{\Delta}_{4_k} \in \mathbb{C}^{n \times n}$ error-matrix on matrix powers, including error propagation from the first to the last power.

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

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

Bound on the error-matrix $\mathbf{\Gamma}_k$

$$\|\mathbf{\Gamma}_k\|_F \leq \frac{1}{\sqrt{n}} \cdot \frac{1}{N-1} \cdot \frac{1}{N+1} \cdot \frac{\varepsilon_4}{\|\mathbf{C}'\|_F \|\mathbf{B}'\|_F}$$

Step 5. Computing L_k

$$L_k := C' P_k B' + \Delta_{5_k},$$

where $\Delta_{5_k} \in \mathbb{C}^{p \times q}$ is the matrix of element-by-element errors for the two matrix multiplications.

Bound on the error-matrix Δ_{5_k}

$$|\Delta_{5_k}| \leq \frac{1}{N+1} \cdot \varepsilon_5.$$

Step 6. Summation

$$S_N = |\mathbf{D}| + \sum_{l=0}^N |\mathbf{L}_l| + \mathbf{\Delta}_6,$$

where the error-matrix $\mathbf{\Delta}_6 \in \mathbb{C}^{p \times q}$ represents the error of $N + 1$ absolute value accumulations.

Bound on the error matrix $\mathbf{\Delta}_{6_k}$

$$\mathbf{\Delta}_{6_k} \leq \frac{1}{N} \varepsilon_6, \quad k = 1 \dots N$$