# Towards reliable code generation for filters

Anastasia Volkova
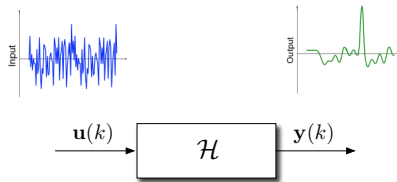
PhD Director: Jean Claude Bajard
PhD Advisors: Thibault Hilaire, Christoph Lauter
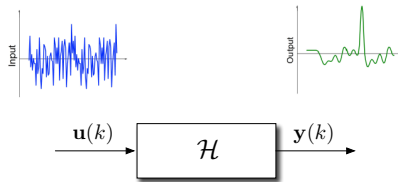
May 11, 2016

SORBONNE UNIVERSITÉS    UPMC SORBONNE UNIVERSITÉS    Cnrs    LIP6
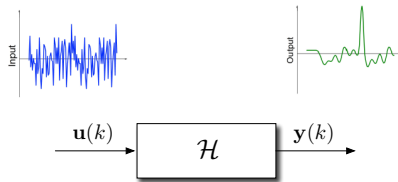
# Context: digital filters

# Context: digital filters



On the one hand

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

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

# Context: digital filters



On the one hand
- LTI filter with Infinite Impulse Response
- Its transfer function:

On the other hand
- Hardware or Software target
- Implementation in Fixed-Point Arithmetic

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

# Context: implementation of LTI filters

# Context: implementation of LTI filters



$$\mathbf{H(z)} = \mathbf{C}(\mathbf{zI} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$$

$$\mathbf{H(z)} = \frac{\sum_{n=1}^{N_b} b_n z^{-(n-1)}}{1 + \sum_{n=1}^{N_a} a_n z^{-n}}$$

- 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.

# Filter-to-code generator



Figure: Automatic filter generator flow.

Before this thesis:

  Stage 1: analytical filter realization representation

  Stage 2: filter quality measures

  Stage 3: fixed-point algorithm (naive approach, computational errors not taken into account)

  Stage 4: Fixed-Point Code Generator
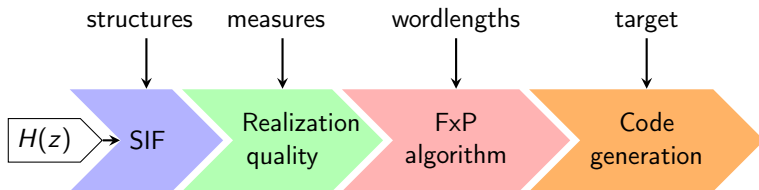
# Filter-to-code generator



Figure: Automatic filter generator flow.

During this thesis:

  Stage 1: analytical filter realization representation

  Stage 2: filter quality measures

  Stage 3: fixed-point algorithm (~~naive~~ **rigorous** approach,
  computational errors ~~not~~ **taken into account**)

  Stage 4: Fixed-Point Code Generator

Representing Lattice Wave Digital Filters
with Specialized Implicit Framework[1].

[1]A.V. et al., "Fixed-Point Implementation of Lattice Wave Digital Filter: Comparison and Error Analysis", in 23rd European Signal Processing Conference, EUSIPCO 2015

Reliable implementation of digital filters
in Fixed-Point Arithmetic

# LTI filters

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

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

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

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

# 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:

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

Determine: the Fixed-Point Formats s.t.

- the least MSBs
- no overflows occur
  - $\rightsquigarrow$ must take into account computational errors

# Reliable Fixed-Point implementation

Input:

- $\mathcal{H} = (\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}, \boldsymbol{D})$
- 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 interval[2]

[2]A.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



Input $\boldsymbol{u}(k)$

$\forall k, \quad |\boldsymbol{u}(k)| \leq \bar{\boldsymbol{u}}$

# Basic brick: the Worst-Case Peak Gain theorem



Input $\boldsymbol{u}(k)$

$\forall k, \quad |\boldsymbol{u}(k)| \leq \bar{\boldsymbol{u}}$

Amplitude

Time

$\boldsymbol{u}(k)$ → $\mathcal{H}$ → $\boldsymbol{y}(k)$

amplification/attenuation

# Basic brick: the Worst-Case Peak Gain theorem

Input $\boldsymbol{u}(k)$

$\forall k, \quad |\boldsymbol{u}(k)| \leq \bar{\boldsymbol{u}}$



Output $\boldsymbol{y}(k)$



$\boldsymbol{u}(k) \quad \boxed{\mathcal{H}} \quad \boldsymbol{y}(k)$

amplification/attenuation

# Basic brick: the Worst-Case Peak Gain theorem



Input $\boldsymbol{u}(k)$

$\forall k, \quad |\boldsymbol{u}(k)| \leq \bar{\boldsymbol{u}}$

Output $\boldsymbol{y}(k)$

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

$\boldsymbol{u}(k)$ → $\mathcal{H}$ → $\boldsymbol{y}(k)$

amplification/attenuation

Worst-Case Peak Gain

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

# Computing the Worst-Case Peak Gain

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

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

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

# Truncation

$$\sum_{k=0}^{\infty} \left| \boldsymbol{C}\boldsymbol{A}^k\boldsymbol{B} \right| \quad \longrightarrow \quad \sum_{k=0}^{N} \left| \boldsymbol{C}\boldsymbol{A}^k\boldsymbol{B} \right|$$

# Truncation

$$\left| \sum_{k=0}^{\infty} \left| CA^k B \right| \quad - \quad \sum_{k=0}^{N} \left| CA^k B \right| \right| \leq \varepsilon_1$$

Compute an approximate lower bound on truncation order $N$ such that the truncation error is smaller than $\varepsilon_1$.

## Lower bound on truncation order N

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

where

$$\lambda - \text{eigenvalues of matrix } A$$

$$R_l - l^{\text{th}} \text{residue matrix computed out of } C, B, \lambda$$

# Powering

$$\sum_{k=0}^{N} \left| \boldsymbol{C} \boldsymbol{A}^k \boldsymbol{B} \right|$$

$$\sum_{k=0}^{N} \left| \boldsymbol{C} \boldsymbol{A}^k \boldsymbol{B} \right|$$

 $\times$  $=$   cancellation

# Powering

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



cancellation

less cancellation

$$\sum_{k=0}^{N} \left| \boldsymbol{C}\boldsymbol{A}^k\boldsymbol{B} \right|$$

 $\times$  $=$  cancellation

 $\times$  $=$  less cancellation

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

$$\sum_{k=0}^{N} \left| \boldsymbol{C} \boldsymbol{A}^k \boldsymbol{B} \right|$$

 $\times$  $=$  cancellation

 $\times$  $=$  less cancellation

$$\boldsymbol{A} = \boldsymbol{X} \boldsymbol{E} \boldsymbol{X}^{-1} \qquad\qquad \boldsymbol{V} \approx \boldsymbol{X} \text{ and } \boldsymbol{T} \approx \boldsymbol{E}$$

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

 $\times$  $=$    cancellation

 $\times$  $=$    less cancellation

$$\mathbf{A} = \mathbf{X} \mathbf{E} \mathbf{X}^{-1} \qquad\qquad \mathbf{V} \approx \mathbf{X} \text{ and } \mathbf{T} \approx \mathbf{E}$$

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

$$\downarrow$$

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

# Powering

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

Given matrix $V$ compute $T$ such that the error of substitution of the product $V T^k V^{-1}$ instead of $A^k$ is less than $\varepsilon_2$.

# Further steps

$$\left| \sum_{k=0}^{N} |CA^kB| \quad - \quad \sum_{k=0}^{N} |CVT^kV^{-1}B| \right| \leq \varepsilon_2$$

Apply the same approach for the other steps:

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

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

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

$$\left| \sum_{k=0}^{N} |L_k| - S_N \right| \leq \varepsilon_6$$

# Further steps

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

Apply the same approach for the other steps:

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

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

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

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

We can determine the output interval of a filter with arbitrary precision.

Determining the Fixed-Point Formats[3]

---

[3]A.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

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

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

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

# Determining the Fixed-Point Formats

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

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

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

We need to find the least $\boldsymbol{m}_y$ such that

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

# Determining the Fixed-Point Formats

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

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

$$\forall k, \quad |\boldsymbol{y}_i(k)| \le \left(\langle\langle\mathcal{H}\rangle\rangle\,\bar{\boldsymbol{u}}\right)_i.$$

We need to find the least $\boldsymbol{m}_y$ such that

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

We have shown that $\boldsymbol{m}_y$ can be computed with

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

# Taking the quantization errors into account

The exact filter $\mathcal{H}$ is:

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

# Taking the quantization errors into account

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

$$\mathcal{H}^\diamond \left\{ \begin{array}{rcl} \boldsymbol{x}^\diamond(k+1) & = & \diamond_{m_x}(\boldsymbol{A}\boldsymbol{x}^\diamond(k) + \boldsymbol{B}\boldsymbol{u}(k)) \\ \boldsymbol{y}^\diamond(k) & = & \diamond_{m_y}(\boldsymbol{C}\boldsymbol{x}^\diamond(k) + \boldsymbol{D}\boldsymbol{u}(k)) \end{array} \right.$$

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:

$$\mathcal{H}^{\diamond} \left\{ \begin{array}{rcl} \boldsymbol{x}^{\diamond}(k+1) & = & \boldsymbol{A}\boldsymbol{x}^{\diamond}(k) + \boldsymbol{B}\boldsymbol{u}(k) + \varepsilon_x(k) \\ \boldsymbol{y}^{\diamond}(k) & = & \boldsymbol{C}\boldsymbol{x}^{\diamond}(k) + \boldsymbol{D}\boldsymbol{u}(k) + \varepsilon_y(k) \end{array} \right.$$

with

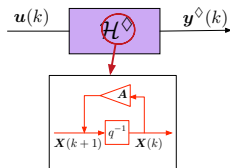$$|\varepsilon_x(k)| \leq 2^{m_x - w_x + 1} \quad \text{and} \quad |\varepsilon_y(k)| \leq 2^{m_y - w_y + 1}.$$

# Taking the quantization errors into account

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

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

with

$$|\varepsilon_x(k)| \leq 2^{\boldsymbol{m}_x - \boldsymbol{w}_x + 1} \quad \text{and} \quad |\varepsilon_y(k)| \leq 2^{\boldsymbol{m}_y - \boldsymbol{w}_y + 1}.$$
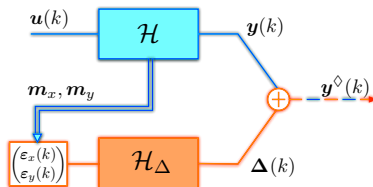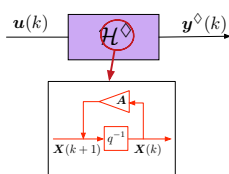
# Taking the quantization errors into account

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

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

with

$$|\varepsilon_x(k)| \le 2^{\boldsymbol{m}_x - \boldsymbol{w}_x + 1} \quad \text{and} \quad |\varepsilon_y(k)| \le 2^{\boldsymbol{m}_y - \boldsymbol{w}_y + 1}.$$

# 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
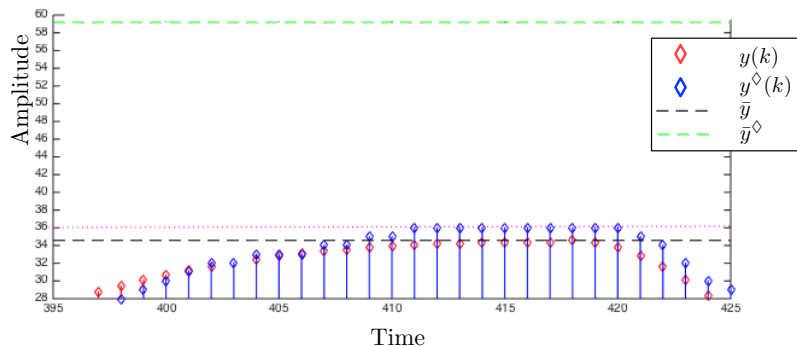
# Numerical examples



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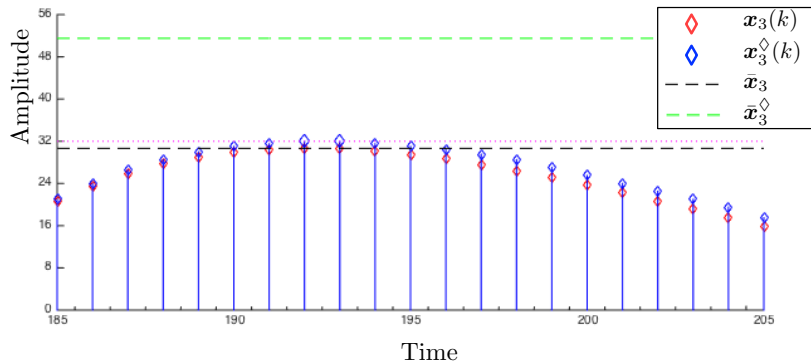
# Numerical examples



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

# Conclusion

- Represented a new structure in the SIF formalism
- Provided reliable evaluation of the WCPG measure
- Applied the WCPG measure to determine the FxPF that guarantee no overflow

Publications:

"Reliable Evaluation of the Worst-Case Peak Gain Matrix in Multiple Precision", ARITH22

"Determining Fixed-Point Formats for a Digital Filter Implementation using the Worst-Case Peak Gain Measure", Asilomar 49

"Fixed-Point Implementation of Lattice Wave Digital Filter: Comparison and Error Analysis", EUSIPCO 2015

# Perspectives

For filter implementation...

- Solve the off-by-one problem for the MSBs
- Get a deeper insight on the behavior of rounding errors
  - ⤳ determine the probability distribution function
- Plug all the stages of the generator into optimization routines
- Accuracy of the algorithms for the design of IIR filters
  - ⤳ take coefficient quantization errors into account in the filter error analysis
  - ⤳ consider coefficients as intervals and make necessary adaptations in the algorithms within the filter

# Perspectives

For filter implementation...

- Solve the off-by-one problem for the MSBs
- Get a deeper insight on the behavior of rounding errors
  - $\rightsquigarrow$ determine the probability distribution function
- Plug all the stages of the generator into optimization routines
- Accuracy of the algorithms for the design of IIR filters
  - $\rightsquigarrow$ take coefficient quantization errors into account in the filter error analysis
  - $\rightsquigarrow$ consider coefficients as intervals and make necessary adaptations in the algorithms within the filter

Mathematical function implementation...

- Draw parallels between filter and elementary function implementation.

# Merci !

## Thank you!

### Спасибо!

#### Дякую!