Fixed-Point implementation of Lattice Wave Digital Filter: comparison and error analysis

Anastasia Volkova, Thibault Hilaire

Sorbonne Universités, University Pierre and Marie Curie, LIP6, Paris, France

EUSIPCO 15 September 2, 2015









Need to deal with

- Discretize functions and coefficients
 - parametric errors
 - computational errors
- Implementation under constraints
 - software implementation
 - hardware implementation

Different filter structures:

- Direct Form I, Direct Form II
- State-space
- Wave, Lattice Wave, ...
- ρ -operator: ρ DFIIt, ρ Modal, ρ State-space...
- LGS, LCW, etc.

Problem:

They are equivalent in *infinite* precision but no more in *finite* precision. The finite precision degradation depends on the realization.

- Represent various realizations
- Evaluate finite precision degradation
- Find an optimal realization Tradeoff:
 - Error
 - Quality
 - Power consumption
 - Speed
 - etc.

- Represent various realizations (in an easy way)
- Evaluate finite precision degradation
- Find an optimal realization Tradeoff:
 - Error
 - Quality
 - Power consumption
 - Speed
 - etc.

- Represent various realizations (in an easy way)
- Evaluate finite precision degradation (a priori/a posteriori)
- Find an optimal realization Tradeoff:
 - Error
 - Quality
 - Power consumption
 - Speed
 - etc.

- Represent various realizations (in an easy way)
- Evaluate finite precision degradation (a priori/a posteriori)
- Find an optimal realization (need to compare realizations) Tradeoff:
 - Error
 - Quality
 - Power consumption
 - Speed
 - etc.

Given transfer function and a target, we want:

- Represent various realizations (in an easy way)
- Evaluate finite precision degradation (a priori/a posteriori)
- Find an optimal realization (need to compare realizations) Tradeoff:
 - Error
 - Quality
 - Power consumption
 - Speed
 - etc.

Specialized Implicit Framework (SIF)

Outline



- 2 Specialized Implicit Framework
- **3** Lattice Wave Digital Filters
- 4 LWDF-to-SIF convertion
- **(5)** Example and comparison



SIF is:

- Macroscopic description
- Based on state-space
- Explicit all the computations and their order
- Any DFG can be transformed to this form
- Analytical derivation of measures

$$\mathcal{H} \begin{cases} \boldsymbol{J}\boldsymbol{t}(k+1) = & \boldsymbol{M}\boldsymbol{x}(k) + \boldsymbol{N}\boldsymbol{u}(k) \\ \boldsymbol{x}(k+1) = \boldsymbol{K}\boldsymbol{t}(k+1) + \boldsymbol{P}\boldsymbol{x}(k) + \boldsymbol{Q}\boldsymbol{u}(k) \\ \boldsymbol{y}(k) = \boldsymbol{L}\boldsymbol{t}(k+1) + \boldsymbol{R}\boldsymbol{x}(k) + \boldsymbol{S}\boldsymbol{u}(k) \end{cases}$$

Denote \mathbf{Z} the matrix containing all the coefficients

$$Z riangleq egin{pmatrix} -J & M & N \ K & P & Q \ L & R & S \end{pmatrix}$$

SIF: measures

Measures

- a priori measures
 - transfer function sensitivity (based on $\frac{\partial H}{\partial Z}$)
 - $\longrightarrow\,$ stochastic measure, takes into account coefficient wordlengths
 - poles or zeros sensitivity (e.g based on $\frac{\partial |\lambda_i|}{\partial Z}$ for a pole λ_i)

 \longrightarrow stochastic measure, takes into account coefficient wordlengths

- RNG, ...
- *a posteriori* measures
 - Signal to Quantization Noise Ratio
 - output error

SIF: measures

Measures

- *a priori* measures
 - transfer function sensitivity (based on $\frac{\partial H}{\partial Z}$)
 - $\longrightarrow\,$ stochastic measure, takes into account coefficient wordlengths
 - poles or zeros sensitivity (e.g based on $\frac{\partial |\lambda_i|}{\partial Z}$ for a pole λ_i)

 $\longrightarrow\,$ stochastic measure, takes into account coefficient wordlengths

- RNG, ...
- a posteriori measures
 - Signal to Quantization Noise Ratio
 - output error

WCPG theorem

Let $\mathcal{H} = \{ \boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}, \boldsymbol{D} \}$ be a BIBO stable MIMO state-space. If $\forall k \ \boldsymbol{u}(k) \leq \bar{\boldsymbol{u}}$ component-wisely, then component-wisely $\forall k \ \boldsymbol{y}(k) \leq \langle \! \langle H \rangle \! \rangle \, \bar{\boldsymbol{u}},$

where $\langle\!\langle H \rangle\!\rangle$ is the Worst-Case Peak Gain matrix of the system and can be computed as

$$\langle\!\langle H
angle\!
angle := |oldsymbol{D}| + \sum_{k=0}^{\infty} \left| oldsymbol{C} oldsymbol{A}^k oldsymbol{B} \right|.$$

Note: we can compute $\langle\!\langle H \rangle\!\rangle$ in arbitrary precision.

Exact filter:

$$\mathcal{H} \left\{ \begin{array}{l} \boldsymbol{Jt} \ (k+1) = & \boldsymbol{Mx} \ (k) + \boldsymbol{Nu}(k) \\ \boldsymbol{x} \ (k+1) = \boldsymbol{Kt} \ (k+1) + \boldsymbol{Px} \ (k) + \boldsymbol{Qu}(k) \\ \boldsymbol{y} \ (k) = \boldsymbol{Lt} \ (k+1) + \boldsymbol{Rx} \ (k) + \boldsymbol{Su}(k) \end{array} \right.$$

Summary

SIF: the rigorous filter error bound

Implemented filter:

$$\mathcal{H}^* \left\{ \begin{array}{ll} \boldsymbol{J} \boldsymbol{t}^*(k+1) = & \boldsymbol{M} \boldsymbol{x}^*(k) + \boldsymbol{N} u(k) + \boldsymbol{\varepsilon}_{\boldsymbol{t}}(k) \\ \boldsymbol{x}^*(k+1) = \boldsymbol{K} \boldsymbol{t}^*(k+1) + \boldsymbol{P} \boldsymbol{x}^*(k) + \boldsymbol{Q} u(k) + \boldsymbol{\varepsilon}_{\boldsymbol{x}}(k) \\ \boldsymbol{y}^*(k) = \boldsymbol{L} \boldsymbol{t}^*(k+1) + \boldsymbol{R} \boldsymbol{x}^*(k) + \boldsymbol{S} u(k) + \boldsymbol{\varepsilon}_{\boldsymbol{y}}(k) \end{array} \right.$$

where $\boldsymbol{\varepsilon}_t(k)$, $\boldsymbol{\varepsilon}_x(k)$ and $\boldsymbol{\varepsilon}_y(k)$ are the computational errors.

Implemented filter:

$$\mathcal{H}^* \left\{ \begin{array}{ll} \boldsymbol{Jt}^*(k+1) = & \boldsymbol{Mx}^*(k) + \boldsymbol{Nu}(k) + \boldsymbol{\varepsilon_t}(k) \\ \boldsymbol{x}^*(k+1) = \boldsymbol{Kt}^*(k+1) + \boldsymbol{Px}^*(k) + \boldsymbol{Qu}(k) + \boldsymbol{\varepsilon_x}(k) \\ \boldsymbol{y}^*(k) = \boldsymbol{Lt}^*(k+1) + \boldsymbol{Rx}^*(k) + \boldsymbol{Su}(k) + \boldsymbol{\varepsilon_y}(k) \end{array} \right.$$

where $\varepsilon_t(k)$, $\varepsilon_x(k)$ and $\varepsilon_y(k)$ are the computational errors. The output error

$$\Delta \boldsymbol{y}(k) \triangleq \boldsymbol{y}^*(k) - \boldsymbol{y}(k)$$

can be seen as the output of a MIMO filter $\mathcal{H}_{\varepsilon}$.

Implemented filter:

$$\mathcal{H}^* \left\{ \begin{array}{ll} \boldsymbol{Jt}^*(k+1) = & \boldsymbol{Mx}^*(k) + \boldsymbol{Nu}(k) + \boldsymbol{\varepsilon_t}(k) \\ \boldsymbol{x}^*(k+1) = \boldsymbol{Kt}^*(k+1) + \boldsymbol{Px}^*(k) + \boldsymbol{Qu}(k) + \boldsymbol{\varepsilon_x}(k) \\ \boldsymbol{y}^*(k) = \boldsymbol{Lt}^*(k+1) + \boldsymbol{Rx}^*(k) + \boldsymbol{Su}(k) + \boldsymbol{\varepsilon_y}(k) \end{array} \right.$$

where $\varepsilon_t(k)$, $\varepsilon_x(k)$ and $\varepsilon_y(k)$ are the computational errors. The output error

$$\Delta \boldsymbol{y}(k) \triangleq \boldsymbol{y}^*(k) - \boldsymbol{y}(k)$$

can be seen as the output of a MIMO filter $\mathcal{H}_{\varepsilon}$.

$$\underbrace{\begin{array}{c} \boldsymbol{u}(k) \\ \boldsymbol{\varepsilon}(k) \\ \boldsymbol{\varepsilon}(k)$$

Implemented filter:

$$\mathcal{H}^* \left\{ \begin{array}{ll} \boldsymbol{Jt}^*(k+1) = & \boldsymbol{Mx}^*(k) + \boldsymbol{Nu}(k) + \boldsymbol{\varepsilon_t}(k) \\ \boldsymbol{x}^*(k+1) = \boldsymbol{Kt}^*(k+1) + \boldsymbol{Px}^*(k) + \boldsymbol{Qu}(k) + \boldsymbol{\varepsilon_x}(k) \\ \boldsymbol{y}^*(k) = \boldsymbol{Lt}^*(k+1) + \boldsymbol{Rx}^*(k) + \boldsymbol{Su}(k) + \boldsymbol{\varepsilon_y}(k) \end{array} \right.$$

where $\varepsilon_t(k)$, $\varepsilon_x(k)$ and $\varepsilon_y(k)$ are the computational errors. The output error

$$\Delta \boldsymbol{y}(k) \triangleq \boldsymbol{y}^*(k) - \boldsymbol{y}(k)$$

can be seen as the output of a MIMO filter $\mathcal{H}_{\varepsilon}$.

$$\underbrace{u(k)}_{\varepsilon(k)} \xrightarrow{\mathcal{H}} \underbrace{y(k)}_{\psi(k)} \xrightarrow{\psi^*(k)} \underbrace{\mathcal{H}}_{\varepsilon} \underbrace{\Delta y(k)}_{\psi^*(k)} \xrightarrow{\psi^*(k)} \underbrace{\mathcal{H}}_{\varepsilon}$$

WCPG theorem on $\mathcal{H}_{\varepsilon}$ gives the output error interval: $\Delta \boldsymbol{y}(k) \leq \langle \langle H_{\varepsilon} \rangle \rangle \, \bar{\varepsilon}$

SIF: code generation

WCPG theorem gives a *rigorous* way to compute Most Significant Bit:

$$\boldsymbol{m}_y = \left\lfloor \log_2\left(\langle\!\langle \boldsymbol{H} \rangle\!\rangle \, \bar{\boldsymbol{u}}\right) \right\rfloor + 1$$

Equivalent technique: WCPG-scaling, it guarantees that no overflows occur.

SIF: code generation

WCPG theorem gives a *rigorous* way to compute Most Significant Bit:

$$\boldsymbol{m}_y = \left\lfloor \log_2\left(\left\langle\!\left\langle \boldsymbol{H} \right\rangle\!\right\rangle \, \bar{\boldsymbol{u}} \right) \right\rfloor + 1$$

Equivalent technique: WCPG-scaling, it guarantees that no overflows occur.

Fixed Point Code Generator (FiPoGen)

Given wordlength and evaluation scheme

• Generates bit-accurate fixed-point algorithms

Given only evaluation scheme

- Optimizes the wordlength under certain criteria (e.g. area)
- Generates bit-accurate fixed-point algorithms

SIF: from transfer function to Fixed-Point code



SIF: from transfer function to Fixed-Point code



SIF: from transfer function to Fixed-Point code



Summary

Lattice Wave Digital Filters



Lattice Wave Digital Filters



Motivation SIF LWDF LWDF-to-SIF Example and comparison Sum

Lattice Wave Digital Filters

Two-port adaptor: Richard's structures



13/22

Motivation SIF LWDF LWDF-to-SIF Example and comparison Summ

Lattice Wave Digital Filters

Two-port adaptor: Richard's structures



13/22

Lattice Wave Digital Filters

Positive sides

- parallelizable
- modular, convenient for VLSI
- often referred to as *stable*

Drawbacks

- Studies of Fixed-Point implementation include complicated infinite-precision optimization
- Comparison is difficult

Objectives

- Represent LWDF in terms of SIF
- Perform *rigorous* error analysis
- Instantly compare with other structures



• SIF representation for subsystems of Type A and Type B



- SIF representation for subsystems of Type A and Type B
- Cascade subsystems into stages



- SIF representation for subsystems of Type A and Type B
- Cascade subsystems into stages
- O Cascade stages into branches



- SIF representation for subsystems of Type A and Type B
- Cascade subsystems into stages
- O Cascade stages into branches
- Cascade branches into low/high pass filter



LWDF-to-SIF convertion: example

Convert DFGs of two adaptors into SIFs:



LWDF-to-SIF convertion: example

Convert DFGs of two adaptors into SIFs:



$$\boldsymbol{Z}_{A} \triangleq \begin{pmatrix} -\boldsymbol{J}_{A} \mid \boldsymbol{M}_{A} \mid \boldsymbol{N}_{A} \\ \boldsymbol{K}_{A} \mid \boldsymbol{P}_{A} \mid \boldsymbol{Q}_{A} \\ \boldsymbol{L}_{A} \mid \boldsymbol{R}_{A} \mid \boldsymbol{S}_{A} \end{pmatrix} = \begin{pmatrix} -1 & 0 & 0 & -1 & 1 \\ \alpha & -1 & 0 & 0 & -1 \\ 0 & -1 & 0 & 0 & 0 \\ \hline 0 & -1 & 1 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$
$$\boldsymbol{Z}_{B} \triangleq \begin{pmatrix} -\boldsymbol{J}_{B} \mid \boldsymbol{M}_{B} \mid \boldsymbol{N}_{B} \\ \boldsymbol{K}_{B} \mid \boldsymbol{P}_{B} \mid \boldsymbol{Q}_{B} \\ \boldsymbol{L}_{B} \mid \boldsymbol{R}_{B} \mid \boldsymbol{S}_{B} \end{pmatrix} = \begin{pmatrix} -1 & 0 & 1 & -1 \\ \alpha & -1 & 1 & 0 \\ \hline 0 & -1 & 0 & 0 \\ \hline -1 & 1 & 0 & 0 \end{pmatrix}$$

 α

1.

Cascading SIFs:

$$Z_1 = \{J_1, K_1, \dots, S_1\}$$

 $Z_2 = \{J_2, K_2, \dots, S_2\}$
Then,
 $\mathbf{u}_{A(k)} \underbrace{\mathcal{H}_A}_{\mathbf{y}_A(k) = \mathbf{u}_B(k)} \underbrace{\mathcal{H}_B}_{\mathbf{y}_B(k)}$

$$egin{aligned} egin{aligned} egin{aligned} egin{aligned} egin{aligned} egin{aligned} -egin{aligned} -egin{aligned} egin{aligned} egin{aligned\\ egin{aligned} e$$

Note: matrix Z is extremely sparse.

Example and comparison

Reference filter: low-pass 5^{th} order Butterworth filter with cutoff frequency 0.1. Structures for the comparison:

- LWDF
- state-space
- ρ -Direct Form II transposed
- Direct Form I

Normalized (*i.e.* all coefficients have the same wordlength) measures:

- transfer function error: $\bar{\sigma}_{\Delta H}^2$
- pole error: $\bar{\sigma}^2_{\Delta|\boldsymbol{\lambda}|}$
- output error: $\overline{\Delta_y}$

Example and comparison

LWDF, Z is 22×22



DFI, Z is 12×12



State-space, Z is 12×12



 ρ DFIIt, Z is 12×12



Example and comparison

Realization	$\operatorname{size}(\mathbf{Z})$	coeff.	$\bar{\sigma}^2_{\Delta H}$	$ar{\sigma}^2_{\Delta oldsymbol{\lambda} }$	$\overline{\Delta_y}$
LWDF	22×22	5	0. 3151	0.56	122.9
state-space	6×6	36	1.15	5.75	23.33
$ ho \mathrm{DFIIt}$	11×11	11	0.09	0.45	94.3
DFI	12×12	11	1.42e + 6	-	7.961

Conclusion and perspectives

Conclusion:

- LWDF converted to SIF
- Normalized sensitivity and output error measures applied
- Comparison with several popular structures presented

Perspectives:

- Use VHDL code generator (FloPoCo) to compare hardware implementations
- Apply ρ -operator to LWDF

Thank you! Questions?