Thibault HILAIRE thibault.hilaire@irisa.fr

Bit Accurate Roundoff Noise Analysis of Fixed-Point Linear Controllers



CAIRN project IRISA/INRIA France IEEE Multi-Conference on Systems and Control San Antonio, Texas 4 September 2008

= /



Implementation of Linear Time Invariant controllers/filtersFinite Word Length context (fixed-point)

< 口 > < 同 >



- Implementation of Linear Time Invariant controllers/filters
- Finite Word Length context (fixed-point)

< 一型



- Implementation of Linear Time Invariant controllers/filters
- Finite Word Length context (fixed-point)

Motivation

- Analysis (accurately) the roundoff noise errors in the implementation
- Compare various realizations and find an optimal one

< 17 ▶



- Implementation of Linear Time Invariant controllers/filters
- Finite Word Length context (fixed-point)

The roundoff will depend on

- the algorithmic relation to compute the output(s) from the input(s)
- the way the computations are implemented (wordlength, roundoff, etc.)

Outline

- Implicit state-space framework
- Roundoff noise analysis
- Fixed-point implementation schemes
- Optimal design
- Sonclusion



=

The need of a unifying framework

Various implementation forms have to be taken into consideration:

- shift-realizations
- δ -realizations
- observer-state-feedback
- direct form I or II
- cascade or parallel realizations
- etc...

Implicit specialized state-space form

$$\begin{pmatrix}J & 0 & 0\\ -K & I & 0\\ -L & 0 & I\end{pmatrix}\begin{pmatrix}T_{k+1}\\ X_{k+1}\\ Y_k\end{pmatrix} = \begin{pmatrix}0 & M & N\\ 0 & P & Q\\ 0 & R & S\end{pmatrix}\begin{pmatrix}T_k\\ X_k\\ U_k\end{pmatrix}$$

T. Hilaire

Implicit specialized state-space form

$$\begin{pmatrix} J & 0 & 0 \\ -K & I & 0 \\ -L & 0 & I \end{pmatrix} \begin{pmatrix} T_{k+1} \\ X_{k+1} \\ Y_k \end{pmatrix} = \begin{pmatrix} 0 & M & N \\ 0 & P & Q \\ 0 & R & S \end{pmatrix} \begin{pmatrix} T_k \\ X_k \\ U_k \end{pmatrix}$$

It corresponds to:

1 $J.T_{k+1} = M.X_k + N.U_k$

$$X_{k+1} = K.T_{k+1} + P.X_k + Q.U_k$$

$$Y_{k} = L.T_{k+1} + R.X_{k} + S.U_{k}$$

Intermediate variables computation

< 17 ▶

Implicit specialized state-space form

$$\begin{pmatrix} J & 0 & 0 \\ -K & I & 0 \\ -L & 0 & I \end{pmatrix} \begin{pmatrix} T_{k+1} \\ X_{k+1} \\ Y_k \end{pmatrix} = \begin{pmatrix} 0 & M & N \\ 0 & P & Q \\ 0 & R & S \end{pmatrix} \begin{pmatrix} T_k \\ X_k \\ U_k \end{pmatrix}$$

It corresponds to:

$$J.T_{k+1} = M.X_k + N.U_k$$

$$X_{k+1} = K.T_{k+1} + P.X_k + Q.U_k$$

$$Y_{k} = L.T_{k+1} + R.X_{k} + S.U_{k}$$

State-space computation

< 同→

Implicit specialized state-space form

$$\begin{pmatrix} J & 0 & 0 \\ -K & I & 0 \\ -L & 0 & I \end{pmatrix} \begin{pmatrix} T_{k+1} \\ X_{k+1} \\ Y_k \end{pmatrix} = \begin{pmatrix} 0 & M & N \\ 0 & P & Q \\ 0 & R & S \end{pmatrix} \begin{pmatrix} T_k \\ X_k \\ U_k \end{pmatrix}$$

It corresponds to:

$$J.T_{k+1} = M.X_k + N.U_k$$

$$X_{k+1} = K.T_{k+1} + P.X_k + Q.U_k$$

$$Y_k = L. T_{k+1} + R. X_k + S. U_k$$

Output(s) computation

< 17 ▶

Implicit specialized state-space form

$$\begin{pmatrix} J & 0 & 0 \\ -K & I & 0 \\ -L & 0 & I \end{pmatrix} \begin{pmatrix} T_{k+1} \\ X_{k+1} \\ Y_k \end{pmatrix} = \begin{pmatrix} 0 & M & N \\ 0 & P & Q \\ 0 & R & S \end{pmatrix} \begin{pmatrix} T_k \\ X_k \\ U_k \end{pmatrix}$$

It is equivalent to the system

$$H: z \mapsto C_Z(zI_n - A_Z)B_Z + D_Z$$

with

$$\begin{pmatrix} A_Z & B_Z \\ C_Z & D_Z \end{pmatrix} \triangleq \begin{pmatrix} K \\ L \end{pmatrix} J^{-1} \begin{pmatrix} M & N \end{pmatrix} + \begin{pmatrix} P & Q \\ R & S \end{pmatrix}$$

T. Hilaire

Bit Accurate Roundoff Noise Analysis of Fixed-Point Linear Controllers

6/34

Intermediate variables

The intermediate variables introduced allow to

- make explicit all the computations done
- show the order of the computations
- express a larger parametrization

Implicit realization

The intermediate variables computation is expressed by

 $J.T_{k+1} = M.X_k + N.U_k$

with J lower triangular with 1 on diagonal, so

- no need to compute J^{-1}
- an intermediate variable may be computed from another one previously computed (in the same stage)
 - \Rightarrow can express realizations like $Y_k = M_1.M_2...M_iU_k$

Intermediate variables

The intermediate variables introduced allow to

- make explicit all the computations done
- show the order of the computations
- express a larger parametrization

Implicit realization

The intermediate variables computation is expressed by

$$J.T_{k+1} = M.X_k + N.U_k$$

with J lower triangular with 1 on diagonal, so

- no need to compute J^{-1}
- an intermediate variable may be computed from another one previously computed (in the same stage)
 - \Rightarrow can express realizations like $Y_k = M_1.M_2...M_iU_k$

A realization with the $\delta\text{-operator}$ is described by :

$$\begin{cases} \delta X_k = A_{\delta} X_k + B_{\delta} U_k \\ Y_k = C_{\delta} X_k + D_{\delta} U_k \end{cases} \qquad \delta \triangleq \frac{q-1}{\Delta} \end{cases}$$

It is computed with

$$\begin{cases} T = A_{\delta}X_k + B_{\delta}U_k \\ X_{k+1} = X_k + \Delta T \\ Y_k = C_{\delta}X_k + D_{\delta}U_k \end{cases}$$

and it corresponds to the following implicit state-space :

$$\begin{pmatrix} I & 0 & 0 \\ -\Delta I & I & 0 \\ 0 & 0 & I \end{pmatrix} \begin{pmatrix} T_{k+1} \\ X_{k+1} \\ Y_k \end{pmatrix} = \begin{pmatrix} 0 & A_{\delta} & B_{\delta} \\ 0 & I & 0 \\ 0 & C_{\delta} & D_{\delta} \end{pmatrix} \begin{pmatrix} T_k \\ X_k \\ U_k \end{pmatrix}$$

A realization with the $\delta\text{-operator}$ is described by :

$$\begin{cases} \delta X_k = A_{\delta} X_k + B_{\delta} U_k \\ Y_k = C_{\delta} X_k + D_{\delta} U_k \end{cases} \qquad \delta \triangleq \frac{q-1}{\Delta}$$

It is computed with

$$\begin{cases} T = A_{\delta}X_k + B_{\delta}U_k \\ X_{k+1} = X_k + \Delta T \\ Y_k = C_{\delta}X_k + D_{\delta}U_k \end{cases}$$

and it corresponds to the following implicit state-space :

$$\begin{pmatrix} I & 0 & 0 \\ -\Delta I & I & 0 \\ 0 & 0 & I \end{pmatrix} \begin{pmatrix} T_{k+1} \\ X_{k+1} \\ Y_k \end{pmatrix} = \begin{pmatrix} 0 & A_{\delta} & B_{\delta} \\ 0 & I & 0 \\ 0 & C_{\delta} & D_{\delta} \end{pmatrix} \begin{pmatrix} T_k \\ X_k \\ U_k \end{pmatrix}$$

A realization with the $\delta\text{-operator}$ is described by :

$$\begin{cases} \delta X_k = A_{\delta} X_k + B_{\delta} U_k \\ Y_k = C_{\delta} X_k + D_{\delta} U_k \end{cases} \qquad \delta \triangleq \frac{q-1}{\Delta}$$

It is computed with

$$\begin{cases} T = A_{\delta}X_k + B_{\delta}U_k \\ X_{k+1} = X_k + \Delta T \\ Y_k = C_{\delta}X_k + D_{\delta}U_k \end{cases}$$

and it corresponds to the following implicit state-space :

$$\begin{pmatrix} I & 0 & 0 \\ -\Delta I & I & 0 \\ 0 & 0 & I \end{pmatrix} \begin{pmatrix} T_{k+1} \\ X_{k+1} \\ Y_k \end{pmatrix} = \begin{pmatrix} 0 & A_{\delta} & B_{\delta} \\ 0 & I & 0 \\ 0 & C_{\delta} & D_{\delta} \end{pmatrix} \begin{pmatrix} T_k \\ X_k \\ U_k \end{pmatrix}$$

One can find the Direct Form II transposed with $\delta\text{-operator}$



with

T. Hilaire



Roundoff Noise Analysis

Fixed point representation

The real numbers are represented by fixed point numbers.

Fixed-point representation

- A number is represented by
 - $2^{-\gamma}.N$ *N* : signed integer with β bits
 - γ : fixed integer (scaling)
- The quantization step $2^{-\gamma}$ is fixed, the dynamic is fixed (and limited)





To Quantize a signal x(k) is equivalent to add a independent white noise e(k). Its first and second-order moments characterize it.

Roundoff Noise Analysis
Quantization

$$x(k)$$
 $Q[]$ $x'(k)$ \equiv $x(k)$ $x'(k)$

To Quantize a signal x(k) is equivalent to add a independent white noise e(k).

Its first and second-order moments characterize it.

The first (μ) and second (σ , ψ) order moments are defined by

$$\mu_{e} \triangleq E \{e(k)\}$$

$$\psi_{e} \triangleq E \{(e(k) - \mu_{e})(e(k) - \mu_{e})^{\top}\}$$

$$\sigma_{e}^{2} \triangleq E \{(e(k) - \mu_{e})^{\top}(e(k) - \mu_{e})\} = tr(\psi_{e})$$



To Quantize a signal x(k) is equivalent to add a independent white noise e(k).

Its first and second-order moments characterize it. *Right shifting of d bits :*

	truncation	best roundoff
μ_{e}	$2^{-\gamma-1}(1-2^{-d})$	$2^{-\gamma - d - 1}$
σ_e^2	$\frac{2^{-2\gamma}}{12}(1-2^{-2d})$	$\frac{2^{-2\gamma}}{12}(1-2^{-2d})$

When implemented, the algorithm used is:

$$\begin{cases} J.T_{k+1} \leftarrow M.X_k + N.U_k \\ X_{k+1} \leftarrow K.T_{k+1} + P.X_k + Q.U_k \\ Y_k \leftarrow L.T_{k+1} + R.X_k + S.U_k \end{cases}$$

Quantizations during computations lead to noise addition ξ_k :

$$\xi_k \triangleq \begin{pmatrix} \xi_{\mathcal{T}_k} \\ \xi_{\mathcal{X}_k} \\ \xi_{\mathcal{Y}_k} \end{pmatrix}$$

T. Hilaire

When implemented, the algorithm used is:

$$\begin{cases} J.T_{k+1} \leftarrow M.X_k + N.U_k + \xi_{T_k} \\ X_{k+1} \leftarrow K.T_{k+1} + P.X_k + Q.U_k + \xi_{X_k} \\ Y_k \leftarrow L.T_{k+1} + R.X_k + S.U_k + \xi_{Y_k} \end{cases}$$

Quantizations during computations lead to noise addition ξ_k :

$$\xi_k \triangleq \begin{pmatrix} \xi_{\mathcal{T}_k} \\ \xi_{\mathcal{X}_k} \\ \xi_{\mathcal{Y}_k} \end{pmatrix}$$

T. Hilaire

Output Roundoff noise power

The output roundoff noise power is defined as the power of the noises added on the output(s):

$${\cal P} \triangleq \sigma_{\xi'}^2$$

The implemented system is equivalent to

Output Roundoff noise power

The output roundoff noise power is defined as the power of the noises added on the output(s):

$$\mathsf{P} \triangleq \sigma_{\xi'}^2$$

The implemented system is equivalent to



Considering the moments of ξ_k through *G*, we've got *P*:

$$P = tr\left(\psi_{\xi}\left(M_{2}^{\top}M_{2} + M_{1}^{\top}W_{o}M_{1}\right)\right) + \mu_{\xi'}^{\top}\mu_{\xi'}$$

where $\mu_{\xi'} = H_{1}(0)\mu_{\xi} = (C_{Z}(I - A_{Z})^{-1}M_{1} + M_{2})\mu_{\xi}$
 ψ_{ξ} and μ_{ξ} depends only on HW/SW considerations, whereas the other terms depends on the realizations.

Considering the moments of ξ_k through *G*, we've got *P*:

$$P = tr\left(\psi_{\xi}\left(M_{2}^{\top}M_{2} + M_{1}^{\top}W_{o}M_{1}\right)\right) + \mu_{\xi'}^{\top}\mu_{\xi'}$$

where $\mu_{\xi'} = H_1(0)\mu_{\xi} = (C_Z(I - A_Z)^{-1}M_1 + M_2)\mu_{\xi}$ ψ_{ξ} and μ_{ξ} depends only on HW/SW considerations, whereas the other terms depends on the realizations.



We supposed the following wordlengths known

- β_Z : coefficient's wordlength
- β_U , β_Y , β_T , β_X : intputs, outputs, intermediate variables and states' wordlength
- β_{ADD} : accumulator's wordlength

 γ_U is also known ($\gamma_U = \beta_U - 1 - \left| \log_2 \overset{\text{max}}{U} \right|$).

It is also supposed that the accumulations (in each scalar product) are done on the same fixed-point format (no shift between two additions).

∃ ► < ∃ ►

We supposed the following wordlengths known

- β_Z : coefficient's wordlength
- β_U , β_Y , β_T , β_X : intputs, outputs, intermediate variables and states' wordlength
- β_{ADD} : accumulator's wordlength

 γ_U is also known ($\gamma_U = \beta_U - 1 - \left| \log_2 \overset{\text{max}}{U} \right|$).

It is also supposed that the accumulations (in each scalar product) are done on the same fixed-point format (no shift between two additions).

Scalar product





Peak value estimation

$$\begin{pmatrix} \gamma_{T} \\ \gamma_{X} \\ \gamma_{Y} \end{pmatrix} = \begin{pmatrix} \beta_{T} \\ \beta_{X} \\ \beta_{Y} \end{pmatrix} - 2.\mathbb{1}_{l+n+p,1} - \left\lfloor \log_{2} \left(\left\| H_{\max} \right\|_{l_{1}} \left\| U \right| \right) \right\rfloor$$

The common fixed-point format of each accumulator γ_{ADD} can be set, in order to represent

- the dynamic of each product without overflow
- the final result without overflow

$$\gamma_{ADD} = \beta_{ADD} - \max_{row} \left(\begin{pmatrix} \beta_T \\ \beta_X \\ \beta_Y \end{pmatrix} - \beta_g - \begin{pmatrix} \gamma_T \\ \gamma_X \\ \gamma_Y \end{pmatrix}, \alpha \right)$$

where

$$\alpha = \max_{row} \left(\beta_{Z} - \gamma_{Z} + \mathbb{1}_{l+n+p,1} \left(\begin{pmatrix} \beta_{T} \\ \beta_{X} \\ \beta_{U} \end{pmatrix} - \begin{pmatrix} \gamma_{T} \\ \gamma_{X} \\ \gamma_{U} \end{pmatrix}^{T} \right) \right)$$

T. Hilaire

Peak value estimation

$$\begin{pmatrix} \gamma_{T} \\ \gamma_{X} \\ \gamma_{Y} \end{pmatrix} = \begin{pmatrix} \beta_{T} \\ \beta_{X} \\ \beta_{Y} \end{pmatrix} - 2.\mathbb{1}_{l+n+p,1} - \left\lfloor \log_{2} \left(\left\| \mathcal{H}_{\max} \right\|_{l_{1}} \left\| \mathcal{U} \right\| \right) \right\rfloor$$

The common fixed-point format of each accumulator $\gamma_{\rm ADD}$ can be set, in order to represent

- the dynamic of each product without overflow
- the final result without overflow

$$\gamma_{ADD} = \beta_{ADD} - \max_{row} \left(\begin{pmatrix} \beta_T \\ \beta_X \\ \beta_Y \end{pmatrix} - \beta_g - \begin{pmatrix} \gamma_T \\ \gamma_X \\ \gamma_Y \end{pmatrix}, \alpha \right)$$

where

$$\alpha = \max_{row} \left(\beta_{Z} - \gamma_{Z} + \mathbb{1}_{l+n+p,1} \left(\begin{pmatrix} \beta_{T} \\ \beta_{X} \\ \beta_{U} \end{pmatrix} - \begin{pmatrix} \gamma_{T} \\ \gamma_{X} \\ \gamma_{U} \end{pmatrix} \right)^{\mathsf{T}} \right)$$

In order to align the results of the products, 2 computational schemes are possible

Computational scheme

- Roundoff After Multiplication : the result of the product is quantized
- Roundoff Before Multiplication : the coefficient is quantized

$\Rightarrow \gamma_Z$ is then deduced

Since the wordlengths, the binary point positions and the quantizations are know, the moments ψ_{ξ} and μ_{ξ} can be (analytically) expressed.

글 > - < 글 >

In order to align the results of the products, 2 computational schemes are possible

Computational scheme

- Roundoff After Multiplication : the result of the product is quantized
- Roundoff Before Multiplication : the coefficient is quantized

$\Rightarrow \gamma_Z$ is then deduced

Since the wordlengths, the binary point positions and the quantizations are know, the moments ψ_{ξ} and μ_{ξ} can be (analytically) expressed.

ヨト イヨト

$$\begin{cases} X(k+1) = \begin{pmatrix} 0.58399 & -0.42019 \\ 0.42019 & 0.1638 \end{pmatrix} X(k) + \begin{pmatrix} 0.64635 \\ -0.23982 \end{pmatrix} \\ Y(k) = \begin{pmatrix} 0.64635 & 0.23982 \end{pmatrix} X(k) + 0.13111U(k) \end{cases}$$

Bit Accurate Roundoff Noise Analysis of Fixed-Point Linear Controllers

3

$$\begin{cases} X(k+1) = \begin{pmatrix} 0.58399 & -0.42019 \\ 0.42019 & 0.1638 \end{pmatrix} X(k) + \begin{pmatrix} 0.64635 \\ -0.23982 \end{pmatrix} \\ Y(k) = \begin{pmatrix} 0.64635 & 0.23982 \end{pmatrix} X(k) + 0.13111U(k) \end{cases}$$

$$\beta_U = \beta_X = \beta_Y = 16$$
$$\beta_Z = 16I_3$$
$$\beta_{ADD} = \begin{pmatrix} 32\\32\\32 \end{pmatrix}$$
$$\overset{\text{max}}{U} = 10 \implies \gamma_U = 1$$
$$\gamma_{ADD} = \begin{pmatrix} 26\\27\\26 \end{pmatrix}$$
$$\gamma_Z = \begin{pmatrix} 15 & 16 & 15\\16 & 17 & 17\\15 & 17 & 17 \end{pmatrix}$$
$$\gamma_X = \begin{pmatrix} 11\\11 \end{pmatrix}$$

Roundoff After Multiplication

Intermediate variables

 $\begin{array}{l} Acc \leftarrow (xn(1)*19136);\\ Acc \leftarrow Acc + (xn(2)*-27537) >> 1;\\ Acc \leftarrow Acc + (u*21179);\\ xnp(1) \leftarrow Acc >> 15;\\ Acc \leftarrow (xn(1)*27537);\\ Acc \leftarrow Acc + (xn(2)*21470) >> 1;\\ Acc \leftarrow Acc + (xn(2)*21470) >> 1;\\ xnp(2) \leftarrow Acc + (xn(2)*31433) >> 1;\\ xnp(2) \leftarrow Acc + (xn(2)*31433) >> 2;\\ Acc \leftarrow Acc + (uu; 17184) >> 2;\\ Acc \leftarrow Acc + (uu; 17184) >> 2;\\ y \leftarrow Acc >> 15;\\ Permutations\\ xn \leftarrow xnp;\\ \end{array}$

▲口▶ ▲圖▶ ▲温▶ ▲温▶ 三連

T. Hilaire

$$\begin{cases} X(k+1) = \begin{pmatrix} 0.58399 & -0.42019 \\ 0.42019 & 0.1638 \end{pmatrix} X(k) + \begin{pmatrix} 0.64635 \\ -0.23982 \end{pmatrix} \\ Y(k) = \begin{pmatrix} 0.64635 & 0.23982 \end{pmatrix} X(k) + 0.13111U(k) \end{cases}$$

$$\begin{aligned} \beta_U &= \beta_X = \beta_Y = 16\\ \beta_Z &= 16I_3\\ \beta_{ADD} &= \begin{pmatrix} 32\\ 32\\ 32 \end{pmatrix}\\ U &= 10 \implies \gamma_U = 11\\ \gamma_{ADD} &= \begin{pmatrix} 26\\ 27\\ 26 \end{pmatrix}\\ \gamma_Z &= \begin{pmatrix} 15 & 16 & 15\\ 16 & 17 & 17\\ 15 & 17 & 17 \end{pmatrix}\\ \gamma_X &= \begin{pmatrix} 11\\ 11 \end{pmatrix} \end{aligned}$$

Roundoff After Multiplication

Intermediate variables

 $\begin{array}{l} Acc \leftarrow (xn(1) * 19136);\\ Acc \leftarrow Acc + (xn(2) * -27537) >> 1;\\ Acc \leftarrow Acc + (x * 21179);\\ xnp(1) \leftarrow Acc >> 15;\\ Acc \leftarrow (xn(1) * 27537);\\ Acc \leftarrow (xn(1) * 27537);\\ Acc \leftarrow Acc + (xn(2) * 21470) >> 1;\\ Acc \leftarrow Acc + (xn(2) * 21473) >> 1;\\ xnp(2) \leftarrow Acc >> 16;\\ \hline {Outputs}\\ Acc \leftarrow (xn(1) * 21179);\\ Acc \leftarrow Acc + (xn(2) * 31433) >> 2;\\ Acc \leftarrow Acc + (x * 17184) >> 2;\\ y \leftarrow Acc >> 15;\\ \hline Permutations\\ xn \leftarrow xnp; \end{array}$

・ロット (雪) (日) (日) (日)

$$\begin{cases} X(k+1) = \begin{pmatrix} 0.58399 & -0.42019 \\ 0.42019 & 0.1638 \end{pmatrix} X(k) + \begin{pmatrix} 0.64635 \\ -0.23982 \end{pmatrix} \\ Y(k) = \begin{pmatrix} 0.64635 & 0.23982 \end{pmatrix} X(k) + 0.13111U(k) \end{cases}$$

$$\beta_U = \beta_X = \beta_Y = 16$$
$$\beta_Z = 16I_3$$
$$\beta_{ADD} = \begin{pmatrix} 32\\32\\32 \end{pmatrix}$$
$$U = 10 \implies \gamma_U = 11$$
$$\gamma_{ADD} = \begin{pmatrix} 26\\27\\26 \end{pmatrix}$$
$$\gamma_Z = \begin{pmatrix} 15 & 15 & 15\\16 & 16 & 16\\15 & 15 & 15 \end{pmatrix}$$
$$\gamma_X = \begin{pmatrix} 11\\11 \end{pmatrix}$$

Roundoff Before Multiplication

・ロト ・雪 ト ・ ヨ ト ・ ヨ ト ・ ヨ

Intermediate variables

 $\begin{array}{l} Acc \leftarrow (xn(1) * 19136);\\ Acc \leftarrow Acc + (xn(2) * -13769);\\ Acc \leftarrow Acc + (u * 21179);\\ xnp(1) \leftarrow Acc >> 15;\\ Acc \leftarrow Acc + (u * 27537);\\ Acc \leftarrow Acc + (xn(2) * 10735);\\ Acc \leftarrow Acc + (u * -15717);\\ xnp(2) \leftarrow Acc >> 16;\\ \hline Outputs\\ Acc \leftarrow (xn(1) * 21179);\\ Acc \leftarrow Acc + (xn(2) * 7858);\\ Acc \leftarrow Acc + (u * 4296);\\ y \leftarrow Acc >> 15;\\ \hline Permutations\\ xn \leftarrow xnp;\\ \end{array}$

4

Optimal fixed-point implementation

Optimal design

It is possible to analytically describe equivalent classes of realization (*Inclusion Principle*)

Equivalent realization

Consider a realization $\mathcal{R}_0.$ All realizations \mathcal{R}_1 such that

$$\begin{pmatrix} -J_1 & M_1 & N_1 \\ K_1 & P_1 & Q_1 \\ L_1 & R_1 & S_1 \end{pmatrix} = \begin{pmatrix} \mathcal{V} & & \\ & \mathcal{U}^{-1} & & \\ & & I_p \end{pmatrix} \begin{pmatrix} -J_0 & M_0 & N_0 \\ K_0 & P_0 & Q_0 \\ L_0 & R_0 & S_0 \end{pmatrix} \begin{pmatrix} \mathcal{W} & & & \\ & \mathcal{U} & & \\ & & I_m \end{pmatrix}$$

are equivalent (with $\mathcal{U} \in \mathbb{R}^{n \times n}$, $\mathcal{Y} \in \mathbb{R}^{l \times l}$ and $\mathcal{W} \in \mathbb{R}^{l \times l}$ non-singular matrices).

 $\mathsf{State-space}:\,(A,B,C,D)\to(\mathcal{T}^{-1}A\mathcal{T},\mathcal{T}^{-1}B,C\mathcal{T},D)$

Optimal design

It is possible to analytically describe equivalent classes of realization (*Inclusion Principle*)

Equivalent realization

Consider a realization $\mathcal{R}_0.$ All realizations \mathcal{R}_1 such that

$$\begin{pmatrix} -J_1 & M_1 & N_1 \\ K_1 & P_1 & Q_1 \\ L_1 & R_1 & S_1 \end{pmatrix} = \begin{pmatrix} \mathcal{Y} & & \\ & \mathcal{U}^{-1} & \\ & & I_p \end{pmatrix} \begin{pmatrix} -J_0 & M_0 & N_0 \\ K_0 & P_0 & Q_0 \\ L_0 & R_0 & S_0 \end{pmatrix} \begin{pmatrix} \mathcal{W} & & \\ & \mathcal{U} & \\ & & I_m \end{pmatrix}$$

are equivalent (with $\mathcal{U} \in \mathbb{R}^{n \times n}$, $\mathcal{Y} \in \mathbb{R}^{l \times l}$ and $\mathcal{W} \in \mathbb{R}^{l \times l}$ non-singular matrices).

Optimal realization problem

The optimal design problem consists in finding the realization \mathcal{R}^{opt} that minimizes $\mathcal J$

$$\mathcal{R}^{\mathsf{opt}} = rgmin_{\mathcal{R}\in\mathscr{R}_H} \mathcal{J}(\mathcal{R})$$

T. Hilaire

We consider the following low-pass butterworth filter

$$H(z) = \frac{0.003622z^2 + 0.007243z + 0.003622}{z^2 - 1.823z + 0.8372}$$

And the following realizations

- Z_1 : direct form II with shift-operator,
- Z₂: roundoff noise-optimal state-space realization,
- Z_3 : roundoff noise-optimal δ -realization.

16 bits for the coefficients and variables. 32 bits for the accumulator.

Roundoff Before Multiplication

Example



The optimizations are done with Adaptative Simulated Annealing method.

realization	Roundoff	Nb. operations
<i>Z</i> ₁	3.914 <i>e</i> - 3	4+ 5×
Z ₂	3.903 <i>e</i> - 7	6+ 9×
<i>Z</i> ₃	3.540 <i>e</i> - 7	$8+11\times$



1 DAC

Conclusions

Conclusion

Some tools are exhibited to help to answer the question: How *optimally* implement filter/controllers ?

- Implicit state-space framework
- Roundoff noise analysis
- Two bit-accurate fixed-point implementation schemes

A Matlab's toolbox (*Finite Wordlength Realization Toolbox*) was developed, with some others Finite WordLength measures (sensitivity):

http://fwrtoolbox.gforge.inria.fr

Conclusions

Conclusion

Some tools are exhibited to help to answer the question: How *optimally* implement filter/controllers ?

- Implicit state-space framework
- Roundoff noise analysis
- Two bit-accurate fixed-point implementation schemes

A Matlab's toolbox (*Finite Wordlength Realization Toolbox*) was developed, with some others Finite WordLength measures (sensitivity):

http://fwrtoolbox.gforge.inria.fr

Any questions ?

Bit Accurate Roundoff Noise Analysis of Fixed-Point Linear Controllers

3

∃ ► < ∃ ► ...</p>



Appendix

Roundoff Before Multiplication scheme.

- 16 bits for the coefficients, states, inputs, outputs, intermediate variables
- Accumulator 32 bits (with 4 guard bits)



Direct Form Directe II

// intermediate variables

Exemple

Balanced state-space form

// intermediate variables Acc $\leftarrow xn(1) * 32370 + xn(2) * -3673 + Acc0 + xn(3) * 42 + xn(4) * 873 + xn(5) * -171 + xn(6) * 51 + u(i) * 32370 + xn(5) + x$ 458: $xnp(1) \leftarrow Acc >> 15$: Acc $\leftarrow xn(1) * 3688 + xn(2) * 31858 + xn(3) * 4405 + xn(4) * 785 + xn(5) * -451 + xn(6) * 74 + u(i) * -605;$ $xnp(2) \leftarrow Acc >> 15;$ Acc \leftarrow xn(1) * 318 + xn(2) * -4416 + xn(3) * 31250 + xn(4) * -3922 + xn(5) * 499 + xn(6) * -120 + u(i) * -790: $xnp(3) \leftarrow Acc >> 15;$ Acc $\leftarrow xn(1) * -900 + xn(2) * 546 + xn(3) * 3833 + xn(4) * 30304 + xn(5) * 1828 + xn(6) * -196 + u(i) * 742;$ $xnp(4) \leftarrow Acc >> 15$: Acc $\leftarrow xn(1) * -551 + xn(2) * 1425 + xn(3) * 1760 + xn(4) * -7483 + xn(5) * 29956 + xn(6) * 1961 + u(i) * xn(5) + xn($ 868: $xnp(5) \leftarrow Acc >> 15$: Acc $\leftarrow xn(1) * -786 + xn(2) * 1182 + xn(3) * 2572 + xn(4) * -4839 + xn(5) * -7995 + xn(6) * 29485 + u(i) * (1) * (1) +$ 956: $xnp(6) \leftarrow Acc >> 15;$ Acc $\leftarrow xn(1) * 14733 + xn(2) * 21060 + xn(3) * -23783 + xn(4) * -22615 + xn(5) * 7488 + xn(6) * -1780 + xn(5) * 21060 + xn(5)$ u(i) * 77: $v(i) \leftarrow Acc >> 15$: // permutations $xn \leftarrow xnp;$

3

・ロト ・聞 と ・ 聞 と ・ 聞 と …

Exemple

Direct form II with δ -operator

```
// intermediate variables
Acc \leftarrow xn(1) * -11383 + xn(2) * -31123 + xn(3) * -22773 + xn(4) * -13468 + xn(5) * -9425 + xn(6) * -1852 + xn(6) + xn
u(i) << 8;
T0 \leftarrow Acc >> 13:
// states
Acc \leftarrow T0 + xn(1) <<2;
xn(1) \leftarrow Acc >>2:
Acc \leftarrow xn(1) + xn(2) <<3;
xn(2) \leftarrow Acc >>3;
Acc \leftarrow xn(2) + xn(3) <<2;
xn(3) \leftarrow Acc >>2:
Acc \leftarrow xn(3) + xn(4) <<2;
xn(4) \leftarrow Acc >>2;
Acc \leftarrow xn(4) + xn(5) <<3;
xn(5) \leftarrow Acc >>3:
Acc \leftarrow xn(5) + xn(6) <<2;
xn(6) \leftarrow Acc >>2: // outputs
Acc \leftarrow xn(1) * 792 + xn(2) * 12559 + xn(3) * 12190 + xn(4) * 29211 + xn(5) * 22483 + xn(6) * 30784 + u(i)
* 19:
y(i) \leftarrow Acc >> 13);
```

< ロ > < 同 > < 回 > < 回 >

Exemple

Direct form II with δ -operator : Roundoff After Multiplication

```
// intermediate variables
Acc \leftarrow (xn(1) * -22767) >> 1 + xn(2) * -31123 + xn(3) * -22773 + (xn(4) * -26936) >> 1 + (xn(5) * -18851)
>> 1 + (xn(6) * -29635) >> 4 + u(i) << 8;
T0 \leftarrow Acc >> 13:
// states
Acc \leftarrow T0 + xn(1) <<2;
xn(1) \leftarrow Acc >>2;
Acc \leftarrow xn(1) + xn(2) <<3;
xn(2) \leftarrow Acc >>3;
Acc \leftarrow xn(2) + xn(3) <<2;
xn(3) \leftarrow Acc >>2:
Acc \leftarrow xn(3) + xn(4) <<2;
xn(4) \leftarrow Acc >>2;
Acc \leftarrow xn(4) + xn(5) <<3;
xn(5) \leftarrow Acc >>3:
Acc \leftarrow xn(5) + xn(6) <<2;
xn(6) \leftarrow Acc >>2: // outputs
Acc \leftarrow (xn(1) * 25342) >> 5 + (xn(2) * 25118) >> 1 + (xn(3) * 24379) >> 1 + xn(4) * 29211 + xn(5) *
22483 + xn(6) * 30784 + (u(i) * 19746) >> 10;
v(i) \leftarrow Acc >>13
```

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >