

Optimal word-length allocation for the fixed-point implementation of linear filters and controllers

Thibault Hilaire, Hacène Ouzia and Benoit Lopez

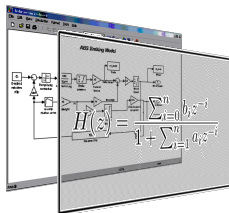
Sorbonne Université, France

Inria, Paris Saclay

ARITH26 2019, June 2019, Kyoto

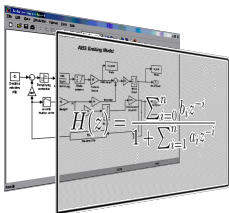


Context



Filters/Controllers
Algorithms

Context

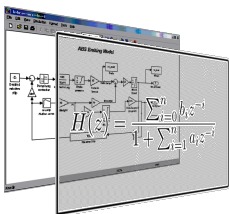


Filters/Controllers
Algorithms



Embedded
Hardware

Context



Filters/Controllers
Algorithms

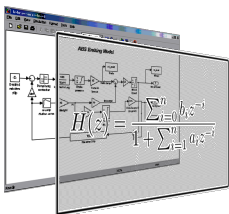
Implementation



Embedded
Hardware

👉 **implementation**: transformation of the mathematical object into finite precision operations to be performed on a specific target

Context



Filters/Controllers
Algorithms

Implementation



Embedded
Hardware

👉 **implementation**: transformation of the mathematical object into finite precision operations to be performed on a specific target

Specially for Fixed-Point arithmetic, the implementation process is:

- time-consuming and error-prone
- provides no guaranty on the the errors
- potentially *non-optimal*

We need a code generator!

Word-length optimization for FPGA/ASIC

We will here focus on

- State-Space systems

Word-length optimization for FPGA/ASIC

We will here focus on

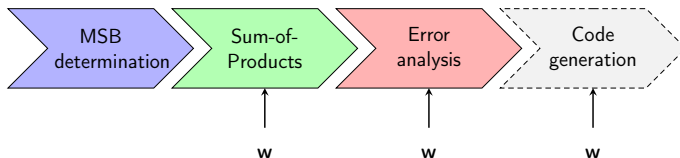
- State-Space systems
- the first parts of the flow, except the code generation



Word-length optimization for FPGA/ASIC

We will here focus on

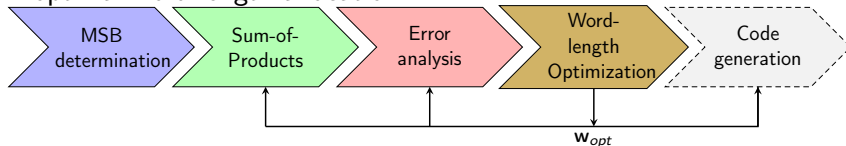
- State-Space systems
- the first parts of the flow, except the code generation
- for FPGAs or ASICs (multiple word-length paradigm)
 w : word-lengths



Word-length optimization for FPGA/ASIC

We will here focus on

- State-Space systems
- the first parts of the flow, except the code generation
- for FPGAs or ASICs (multiple word-length paradigm)
 w : word-lengths
- optimal word-length allocation



We also want it to be reliable:

- ☞ classical Signal Processing approach models errors as *noises* and perform statistical error analysis
- ☞ we want to use *worst case analysis* (rigorous bounds)

Outline

- 1 Introduction
 - Context
 - State-Space systems
- 2 Fixed-Point implementation
 - Determining the Fixed-Point Formats
 - Sum-of-Products by real Constants
- 3 Word-length optimization under accuracy constraint
 - Error analysis
 - Word-length allocation problem
 - Examples
- 4 Conclusions and Perspectives

State-Space

We consider Linear Time Invariant filters (or controllers) expressed with State-Space systems:

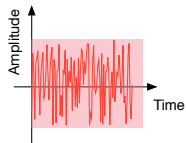
$$\mathcal{H} \begin{cases} \mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \mathbf{y}(k) = \mathbf{C}\mathbf{x}(k) + \mathbf{D}\mathbf{u}(k) \end{cases}$$

where

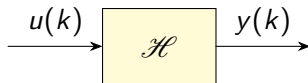
- $\mathbf{u}(k)$ and $\mathbf{y}(k)$ are the input and output vector *at time* k ;
- $\mathbf{x}(k)$ is the state vector ;
- \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} are matrices defining the system.

Worst-Case Peak Gain

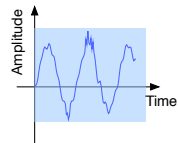
Input $u(k)$



Stable filter



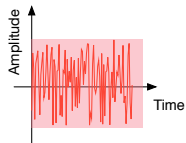
Output $y(k)$



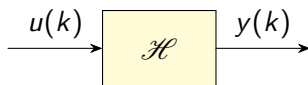
¹Volkova, H., and Lauter, "Reliable evaluation of the Worst-Case Peak Gain matrix in multiple precision," in 22nd IEEE Symposium on Computer Arithmetic, 2015

Worst-Case Peak Gain

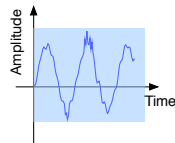
Input $u(k)$



Stable filter



Output $y(k)$

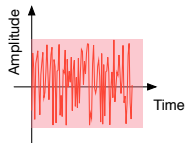


$$\forall k, |u(k)| \leq \bar{u}$$

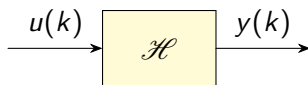
¹Volkova, H., and Lauter, "Reliable evaluation of the Worst-Case Peak Gain matrix in multiple precision," in 22nd IEEE Symposium on Computer Arithmetic, 2015

Worst-Case Peak Gain

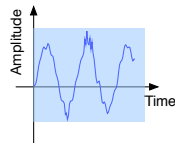
Input $u(k)$



Stable filter



Output $y(k)$



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

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

Worst-Case Peak Gain: $\langle\langle \mathcal{H} \rangle\rangle = |D| + \sum_{k=0}^{\infty} |CA^k B|$

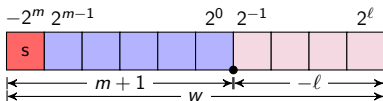
The WCPG matrix can be computed at any arbitrary precision¹.

¹Volkova, H., and Lauter, "Reliable evaluation of the Worst-Case Peak Gain matrix in multiple precision," in 22nd IEEE Symposium on Computer Arithmetic, 2015

Outline

- 1 Introduction
- 2 Fixed-Point implementation
 - Determining the Fixed-Point Formats
 - Sum-of-Products by real Constants
- 3 Word-length optimization under accuracy constraint
- 4 Conclusions and Perspectives

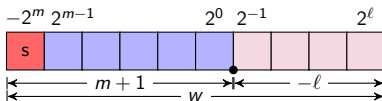
Fixed-Point Arithmetic



$$x = M \cdot 2^l$$

M : mantissa (two's complement) $\in [-2^{w-1}; 2^{w-1} - 1]$

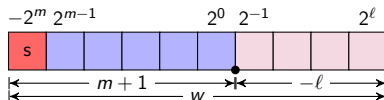
Fixed-Point Arithmetic



$$w = m - \ell + 1$$

w		wordlength
m		Most Significant Bit
ℓ		Least Significant Bit

Fixed-Point Arithmetic



$$w = m - \ell + 1$$

w		wordlength
m		Most Significant Bit
ℓ		Least Significant Bit

The MSB/LSB (or MSB/word-length) must be chosen carefully:

- we choose MSB such that no overflow occurs:

$$\forall k, x(k) \in [-2^m; 2^m - 2^\ell]$$

- we choose LSB such that we achieve a certain accuracy

Determining the MSB position

We consider the vector of our variables $\zeta(k) = \begin{pmatrix} \mathbf{x}(k) \\ \mathbf{y}(k) \end{pmatrix}$.

We want to determine their MSB \mathbf{m} such that

$$\forall k, \quad \zeta_i(k) \in [-2^{\mathbf{m}_i}, 2^{\mathbf{m}_i}(1 - 2^{1-\mathbf{w}_i})].$$

Determining the MSB position

We consider the vector of our variables $\zeta(k) = \begin{pmatrix} \mathbf{x}(k) \\ \mathbf{y}(k) \end{pmatrix}$.

We want to determine their MSB \mathbf{m} such that

$$\forall k, \quad \zeta_i(k) \in [-2^{\mathbf{m}_i}, 2^{\mathbf{m}_i}(1 - 2^{1-\mathbf{w}_i})].$$

Our system is rewritten as

$$\mathcal{H}_\zeta \begin{cases} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \zeta(k) &= \mathbf{M}_1\mathbf{x}(k) + \mathbf{M}_2\mathbf{u}(k), \end{cases}$$

Determining the MSB position

We consider the vector of our variables $\zeta(k) = \begin{pmatrix} \mathbf{x}(k) \\ \mathbf{y}(k) \end{pmatrix}$.

We want to determine their MSB \mathbf{m} such that

$$\forall k, \quad \zeta_i(k) \in [-2^{\mathbf{m}_i}, 2^{\mathbf{m}_i}(1 - 2^{1-\mathbf{w}_i})].$$

Our system is rewritten as

$$\mathcal{H}_\zeta \begin{cases} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \zeta(k) &= \mathbf{M}_1\mathbf{x}(k) + \mathbf{M}_2\mathbf{u}(k), \end{cases}$$

And then the least MSB is obtained with

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

Determining the MSB position

We consider the vector of our variables $\zeta(k) = \begin{pmatrix} \mathbf{x}(k) \\ \mathbf{y}(k) \end{pmatrix}$.

We want to determine their MSB \mathbf{m} such that

$$\forall k, \quad \zeta_i(k) \in [-2^{\mathbf{m}_i}, 2^{\mathbf{m}_i}(1 - 2^{1-\mathbf{w}_i})].$$

Our system is rewritten as

$$\mathcal{H}_\zeta \begin{cases} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \zeta(k) &= \mathbf{M}_1\mathbf{x}(k) + \mathbf{M}_2\mathbf{u}(k), \end{cases}$$

And then the least MSB is obtained with

$$\mathbf{m}(\mathbf{w}) = \lceil \log_2(\bar{\zeta}) \rceil + \delta(\mathbf{w}),$$

Determining the MSB position

We consider the vector of our variables $\zeta(k) = \begin{pmatrix} \mathbf{x}(k) \\ \mathbf{y}(k) \end{pmatrix}$.

We want to determine their MSB \mathbf{m} such that

$$\forall k, \quad \zeta_i(k) \in [-2^{\mathbf{m}_i}, 2^{\mathbf{m}_i}(1 - 2^{1-\mathbf{w}_i})].$$

Our system is rewritten as

$$\mathcal{H}_\zeta \begin{cases} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \zeta(k) &= \mathbf{M}_1\mathbf{x}(k) + \mathbf{M}_2\mathbf{u}(k), \end{cases}$$

And then the least MSB is obtained with

$$\mathbf{m}(\mathbf{w}) = \lceil \log_2(\bar{\zeta}) \rceil + \delta(\mathbf{w}),$$

$$\delta_i(\mathbf{w}) = \begin{cases} 0 & \text{if } \mathbf{w}_i \geq \tilde{\mathbf{w}}_i \\ 1 & \text{if } \mathbf{w}_i < \tilde{\mathbf{w}}_i \end{cases}, \quad \tilde{\mathbf{w}} \triangleq \mathbf{1} + \lceil \log_2(\bar{\zeta}) \rceil - \lfloor \log_2(2^{\lceil \log_2(\bar{\zeta}) \rceil} - \bar{\zeta}) \rfloor$$

$\tilde{\mathbf{w}}$ can be seen as a threshold. How many bits such that $\zeta_i \notin [2^{\mathbf{m}_i}(1 - 2^{\mathbf{w}_i}); 2^{\mathbf{m}_i}]$

Sum-of-Products by real Constants

In the State-Space evaluation, the Sum of Products by real Constants is the basic brick operation $r = \sum_{i=1}^N c_i v_i$

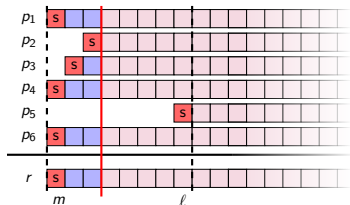
We want to compute r with last bit accuracy with format (m, ℓ)

²Volkova, Istoan, De Dinechin and H. "Towards hardware IIR filters computing just right: Direct form I case study," IEEE Transactions on Computers, 2019

Sum-of-Products by real Constants

In the State-Space evaluation, the Sum of Products by real Constants is the basic brick operation $r = \sum_{i=1}^N c_i v_i$

We want to compute r with last bit accuracy with format (m, ℓ)

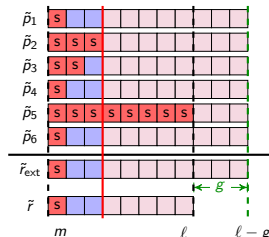


²Volkova, Istoan, De Dinechin and H. "Towards hardware IIR filters computing just right: Direct form I case study," IEEE Transactions on Computers, 2019

Sum-of-Products by real Constants

In the State-Space evaluation, the Sum of Products by real Constants is the basic brick operation $r = \sum_{i=1}^N c_i v_i$

We want to compute r with last bit accuracy with format (m, ℓ)



It can be done using g extra guard bits.

If g is large enough², the error $|r - \tilde{r}|$ is then bounded by 2^ℓ

☞ Can be done using Table-based Constant Multiplier for FPGAs

²Volkova, Istoan, De Dinechin and H. "Towards hardware IIR filters computing just right: Direct form I case study," IEEE Transactions on Computers, 2019

Outline

- 1 Introduction
- 2 Fixed-Point implementation
- 3 Word-length optimization under accuracy constraint
 - Error analysis
 - Word-length allocation problem
 - Examples
- 4 Conclusions and Perspectives

Error analysis – 1

The exact filter \mathcal{H} is:

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

Error analysis – 1

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

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

where $\varepsilon(k) = \begin{pmatrix} \varepsilon_x(k) \\ \varepsilon_y(k) \end{pmatrix}$ collects the errors due to the Sum-of-Products.

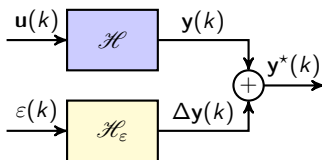
Error analysis – 1

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

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

where $\varepsilon(k) = \begin{pmatrix} \varepsilon_x(k) \\ \varepsilon_y(k) \end{pmatrix}$ collects the errors due to the Sum-of-Products.

It can be shown that the implemented system \mathcal{H}^* can be seen as



with \mathcal{H}_ε the *error filter* (state-space).

Error analysis – 2

Using last-bit accuracy Sum-of-Products, then the error $\varepsilon(k)$ is bounded by

$$\bar{\varepsilon} = 2^{\mathbf{m-w+1}}$$

Error analysis – 2

Using last-bit accuracy Sum-of-Products, then the error $\varepsilon(k)$ is bounded by

$$\bar{\varepsilon} = 2^{\mathbf{m}-\mathbf{w}+\mathbf{1}}$$

So, the output error $\Delta\mathbf{y}(k)$ is bounded by

$$\Delta\bar{\mathbf{y}} = \langle\langle\mathcal{H}_\varepsilon\rangle\rangle \bar{\varepsilon}$$

And finally

$$\Delta\bar{\mathbf{y}} = \langle\langle\mathcal{H}_\varepsilon\rangle\rangle 2^{\lceil \log_2(\langle\langle\mathcal{H}_\zeta\rangle\rangle \bar{\mathbf{u}}) \rceil + \mathbf{1} - \mathbf{w} + \delta(\mathbf{w})}$$

Error analysis – 2

Using last-bit accuracy Sum-of-Products, then the error $\varepsilon(k)$ is bounded by

$$\bar{\varepsilon} = 2^{\mathbf{m}-\mathbf{w}+\mathbf{1}}$$

So, the output error $\Delta\mathbf{y}(k)$ is bounded by

$$\Delta\bar{\mathbf{y}} = \langle\langle\mathcal{H}_\varepsilon\rangle\rangle \bar{\varepsilon}$$

And finally

$$\Delta\bar{\mathbf{y}} = \langle\langle\mathcal{H}_\varepsilon\rangle\rangle 2^{\lceil \log_2(\langle\langle\mathcal{H}_\zeta\rangle\rangle \bar{\mathbf{u}}) \rceil + \mathbf{1} - \mathbf{w} + \delta(\mathbf{w})}$$

The output error bound depends on the word-lengths \mathbf{w} .

Word-length allocation problem – 1

We want to minimize the total word-length involved *while* guarantying a certain accuracy ϵ .

Word-length allocation problem – 1

We want to minimize the total word-length involved *while* guarantying a certain accuracy ϵ .

It leads to the following optimization problem

$$\begin{aligned} \mathbf{w}_{opt} &= \arg \min \sum_i \mathbf{w}_i \\ \text{subject to} \\ \Delta \bar{\mathbf{y}}_j &\leq \epsilon_j \end{aligned}$$

Word-length allocation problem – 1

We want to minimize the total word-length involved *while* guarantying a certain accuracy ϵ .

It leads to the following optimization problem

$$\mathbf{w}_{opt} = \arg \min \sum_i \mathbf{w}_i$$

subject to

$$\sum_j \mathbf{E}_{ij} 2^{-\mathbf{w}_j + \delta_j} \leq \epsilon_j$$

$$\text{with } \mathbf{E}_{ij} \triangleq \langle\langle \mathcal{H}_\epsilon \rangle\rangle_{ij} 2^{\left\lceil \log_2 \left(\langle\langle \mathcal{H}_\zeta \rangle\rangle_{ij} \bar{\mathbf{u}} \right) \right\rceil}$$

Word-length allocation problem – 1

We want to minimize the total word-length involved *while* guarantying a certain accuracy ϵ .

It leads to the following optimization problem

$$\begin{aligned} \mathbf{w}_{opt} &= \arg \min \sum_i \mathbf{w}_i \\ \text{subject to} \\ \sum_j \mathbf{E}_{ij} 2^{-\mathbf{w}_j + \delta_j} &\leq \epsilon_j \end{aligned}$$

Recall that $\delta_j = \begin{cases} 0 & \text{if } w_j \geq \tilde{w}_j \\ 1 & \text{if } w_j < \tilde{w}_j \end{cases}$

And this condition is equivalent to (w_j and $\tilde{w}_j \in [2, u]$)

$$(2 - \tilde{w}_j) \delta_j + 1 \leq w_j - \tilde{w}_j + 1 \leq (1 - \delta_j) (u - \tilde{w}_j + 1), \quad \delta_j \in \{0, 1\}$$

Word-length allocation problem – 1

We want to minimize the total word-length involved *while* guarantying a certain accuracy ϵ .

It leads to the following optimization problem

$$\mathbf{w}_{opt} = \arg \min \sum_i \mathbf{w}_i$$

subject to

$$\sum_j \mathbf{E}_{ij} 2^{-\mathbf{w}_j + \delta_j} \leq \epsilon_j$$

$$(2 - \tilde{\mathbf{w}}_j) \delta_j + 1 \leq \mathbf{w}_j - \tilde{\mathbf{w}}_j + 1$$

$$\mathbf{w}_j - \tilde{\mathbf{w}}_j + 1 \leq (1 - \delta_j) (\mathbf{u}_j - \tilde{\mathbf{w}}_j + 1)$$

$$2 \leq \mathbf{w}_j \leq \mathbf{u}_j$$

$$\mathbf{w}_j \in \mathbb{Z}, \delta_j \in \{0, 1\}$$

Word-length allocation problem – 2

☞ Separable convex non-linear integer optimization problem
Generally solved using branch-and-bound and outer approximation methods.
Existing solvers (like Bonmin, Artylis-Knitro, etc.) can be used

Word-length allocation problem – 2

☞ Separable convex non-linear integer optimization problem
Generally solved using branch-and-bound and outer approximation methods.

Existing solvers (like Bonmin, Artylis-Knitro, etc.) can be used

Two sub-optimal problem can be defined and solved:

- Uniform word-lengths

$$w_{uni} = \max \left(\lceil \log_2 (\mathbf{E}\mathbf{1}) - \log_2(\epsilon) \rceil \right) + 1$$

Word-length allocation problem – 2

☞ Separable convex non-linear integer optimization problem
Generally solved using branch-and-bound and outer approximation methods.

Existing solvers (like Bonmin, Artylis-Knitro, etc.) can be used

Two sub-optimal problem can be defined and solved:

- Uniform word-lengths

$$w_{uni} = \max \left(\lceil \log_2 (\mathbf{E}\mathbf{1}) - \log_2(\epsilon) \rceil \right) + 1$$

- Equitably distributed budget error

Constraints $\sum_{j=1}^N \mathbf{E}_{ij} 2^{-w_j + \delta_j} \leq \epsilon_i$ are transformed in N stricter but simpler constraints $\mathbf{E}_{ij} 2^{-w_j + \delta_j} \leq \frac{\epsilon_i}{N}$

Word-length allocation problem – 2

☞ Separable convex non-linear integer optimization problem
Generally solved using branch-and-bound and outer approximation methods.

Existing solvers (like Bonmin, Artylis-Knitro, etc.) can be used

Two sub-optimal problem can be defined and solved:

- Uniform word-lengths

$$w_{uni} = \max \left(\lceil \log_2 (\mathbf{E}\mathbf{1}) - \log_2(\epsilon) \rceil \right) + 1$$

- Equitably distributed budget error

Constraints $\sum_{j=1}^N \mathbf{E}_{ij} 2^{-w_j + \delta_j} \leq \epsilon_i$ are transformed in N stricter but simpler constraints $\mathbf{E}_{ij} 2^{-w_j + \delta_j} \leq \frac{\epsilon_i}{N}$

☞ Similar results, if the cost function is the total number of bits involved in the computations

Example – 1

A 10th order active controller of longitudinal oscillation³.
Designed to remove the unpleasant oscillations of the vehicle
(acting on the engine torque).

☞ $\bar{u} = 10$ and $\epsilon = 2^{-6}$

	oscillation controller										y	$f(\mathbf{w})$
	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}		
m	9	8	9	9	8	8	8	6	6	2	4	
\tilde{w}	3	3	3	3	3	3	3	4	8	3	3	
optimal	17	15	17	17	15	15	15	11	11	3	13	149
uniform	16	16	16	16	16	16	16	16	16	16	16	176
eq. distributed	17	15	18	17	15	15	15	12	12	4	14	154

$$\bar{\zeta}_9 \approx 63.412$$

³Lefebvre, Chevrel, and Richard, "An H_∞ based control design methodology dedicated to the active control of longitudinal oscillations," IEEE Trans. on Control Systems Technology, 2003

Example – 2

A random stable 4th order State-Space, 5 inputs, 7 outputs

	random controller											
	x_1	x_2	x_3	x_4	y_1	y_2	y_3	y_4	y_5	y_6	y_7	$f(\mathbf{w})$
\mathbf{m}	6	6	5	6	16	15	16	16	16	15	16	
$\tilde{\mathbf{w}}$	3	4	3	6	4	3	6	3	4	3	6	
optimal	13	11	14	12	4	3	5	4	4	3	5	78
uniform	31	31	31	31	31	31	31	31	31	31	31	341
eq. distributed	33	32	34	32	26	25	26	26	26	25	26	311

Example – 2

A random stable 4th order State-Space, 5 inputs, 7 outputs

	random controller											
	x_1	x_2	x_3	x_4	y_1	y_2	y_3	y_4	y_5	y_6	y_7	$f(\mathbf{w})$
\mathbf{m}	6	6	5	6	16	15	16	16	16	15	16	
$\tilde{\mathbf{w}}$	3	4	3	6	4	3	6	3	4	3	6	
optimal	13	11	14	12	4	3	5	4	4	3	5	78
uniform	31	31	31	31	31	31	31	31	31	31	31	341
eq. distributed	33	32	34	32	26	25	26	26	26	25	26	311

For optimal: $\mathbf{w}_3 < \tilde{\mathbf{w}}_3$ and $\mathbf{w}_7 < \tilde{\mathbf{w}}_7$ (so $\delta_3 = \delta_7 = 0$)

Conclusion

- Reliable Fixed-Point implementation of State-Space systems
 - MSB determination
 - Last-bit accuracy sum-of-Products
- Error analysis
- Word-length allocation problem, solved with three heuristics

Conclusion

- Reliable Fixed-Point implementation of State-Space systems
 - MSB determination
 - Last-bit accuracy sum-of-Products
- Error analysis
- Word-length allocation problem, solved with three heuristics

Perspectives:

- Now consider the code generation, using dedicated tools (FloPoCo, etc.)
- Allow a more realistic cost function
- Extend this work to full class of linear filters/controllers, in order to compare various algorithms and implementations

Thank you
Any questions ?