# Determining Fixed-Point Formats using the Worst-Case Peak Gain measure

#### Anastasia Volkova, Thibault Hilaire, Christoph Lauter

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

#### ASILOMAR 49 November 10, 2015













Filter implementation flow:

• Transfer function generation





- Transfer function generation
- State-space, DFI, DFII, ...



- Transfer function generation
- State-space, DFI, DFII, ...
- Software or Hardware implementation



- Transfer function generation
  - ! Coefficient quantization
- State-space, DFI, DFII, ...
- Software or Hardware implementation



- Transfer function generation
  - ! Coefficient quantization
- State-space, DFI, DFII, ...
  - ! Large variety of structures with no common quality criteria
- Software or Hardware implementation



- Transfer function generation
  - ! Coefficient quantization
- State-space, DFI, DFII, ...
  - ! Large variety of structures with no common quality criteria
- Software or Hardware implementation
  - ! Constraints: power consumption, area, error, speed, etc.

# Motivation: Automatized filter implementation flow



# Motivation: Automatized filter implementation flow



Focus on Fixed-Point realization.

# Motivation: Automatized filter implementation flow



Focus on Fixed-Point realization. Optimization of wordlengths:

- Take more pay more
- Take less overflow risk

What we want in the end:

- Rigorous algorithm for Fixed-Point Formats (FxPF) determination
- Integration into automatized code generator for filters
- Multiple wordlength paradigm

 $\mathcal{H} := (A, B, C, D)$  is a Bounded Input Bounded Output LTI filter in state-space representation:

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

with q inputs, n states and p outputs and state matrices









### Two's complement Fixed-Point arithmetic



$$t = -2^{m}t_{m} + \sum_{i=\ell}^{m-1} 2^{i}t_{i}$$

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

#### Two's complement Fixed-Point arithmetic



$$t = -2^{m}t_{m} + \sum_{i=\ell}^{m-1} 2^{i}t_{i}$$

- quantization step:  $2^{\ell}$
- t represented by integer  $T = t \cdot 2^{\ell}$
- $\bullet \ T \in [-2^m; 2^m 2^\ell] \cap \mathbb{Z}$

#### Two's complement Fixed-Point arithmetic



$$t = -2^{m}t_{m} + \sum_{i=\ell}^{m-1} 2^{i}t_{i}$$

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

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

#### Problem statement

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

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

Given

- wordlength constraints  $\boldsymbol{w}_x$  and  $\boldsymbol{w}_y$  for each state and output variable
- ${\scriptstyle \bullet}\,$  input domain  $\bar{u}$

we need to determine the minimal FPF for all variables of filter  $\mathcal{H}$ , i.e. find the least  $m_x$  and  $m_y$  such that

$$\forall k, \quad \boldsymbol{x}_{i}(k) \in [-2^{\boldsymbol{m}_{x_{i}}}; 2^{\boldsymbol{m}_{x_{i}}} - 2^{\boldsymbol{m}_{x_{i}}-\boldsymbol{w}_{x_{i}}+1}]$$
  
$$\forall k, \quad \boldsymbol{y}_{i}(k) \in [-2^{\boldsymbol{m}_{y_{i}}}; 2^{\boldsymbol{m}_{y_{i}}} - 2^{\boldsymbol{m}_{y_{i}}-\boldsymbol{w}_{y_{i}}+1}].$$

# Modification of ${\cal H}$

Let 
$$\boldsymbol{\zeta}(k) := \begin{pmatrix} \boldsymbol{x}(k) \\ \boldsymbol{y}(k) \end{pmatrix}$$
 be a vertical concatenation of state and output vectors.

### Modification of ${\cal H}$

Let  $\boldsymbol{\zeta}(k) := \begin{pmatrix} \boldsymbol{x}(k) \\ \boldsymbol{y}(k) \end{pmatrix}$  be a vertical concatenation of state and output vectors.

Then the state-space takes the following form:

$$\mathcal{H}_{\zeta} \left\{ \begin{array}{rcl} \boldsymbol{x}(k+1) &=& \boldsymbol{A}\boldsymbol{x}(k) &+& \boldsymbol{B}\boldsymbol{u}(k) \\ \boldsymbol{\zeta}(k) &=& \begin{pmatrix} \boldsymbol{I} \\ \boldsymbol{C} \end{pmatrix} \boldsymbol{x}(k) &+& \begin{pmatrix} \boldsymbol{0} \\ \boldsymbol{D} \end{pmatrix} \boldsymbol{u}(k) \end{array} \right.$$

#### Modification of ${\cal H}$

Let  $\boldsymbol{\zeta}(k) := \begin{pmatrix} \boldsymbol{x}(k) \\ \boldsymbol{y}(k) \end{pmatrix}$  be a vertical concatenation of state and output vectors.

Then the state-space takes the following form:

$$\mathcal{H}_{\zeta} \left\{ \begin{array}{rcl} \boldsymbol{x}(k+1) &=& \boldsymbol{A}\boldsymbol{x}(k) &+& \boldsymbol{B}\boldsymbol{u}(k) \\ \boldsymbol{\zeta}(k) &=& \begin{pmatrix} \boldsymbol{I} \\ \boldsymbol{C} \end{pmatrix} \boldsymbol{x}(k) &+& \begin{pmatrix} \boldsymbol{0} \\ \boldsymbol{D} \end{pmatrix} \boldsymbol{u}(k) \end{array} \right.$$

We seek to determine the least  $m_{\zeta}$  such that

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

# Computing the MSB using the WCPG theorem

Applying the WCPG theorem on filter  $\mathcal{H}_{\zeta}$  gives

 $\forall k, \quad |\boldsymbol{\zeta}_i(k)| \leqslant (\langle\!\langle H_{\boldsymbol{\zeta}} \rangle\!\rangle \, \bar{\boldsymbol{u}})_i$ 

# Computing the MSB using the WCPG theorem

Applying the WCPG theorem on filter  $\mathcal{H}_{\zeta}$  gives

 $\forall k, |\boldsymbol{\zeta}_i(k)| \leq (\langle\!\langle H_{\boldsymbol{\zeta}} \rangle\!\rangle \, \bar{\boldsymbol{u}})_i$ 

Therefore, the smallest  $\boldsymbol{m}_{\zeta_i}$  satisfying

$$(\langle\!\langle H_{\zeta} \rangle\!\rangle \, \bar{\boldsymbol{u}})_i \leqslant 2^{\boldsymbol{m}_{\zeta_i}} - 2^{\boldsymbol{m}_{\zeta_i} - \boldsymbol{w}_i + 1},$$

will satisfy the wordlength constraints.

# Computing the MSB using the WCPG theorem

Applying the WCPG theorem on filter  $\mathcal{H}_{\zeta}$  gives

 $\forall k, |\boldsymbol{\zeta}_i(k)| \leq (\langle\!\langle H_{\boldsymbol{\zeta}} \rangle\!\rangle \, \bar{\boldsymbol{u}})_i$ 

Therefore, the smallest  $\boldsymbol{m}_{\zeta_i}$  satisfying

$$(\langle\!\langle H_{\zeta} \rangle\!\rangle \, \bar{\boldsymbol{u}})_i \leqslant 2^{\boldsymbol{m}_{\zeta_i}} - 2^{\boldsymbol{m}_{\zeta_i} - \boldsymbol{w}_i + 1},$$

will satisfy the wordlength constraints.

We can compute  $\boldsymbol{m}_{\zeta_i}$  with

$$\boldsymbol{m}_{\zeta_i} = \left\lceil \log_2 \left( \left\langle\!\left\langle H_{\zeta} \right\rangle\!\right\rangle \bar{\boldsymbol{u}} \right)_i - \log_2 \left( 1 - 2^{1 - \boldsymbol{w}_{\zeta_i}} \right) \right\rceil.$$

Taking the quantization errors into account

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

$$\mathcal{H}_{\zeta} \left\{ egin{array}{ll} m{x} & (k+1) &= & m{A} m{x} & (k) + m{B} m{u}(k) \ m{\zeta} & (k) &= & egin{array}{ll} m{I} \ m{C} \ m{x} & + egin{array}{ll} m{0} \ m{D} \ m{u}(k) \end{array} 
ight.$$

#### Taking the quantization errors into account

The actually implemented filter  $\mathcal{H}^{\Diamond}_{\mathcal{L}}$  is:

$$\mathcal{H}_{\zeta}^{\diamond} \left\{ \begin{array}{rcl} \boldsymbol{x}^{\diamond}(k+1) &=& \diamondsuit_{\boldsymbol{m}_{\boldsymbol{x}}} \left( \boldsymbol{A} \boldsymbol{x}^{\diamond}(k) + \boldsymbol{B} \boldsymbol{u}(k) \right) \\ \boldsymbol{\zeta}^{\diamond}(k) &=& \diamondsuit_{\boldsymbol{m}_{\zeta}} \left( \begin{pmatrix} \boldsymbol{I} \\ \boldsymbol{C} \end{pmatrix} \boldsymbol{x}^{\diamond}(k) + \begin{pmatrix} \boldsymbol{0} \\ \boldsymbol{D} \end{pmatrix} \boldsymbol{u}(k) \right) \end{array} \right.$$

where  $\Diamond_m$  is some operator ensuring faithful rounding:

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

#### Taking the quantization errors into account

The actually implemented filter  $\mathcal{H}^{\Diamond}_{\mathcal{L}}$  is:

$$\mathcal{H}_{\zeta}^{\diamond} \left\{ \begin{array}{rcl} \boldsymbol{x}^{\diamond}(k+1) &=& \diamondsuit_{\boldsymbol{m}_{\boldsymbol{x}}} \left( \boldsymbol{A} \boldsymbol{x}^{\diamond}(k) + \boldsymbol{B} \boldsymbol{u}(k) \right) \\ \boldsymbol{\zeta}^{\diamond}(k) &=& \diamondsuit_{\boldsymbol{m}_{\boldsymbol{\zeta}}} \left( \begin{pmatrix} \boldsymbol{I} \\ \boldsymbol{C} \end{pmatrix} \boldsymbol{x}^{\diamond}(k) + \begin{pmatrix} \boldsymbol{0} \\ \boldsymbol{D} \end{pmatrix} \boldsymbol{u}(k) \right) \end{array} \right.$$

where  $\Diamond_m$  is some operator ensuring faithful rounding:

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

It holds

$$\mathcal{H}_{\zeta}^{\Diamond} \left\{ \begin{array}{rcl} \boldsymbol{x}^{\Diamond}(k+1) &=& \boldsymbol{A}\boldsymbol{x}^{\Diamond}(k) + \boldsymbol{B}\boldsymbol{u}(k) + \boldsymbol{\varepsilon}_{x}(k) \\ \boldsymbol{\zeta}^{\Diamond}(k) &=& \begin{pmatrix} \boldsymbol{I} \\ \boldsymbol{C} \end{pmatrix} \boldsymbol{x}^{\Diamond}(k) + \begin{pmatrix} \boldsymbol{0} \\ \boldsymbol{D} \end{pmatrix} \boldsymbol{u}(k) + \begin{pmatrix} \boldsymbol{\varepsilon}_{x}(k) \\ \boldsymbol{\varepsilon}_{y}(k) \end{pmatrix} \right.$$

with

$$|\boldsymbol{\varepsilon}_x(k)| \leq 2^{\boldsymbol{m}_x - \boldsymbol{w}_x + 1}$$
 and  $|\boldsymbol{\varepsilon}_y(k)| \leq 2^{\boldsymbol{m}_y - \boldsymbol{w}_y + 1}$ .

$$\begin{array}{c|c} u(k) & \mathcal{H}^{\Diamond}_{\zeta} & \stackrel{\zeta^{\Diamond}(k)}{\longrightarrow} \end{array}$$









! Filter  $\mathcal{H}_{\Delta}$  depends on the MSBs of filter  $\mathcal{H}_{\zeta}$ 



! Filter  $\mathcal{H}_{\Delta}$  depends on the MSBs of filter  $\mathcal{H}_{\zeta}$ 



! Filter  $\mathcal{H}_{\Delta}$  depends on the MSBs of filter  $\mathcal{H}_{\zeta}$ 

# Two step approach

- Step 1 Determine the MSBs  $m_{\zeta}$ for the exact filter  $\mathcal{H}_{\zeta}$ , applying the WCPG theorem;



 $\mathcal{H}_{\Delta}$ 

 $\Delta_{\zeta}(k)$ 

 $\begin{pmatrix} \varepsilon_x(k) \\ \varepsilon_y(k) \end{pmatrix}$ 

$$\boldsymbol{m}_{\zeta_{i}} = \left[\log_{2}\left(\left\langle\!\left\langle \mathcal{H}_{\zeta}\right\rangle\!\right\rangle \bar{\boldsymbol{u}}\right)_{i} + \log_{2}\left(1 - 2^{1 - \boldsymbol{w}_{i}}\right)\right]$$

$$\boldsymbol{m}_{\zeta_{i}} = \left\lceil \log_{2} \left( \langle\!\langle \mathcal{H}_{\zeta} \rangle\!\rangle \, \boldsymbol{\bar{u}} \right)_{i} + \log_{2} \left( 1 - 2^{1 - \boldsymbol{w}_{i}} \right) \right\rceil}$$
$$\widehat{\left\langle \langle \mathcal{H}_{\zeta} \rangle \right\rangle} + \boldsymbol{\varepsilon}_{WCPG}$$

$$\boldsymbol{m}_{\zeta_{i}} = \left\lceil \log_{2} \left( \langle\!\langle \mathcal{H}_{\zeta} \rangle\!\rangle \, \boldsymbol{\bar{u}} \right)_{i} + \log_{2} \left( 1 - 2^{1 - \boldsymbol{w}_{i}} \right) \right\rceil$$
$$\widehat{\left\langle \langle \mathcal{H}_{\zeta} \rangle \right\rangle} + \boldsymbol{\varepsilon}_{WCPG}$$

$$egin{aligned} \widehat{m{m}}_{\zeta_i} &= \left\lceil \log_2 \left( \left<\!\!\left< \mathcal{H}_{\zeta} \right>\!\!\right> \bar{m{u}} 
ight)_i + \log_2 \left( 1 - 2^{1-m{w}_i} 
ight) \\ &+ \log_2 \left( 1 + rac{\left(m{arepsilon}_{WCPG} \cdot ar{m{u}} 
ight)_i}{\left( \left<\!\!\left< \mathcal{H}_{\zeta} \right>\!\!\right> \bar{m{u}} 
ight)_i} 
ight) 
ight
ceil \end{aligned}$$

$$\begin{split} \boldsymbol{m}_{\zeta_{i}} &= \left\lceil \log_{2} \left( \langle\!\langle \mathcal{H}_{\zeta} \rangle\!\rangle \, \bar{\boldsymbol{u}} \right)_{i} + \log_{2} \left( 1 - 2^{1 - \boldsymbol{w}_{i}} \right) \right\rceil \\ & \widehat{\boldsymbol{w}}_{\zeta_{i}} = \left\lceil \log_{2} \left( \langle\!\langle \mathcal{H}_{\zeta} \rangle\!\rangle \, \bar{\boldsymbol{u}} \right)_{i} + \log_{2} \left( 1 - 2^{1 - \boldsymbol{w}_{i}} \right) \\ &+ \log_{2} \left( 1 + \frac{(\boldsymbol{\varepsilon}_{WCPG} \cdot \bar{\boldsymbol{u}})_{i}}{\left( \langle\!\langle \mathcal{H}_{\zeta} \rangle\!\rangle \, \bar{\boldsymbol{u}} \right)_{i}} \right) \right\rceil \\ &= \boldsymbol{m}_{\zeta_{i}} + [\ldots] \\ & \in \{0, 1\} \end{split}$$

$$\begin{split} \boldsymbol{m}_{\zeta_i} &= \left\lceil \log_2\left(\langle\!\langle \mathcal{H}_{\zeta} \rangle\!\rangle \, \bar{\boldsymbol{u}} \right)_i + \log_2\left(1 - 2^{1 - \boldsymbol{w}_i}\right) \right\rceil \\ &\overbrace{\langle\langle \mathcal{H}_{\zeta} \rangle\rangle} + \boldsymbol{\varepsilon}_{WCPG} \\ \widehat{\boldsymbol{m}}_{\zeta_i} &= \left\lceil \log_2\left(\langle\!\langle \mathcal{H}_{\zeta} \rangle\!\rangle \, \bar{\boldsymbol{u}} \right)_i + \log_2\left(1 - 2^{1 - \boldsymbol{w}_i}\right) \\ &+ \log_2\left(1 + \frac{(\boldsymbol{\varepsilon}_{WCPG} \cdot \bar{\boldsymbol{u}})_i}{\left(\langle\!\langle \mathcal{H}_{\zeta} \rangle\!\rangle \, \bar{\boldsymbol{u}} \right)_i}\right) \right\rceil \\ &= \boldsymbol{m}_{\zeta_i} + [\dots] \\ &\overbrace{\in \{0, 1\}} \end{split}$$

Adjust error term  $\varepsilon_{WCPG}$  in order to be at most off by one.

# Algorithm

- Step 1 Determine the MSBs  $m_{\zeta}$  for the exact filter  $\mathcal{H}_{\zeta}$ , applying the WCPG theorem;
- Step 2 Compute the error-filter  $\mathcal{H}_{\Delta}$ , induced by the format  $\boldsymbol{m}_{\zeta}$  and deduce the MSB  $\boldsymbol{m}_{\zeta}^{\Diamond}$  of the implemented filter;

$$\begin{array}{ll} \text{Step 3 If } \boldsymbol{m}_{\zeta_i}^{\Diamond} == \boldsymbol{m}_{\zeta_i} \text{ then return } \boldsymbol{m}_{\zeta_i}^{\Diamond} \\ \text{otherwise } \boldsymbol{m}_{\zeta_i} \leftarrow \boldsymbol{m}_{\zeta_i} + 1 \text{ and go} \\ \text{to Step 2.} \end{array}$$





# Numerical examples

Example:

- Random filter with 6 states, 1 input, 3 outputs
- $\bar{u} = 3.7776$ , wordlengths set to 7 bits

$$m_{\zeta} = (4, 4, 4, 4, 2, 3, 6, 5, 5)$$
  
 $m_{\zeta}^{\diamond} = (4, 5, 4, 4, 2, 3, 6, 5, 5)$ 

# Numerical examples

Example:

- Random filter with 6 states, 1 input, 3 outputs
- $\bar{u} = 3.7776$ , wordlengths set to 7 bits

$$m_{\zeta} = (4, 4, 4, 4, 2, 3, 6, 5, 5)$$
  
 $m_{\zeta}^{\Diamond} = (4, 5, 4, 4, 2, 3, 6, 5, 5)$   
After 3 iterations:  
 $m_{\zeta}^{\Diamond} = (4, 5, 5, 4, 3, 4, 6, 5, 5)$ 

### Numerical examples



# Conclusion and Perspectives

Conclusion

- Rigorous procedure for Fixed-Point Formats determination
- Filter computation errors are taken into account, ensuring that no overflow occurs
- Multiple-wordlength paradigm

Perspectives

- Integrate into optimization procedures in automatic workflow
- Solve off-by-one problem

Thank you! Questions?