# RELIABLE FIXED-POINT IMPLEMENTATION OF LINEAR DATA-FLOWS

Thibault HILAIRE, Anastasia VOLKOVA and Maminionja RAVOSON

Sorbonne Universités, UPMC Univ Paris 06, UMR 7606, LIP6, F-75005, Paris, France

## MOTIVATION

Implementation of Linear Time-Invariant systems in finite precision:

- Coefficient quantization, computational errors and their propagation
- Software/Hardware implementation under constraints
- Large variety of possible structures (e.g. Direct Forms, Lattices, etc.) described by different means (graphical, analytical)
- Most error analysis are based on statistical models ⇒ no guarantee on finite-precision implementation

## CONTRIBUTIONS

Automatic Fixed-Point code Generator for filters for which we:

- Use unified *analytical* representation of linear data-flows
- Can describe any analytical and graphical representation with our framework
- Adopted numerous classical and developed new quality measures
- Provide **fully rigorous** and **reliable** implementation in Fixed-Point arithmetic
- Generate C and VHDL code (for ASICs and FPGAs)

## PERSPECTIVES

We plan to:

- Perform a reliable verification of specification-to-implementation correspondance
- Optimize Software/Hardware implementation for various constraints using our rigorous approach
- Wrap the generator with optimization routines to create a complete and efficient filter-to-code tool

## SIF

**S**pecialized **I**mplicit **F**orm is an analytical matrix-based representation of input/output relationship:

$$\begin{pmatrix} J & 0 & 0 \\ -K & I_n & 0 \\ -L & 0 & I_p \end{pmatrix} \begin{pmatrix} t(k+1) \\ x(k+1) \\ y(k) \end{pmatrix} = \begin{pmatrix} 0 & M & N \\ 0 & P & Q \\ 0 & R & S \end{pmatrix} \begin{pmatrix} t(k) \\ x(k) \\ u(k) \end{pmatrix}$$

- $u(k)$ - inputs    ◇ $t(k)$ - temp. variables
- $y(k)$ - outputs   ◇ $x(k)$ - states

Some properties of SIF:

- Easy algebraic computations
- Can describe any linear data flow
- Unifies analysis and implementation

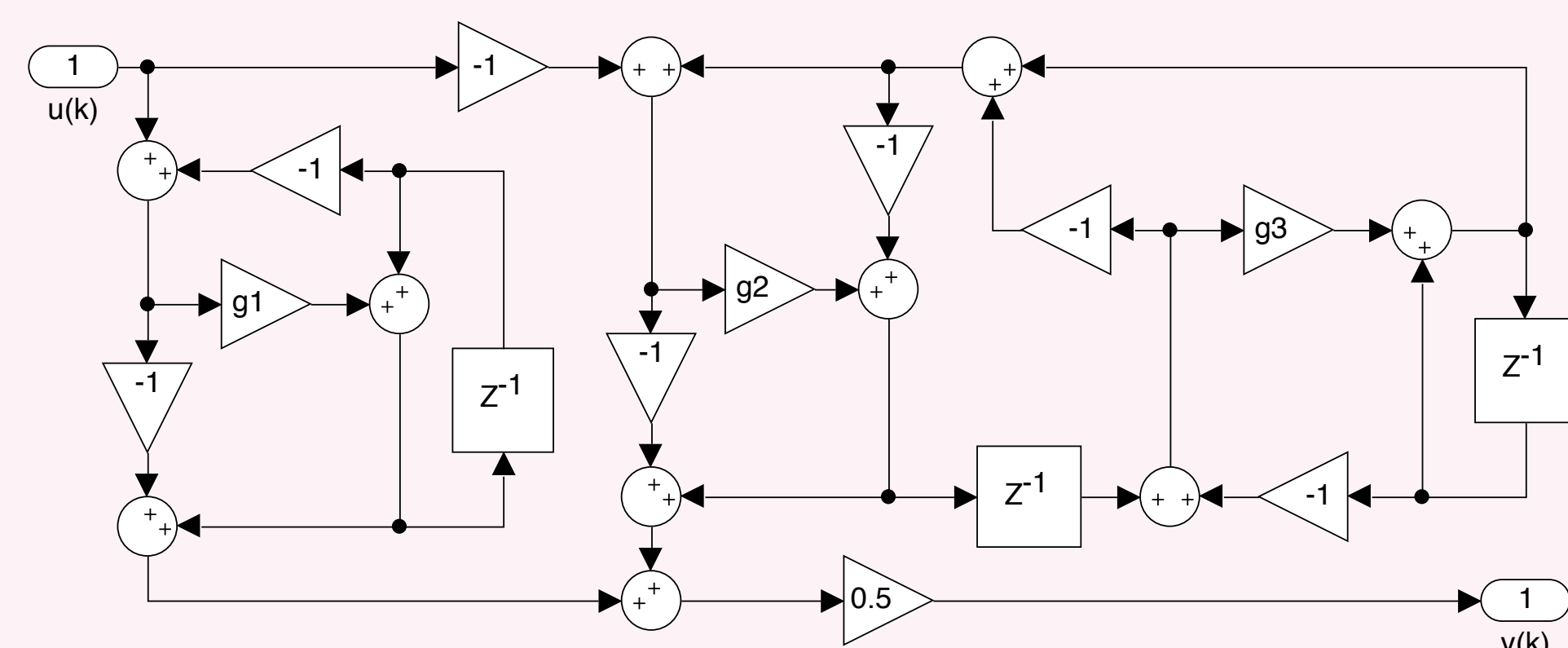**Important:** order of computations is preserved in the lower-triangular matrix $J$. For example,

$$y = M_2(M_1 x)$$

is described using temporary variable $t$ with

$$\begin{pmatrix} I & 0 \\ -M_2 & I \end{pmatrix} \begin{pmatrix} t \\ y \end{pmatrix} = \begin{pmatrix} M_1 \\ 0 \end{pmatrix} x$$

## LINEAR DATA FLOW



**EXAMPLE OF A LINEAR DATA FLOW: LATTICE WAVE DIGITAL FILTER**



## SIMULINK-TO-SIF

**Step 1:** Label variables
- states ⟵ delays
- SIF coefficients ⟵ gain elements
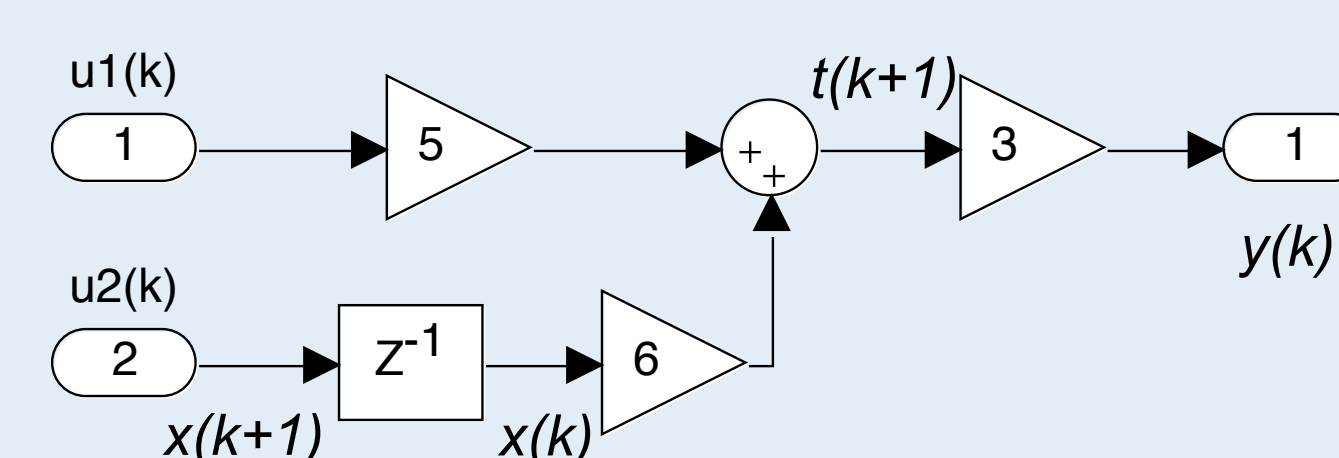- temporary variables ⟵ outputs of gain and sum operators

**Step 2:** Form the Sums-of-Products
- merge all directly cascading blocks
- regroup successive sums
- flatten design if subsystem is present

**Step 3:** Preserve the order of computations
- topological sort of temporary variables
- remove unnecessary variables

For example, a data-flow



Corresponds to SIF equations:

$$\begin{cases} t(k+1) & = 6 \cdot x(k) + 5 \cdot u_1(k) \\ x(k+1) & = 1 \cdot u_2(k) \\ y & = 3 \cdot t(k+1) \end{cases}$$

## QUALITY MEASURES

**Classical measures**

- Based on sensitivity wrt. the coefficients:
  → transfer function and pole/zero sensitivities
- Signal to Quantization Noise Ratio
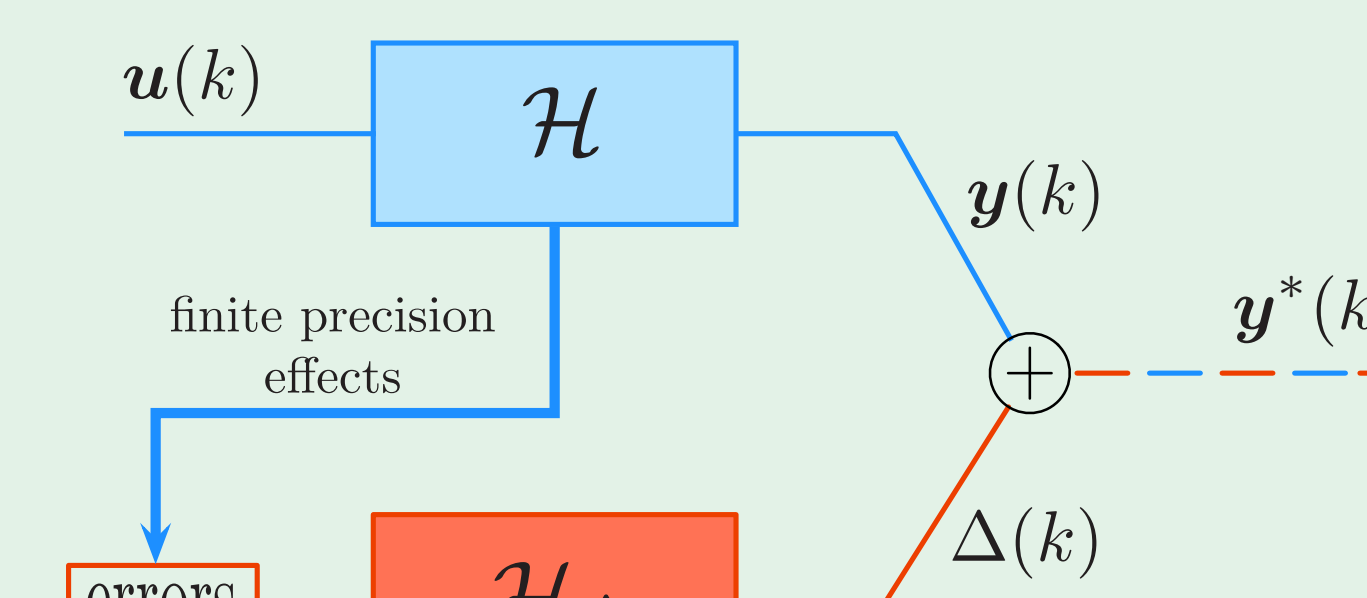  → errors modeled as noises ($\|\mathcal{H}_\Delta\|_2$)

**New measures: rigorous approach**

**Theorem 1.** *For a stable LTI system $\mathcal{H}$, if $\forall k \ |u(k)| \leq \bar{u}$, then the output is bounded*

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

*where $\langle\langle \mathcal{H} \rangle\rangle$ is the Worst-Case Peak Gain (WCPG) matrix of the system.*

**Theorem 2.** *The finite precision output $y^*(k)$ is bounded by the outputs of the exact and special error filters:* $|y^*(k)| \leq |y(k)| + |\Delta(k)|$



$\mathcal{H}_\Delta$ shows the propagation of finite precision effects.

## FIXED-POINT IMPLEMENTATION

**Most Significant Bit determination**
Reliable approach based on mathematical proofs:

**Step 1:** Determine output interval using our rigorous evaluation of the WCPG in arbitrary precision.

**Step 2:** Deduce the Fixed-Point implementation parameters while taking into account the propagation of computational errors through the filter.

**Important:** we prove that no overflow occurs and MSB positions are overestimated at most by one.

**Sum-of-Products**
The computations involves sums of products by constants. They can be performed with **faithful rounding** using $\lceil \log_2 N \rceil$ guard bits.

**Least Significant Bit determination**
The output error is analytically determined from the word-lengths $w$

$$\overline{\Delta} = 2 \langle\langle \mathcal{H}_\Delta \rangle\rangle \left( \lceil \langle\langle \mathcal{H} \rangle\rangle \bar{u} \rceil_2 \times 2^{-w} \right)$$

The word-length optimization problem can be solved with a Mixed-Integer Non Linear Programming solver.

## CODE GENERATION

Our tool can then generate code:
- C code for $\mu$Cs and DSPs
- VHDL for ASICs (using Stratus[a])
- VHDL for FPGAs (using FloPoCo[b])

---

[a] https://soc-extras.lip6.fr/en/coriolis/
[b] Flopoco: http://flopoco.gforge.inria.fr/