Project Course 2010	Floating-Point Arithmetic
Finite-Wordlength Implementation of Controllers	Hardware-supported on modern high-end processors (FPUs)
 Computer arithmetic Floating-point arithmetic Fixed-point arithmetic Controller realizations 	Number representation: $\pm f \times 2^{\pm e}$ • <i>f</i> : mantissa, significand, fraction • 2: base • <i>e</i> : exponent
 Based on material from K. J. Åström, B. Wittenmark: <i>Computer-Controlled Systems</i>, 3rd Ed., 1997. B. Lincoln: "Fixed-point controller implementation", Lecture in Embedded Systems, 2004. A. Cervin, KE. Årzén: "Implementation Aspects 2", Lecture in Real-Time Systems, 2008. 	The binary point is variable (floating) and depends on the value of the exponent Dynamic range and resolution Fixed number of significant digits
IEEE 754 Binary Floating-Point Standard	What is the output of this program?
IEEE 754 Binary Floating-Point Standard Used by almost all FPUs; implemented in software libraries	What is the output of this program? #include <stdio.h></stdio.h>
IEEE 754 Binary Floating-Point Standard Used by almost all FPUs; implemented in software libraries Single precision (Java/C float):	<pre>What is the output of this program? #include <stdio.h> main() {</stdio.h></pre>
 IEEE 754 Binary Floating-Point Standard Used by almost all FPUs; implemented in software libraries Single precision (Java/C float): 32-bit word divided into 1 sign bit, 8-bit biased exponent, and 23-bit mantissa (≈ 7 decimal digits) Range: 2⁻¹²⁶ - 2¹²⁸ 	<pre>What is the output of this program? #include <stdio.h> main() { float a[] = { 10000.0, 1.0, 10000.0 }; float b[] = { 10000.0, 1.0, -10000.0 }; float sum = 0.0; int i:</stdio.h></pre>
IEEE 754 Binary Floating-Point StandardUsed by almost all FPUs; implemented in software librariesSingle precision (Java/C float):32-bit word divided into 1 sign bit, 8-bit biased exponent, and 23-bit mantissa (\approx 7 decimal digits)Range: $2^{-126} - 2^{128}$ Double precision (Java/C double):	<pre>What is the output of this program? #include <stdio.h> main() { float a[] = { 10000.0, 1.0, 10000.0 }; float b[] = { 10000.0, 1.0, -10000.0 }; float sum = 0.0; int i; fam (i=0; i<2; i+1)</stdio.h></pre>
 IEEE 754 Binary Floating-Point Standard Used by almost all FPUs; implemented in software libraries Single precision (Java/C float): 32-bit word divided into 1 sign bit, 8-bit biased exponent, and 23-bit mantissa (≈ 7 decimal digits) Range: 2⁻¹²⁶ - 2¹²⁸ Double precision (Java/C double): 64-bit word divided into 1 sign bit, 11-bit biased exponent, and 52-bit mantissa (≈ 15 decimal digits) 	<pre>What is the output of this program? #include <stdio.h> main() { float a[] = { 10000.0, 1.0, 10000.0 }; float b[] = { 10000.0, 1.0, -10000.0 }; float sum = 0.0; int i; for (i=0; i<3; i++) sum += a[i]*b[i]; </stdio.h></pre>
 IEEE 754 Binary Floating-Point Standard Used by almost all FPUs; implemented in software libraries Single precision (Java/C float): 32-bit word divided into 1 sign bit, 8-bit biased exponent, and 23-bit mantissa (≈ 7 decimal digits) Range: 2⁻¹²⁶ - 2¹²⁸ Double precision (Java/C double): 64-bit word divided into 1 sign bit, 11-bit biased exponent, and 52-bit mantissa (≈ 15 decimal digits) Range: 2⁻¹⁰²² - 2¹⁰²⁴ 	<pre>What is the output of this program? #include <stdio.h> main() { float a[] = { 10000.0, 1.0, 10000.0 }; float b[] = { 10000.0, 1.0, -10000.0 }; float sum = 0.0; int i; for (i=0; i<3; i++) sum += a[i]*b[i]; printf("sum = %f\n", sum); }</stdio.h></pre>

1

Remarks: • The result depends on the order of the operations • Finite-wordlength operations are neither associative nor distributive	 Arithmetic in Embedded Systems Small microprocessors used in embedded systems typically do not have hardware support for floating-point arithmetic Options: Software emulation of floating-point arithmetic compiler/library supported large code size, slow Fixed-point arithmetic often manual implementation fast and compact
<pre>Fixed-Point Arithmetic Represent all numbers (parameters, variables) using integers Use binary scaling to make all numbers fit into one of the integer data types, e.g. 8 bits (char, int8_t): [-128, 127] 16 bits (short, int16_t): [-32768, 32767] 32 bits (long, int32_t): [-2147483648, 2147483647]</pre>	 Challenges Must select data types to get sufficient numerical precision Must know (or estimate) the minimum and maximum value of every variable in order to select appropriate scaling factors Must keep track of the scaling factors in all arithmetic operations Must handle potential arithmetic overflows



© Dept. of Automatic Control, Lund

• *n* too large \Rightarrow risk of overflow

127

-128

-127 :

-1

01111111

10000000

:



	1
Example: MultiplicationTwo numbers in Q5.2 format are multiplied: $x = 6.25 \Rightarrow X = 25$ $y = 4.75 \Rightarrow Y = 19$ Intermediate result: $X \cdot Y = 475$ Final result: $Z = 475/2^2 = 118 \Rightarrow z = 29.5$ (exact result is 29.6875)	$ \begin{array}{c} 0 & 0 & 0 & 1 & 1 & 0 & 0 & $
Multiplication of Operands with Different Q-format In general, multiplication of two fixed-point numbers $Qm_1.n_1$ and $Qm_2.n_2$ gives an intermediate result in the format $Qm_1+m_2.n_1+n_2$ which may then be right-shifted $n_1+n_2-n_3$ steps and stored in the format $Qm_3.n_3$	Implementation of Multiplication in CAssume Q4.3 operands and result#include <inttypes.h> /* define int8_t, etc. (Linux only) */ #define n 3 /* number of fractional bits */ int8_t X, Y, Z; /* Q4.3 operands and result */ int16_t temp; /* Q9.6 intermediate result */ temp = (int16_t)X * Y; /* cast operands to 16 bits and multiply */ temp = temp >> n; /* divide by 2^n */ Z = temp; /* truncate and assign result */</inttypes.h>
Common case: $n_2 = n_3 = 0$ (one real operand, one integer operand, and integer result). Then $Z = (X \cdot Y)/2^{n_1}$	

<pre>Implementation of Multiplication in C with Rounding and Saturation #include <inttypes.h> /* defines int8_t, etc. (Linux only) */ #define n 3 /* number of fractional bits */ int8_t X, Y, Z; /* Q4.3 operands and result */ int16_t temp; /* Q9.6 intermediate result */ temp = (int16_t)X * Y; /* cast operands to 16 bits and multiply */ temp = temp + (1 << n-1); /* add 1/2 to give correct rounding */ temp = temp >> n; /* divide by 2^n */ if (temp > INT8_MAX) /* saturate the result before assignment */ Z = INT8_MAX; else if (temp < INT8_MIN) Z = INT8_MIN; else Z = temp;</inttypes.h></pre>	<pre>Implementation of Division in C with Rounding #include <inttypes.h> /* define int8_t, etc. (Linux only) */ #define n 3 /* number of fractional bits */ int8_t X, Y, Z; /* Q4.3 operands and result */ int16_t temp; /* Q9.6 intermediate result */ temp = (int16_t)X << n; /* cast operand to 16 bits and shift */ temp = temp + (Y >> 1); /* Add Y/2 to give correct rounding */ temp = temp / Y; /* Perform the division (expensive!) */ Z = temp; /* Truncate and assign result */</inttypes.h></pre>
Realizations	Direct Form
A linear controller	The input-output form can be directly implemented as
$H(z)=rac{b_0+b_1z^{-1}+\ldots+b_nz^{-n}}{1+a_1z^{-1}+\ldots+a_nz^{-n}}$	$u(k)=\sum_{i=0}^n b_i y(k-i)-\sum_{i=1}^n a_i u(k-i)$
can be realized in a number of different ways with equivalent input- output behavior, e.g.	 Nonminimal (all old inputs and outputs are used as states) Very sensitive to roundoff in coefficients
Direct form	Avoid!
Companion (canonical) form	
Series (cascade) or parallel form	

Companion Forms Pole Sensitivity E.g. controllable or observable canonical form How sensitive are the poles to errors in the coefficients? Assume characteristic polynomial with distinct roots. Then $x(k+1) = \begin{pmatrix} -a_1 & -a_2 & \cdots & -a_{n-1} & -a_n \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \vdots & & & & \\ 0 & 0 & & 1 & 0 \end{pmatrix} x(k) + \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} y(k)$ $A(z) = 1 - \sum_{k=1}^n a_k z^{-k} = \prod_{j=1}^n (1 - p_j z^{-1})$ Pole sensitivity: $u(k) = \left(\begin{array}{cccc} b_1 & b_2 & \cdots & b_n\end{array}\right) x(k)$ $\frac{\partial p_i}{\partial a_k}$ • Same problem as for the Direct form • Very sensitive to roundoff in coefficients Avoid! **Better: Series and Parallel Forms** The chain rule gives Divide the transfer function of the controller into a number of first- $\frac{\partial A(z)}{\partial p_i} \frac{\partial p_i}{\partial q_i} = \frac{\partial A(z)}{\partial q_i}$ or second-order subsystems:

Evaluated in $z = p_i$ we get

$$rac{\partial p_i}{\partial a_k} = rac{p_i^{n-k}}{\prod_{j=1, j
eq i}^n (p_i - p_j)}$$

- Having poles close to each other is bad
- For stable filter, a_n is the most sensitive parameter



• Try to balance the gain such that each subsystem has about the same amplification



© Dept. of Automatic Control, Lund



Jackson's Rules for Series Realizations

How to pair and order the poles and zeros?

Jackson's rules (1970):

- Pair the pole closest to the unit circle with its closest zero. Repeat until all poles and zeros are taken.
- Order the filters in increasing or decreasing order based on the poles closeness to the unit circle.

This will push down high internal resonance peaks.

Well-Conditioned Parallel Realizations

Assume n_r distinct real poles and n_c distinct complex-pole pairs Modal (diagonal/parallel/coupled) form:

$$z_i(k+1) = \lambda_i z_i(k) + \beta_i y(k) \qquad \qquad i = 1, \dots, n_r$$

$$egin{aligned} v_i(k+1) &= egin{pmatrix} \sigma_i & \omega_i \ -\omega_i & \sigma_i \end{pmatrix} v_i(k) + egin{pmatrix} \gamma_{i1} \ \gamma_{i2} \end{pmatrix} y(k) & i=1,\ldots,n_c \ u(k) &= Dy(k) + \sum_{i=1}^{n_r} \gamma_i z_i(k) + \sum_{i=1}^{n_c} \delta_i^T v_i(k) \end{aligned}$$

Matlab: sysm = canon(sys, 'modal')

Multiple eigenvalues require Jordan blocks

Possible Pole Locations for Direct vs Modal Form



Short Sampling Interval Modification

In the state update equation

 $x(k+1) = \Phi x(k) + \Gamma y(k)$

the system matrix Φ will be close to *I* if *h* is small. Round-off errors in the coefficients of Φ can have drastic effects.

Better: use the modified equation

 $x(k+1) = x(k) + (\Phi - I)x(k) + \Gamma y(k)$

- Both ΦI and Γ are roughly proportional to h
 - Less round-off noise in the calculations
- Also known as the $\delta\text{-form}$

A Unifying Framework for Finite Wordlength Realizations

(Hilaire, Chevrel, Whidborn, IEEE Trans. Circuits & Systems, 2007)

Specialized implicit state-space form of a filter with input U, intermediate variable T, state X, output Y:

JT(k + 1) = MX(k) + NU(k) X(k + 1) = KT(k + 1) + PX(k) + QU(k)Y(k) = LT(k + 1) + RX(k) + SU(k)

Covers δ -form, direct forms, cascade/parallel forms, lattice filters, mixed $q/\delta, \ldots$

Sensitivity towards coefficient quantization can be analyzed

Short Sampling Interval and Integral Action

Fast sampling and slow integral action can give roundoff problems:

$$I(k+1) = I(k) + \underbrace{e(k) \cdot h/T_i}_{\approx 0}$$

Possible solutions:

- Use a dedicated high-resolution variable (e.g. 32 bits) for the I-part
- Update the I-part at a slower rate

General problem for filters with very different time constants

Pulse-Width Modulation (PWM)

Poor D-A resolution (e.g. 1 bit) can often be handled by fast switching between levels + low-pass filtering

The new control variable is the duty-cycle of the switched signal

