Matula Kommerup 22,20 July 28

1. Rationally biased rounding can be beneficial.

The larger gaps between representable values in floating slash representation occur only where one of the boundary values is a particularly simple rational. Rounding by the number theoretic concept of best rational approximation appropriately chooses the simple rational. For problems where simple rational values may occur and be meaningful, this feature could be very useful. E.g. linear programming and combinatorial optimization are two of many areas where simple rational output may often arise.

Two references providing more information on this feature are:

Matula, D.W. and Kornerup, P.: "Approximate Rational Arithmetic Systems: Analysis of Recovery of Simple Fractions During Expression Evaluation", Symbolic and Algebraic Computation, E.W. Ng, ed., Lecture Notes in Computer Science 72, Springer-Verlag, Berlin, 1979, 383-397.

Matula, D.W. and Ferguson, W.: "Rationally Biased Arithmetic", Proc. 7th Sym. on Comp. Arith., IEEE Cat #85CH2146-9, 1985, 194-202.

2. A precision fill feature can be added to floating slash representation.

The standard separate numerator and denominator word representation of rational fractions provides limitations in both range and precision features for many desired approximate real applications. Thus for approximation of reals in floating slash we must consider both a range problem and a precision problem.

Range problem: Our previously described extended range floating slash format allows scaled up numerators (implicit denominator unity) and scaled up denominators (implicit numerator unity) to achieve ranges comparable to typical floating point ranges. This provides a reasonably convenient solution to the range problem for rational representation.

Precision problem: We wish to point out a "precision fill" feature can be created in floating slash representation to achieve reasonably uniform relative spacing. This is possible by treating the presence of a non unit valued GCD as an indication that more accuracy is to be provided by interpreting the value of the GCD as giving leading bit information on the next partial quotient in the expansion of the approximate value. I.e. 1000/2000 in such a "denormalized" interpretation would have the GCD of 1000 indicate a small deviation from 1/2 whose precise meaning requires more discussion. We simply summarize here that the current redundancy of a simple fraction (GCD values available) for each fraction in floating slash is of just the right order to appropriately allow filling the gaps on both sides of the simple rational. "Appropriate" means that n bit floating slash can achieve accuracy of about say one part in $2^{(n-k)}$ for k a constant bounded by two or three in an appropriate development of this feature. This idea is currently under investigation by these authors as part of further developments in our investigation of floating slash arithmetic.

3. An alternative continued fraction based rational system is also available.

The use of a scaled version of k-bit LCF bitstrings provides an alternative to floating slash representation with a very similar and competitive arithmetic unit architecture.

Such systems have far less gap variation between k-bit representable values than in standard floating slash and still allow for exact representation of all appropriately simple rational values. We suggest that scaling to achieve a large range is a reasonably achievable extention to the unscaled k-bit LCF format. The unscaled LCF representation and corresponding arithmetic unit are described in the following two references:

Kornerup, P. and Matula, D.W.: "Finite Precision Lexicographic Continued Fraction Number Systems", Proc. 7th Sym. on Comp. Arith., IEEE Cat #85CH2146-9, 1985, 207-214.

"An On-Line Arithmetic Unit for Bit-Pipelined Rational Arithmetic", J. Parallel and Distributed Comp., 5, 1988, 310-330.

4. Further references.

For further information on floating slash representation we note the following two references:

Kornerup, P. and Matula, D.W.: "Finite Precision Rational Arithmetic: An Arithmetic Unit", IEEE Trns. on Comp., C-32, 1983, 378-388.

Matula, D.W. and Kornerup, P.: "Finite Precision Rational Arithmetic: Slash Number Systems", IEEE Trans. on Comp., 1985, 3-18.

> David W. Matula Peter Kornerup Aarhus, 22 July 1988

APPROXIMATE RATIONAL ARITHMETIC SYSTEMS: ANALYSIS OF RECOVERY OF SIMPLE FRACTIONS DURING EXPRESSION EVALUATION*

David W. Matula Peter Kornerup

Recovery of Exact Simple Fractions: A Computational Example

<u>Recovery of Exactness</u>: Evaluation in fixed-slash or floating-slash approximate rational arithmetic of a rational expression whose true result is a relatively simple fraction will generally yield that exact simple fraction as the final result even though intermediate results are rounded.

We morivate this feature by the following example.

Example: Table 1 illustrates the computation of the determinant

	10 13	20 17	13	
D = det	11 19	$\frac{7}{11}$	77 95	
	<u>69</u> 91	$\frac{4}{17}$	<u>56</u> 65	

= $\left(\left(\frac{10}{13} \cdot \frac{7}{11}\right) \frac{56}{65} + \left(\frac{11}{19} \cdot \frac{4}{17}\right) \frac{1}{13}\right) + \left(\frac{69}{91} \cdot \frac{20}{17}\right) \frac{77}{95}\right) - \left(\left(\frac{69}{91} \cdot \frac{7}{11}\right) \frac{1}{13} + \left(\frac{11}{19} \cdot \frac{20}{17}\right) \frac{56}{65}\right) + \left(\frac{10}{13} \cdot \frac{4}{17}\right) \frac{77}{95}\right) = \frac{5}{13}$ in the order indicated by the parenthesis. The approximate computation where all intermediate fractions are rounded by mediant rounding to the last convergent not exceeding 3 decimal digits in numerator or denominator (e.g. Φ_{ggg}) is shown in Table I

along with the absolute and relative errors accumulated at each stage. Note that the final step of the computation involves the rounding of $\frac{320}{277} - \frac{84}{109} = \frac{11612}{30193}$, whose convergents are $0, \frac{1}{2}, \frac{1}{3}, \frac{2}{5}, \frac{3}{8}, \frac{5}{13}, \frac{1288}{3349}, \frac{2581}{6711}, \frac{11612}{30193}$, so that $\phi_{999}(\frac{11612}{30193}) = \frac{5}{13}$ and the true result is recovered by the final rounding.

The number theoretic foundations of mediant rounding serve to explain our success in this example. From Theorem 1 (vi) and the definition of ϕ_n in (5) it is evident that the absolute error in rounding to $\frac{p_1}{q_1}$ is at most $1/q_1q_{1+1} \leq 1/q_1(n+1)$. From basic properties of Farey series and mediants [Hardy and Wright, 1960, Ch. III] it is further noted that the interval rounding to $\frac{p_1}{q_1}$ must extend at least $1/q_1(n+q_1)$ on either side of $\frac{p_1}{q_1}$. Thus the interval rounding to $\frac{5}{13}$ must extend at least $1/(13\times1012)=7.60\times10^{-5}$ on either side of 5/13, and the accumulated computation error had not grown sufficiently , to escape this relatively large interval rounding to the value 5/13 in our example.

	exact	computed value	accumulated abs. error	accumulated
$e_1 = \frac{10}{13} \cdot \frac{7}{11}$	70 143	$\frac{70}{143}$	0	0
$r_2 = r_1 \cdot \frac{56}{65}$	784 1859	<u>229</u> 543	9.9 · 10 ⁻⁷	$2.2 \cdot 10^{-6}$
$t_3 = \frac{11}{19} \cdot \frac{4}{17}$	<u>44</u> 323	<u>44</u> 323	0	0
$t_4 = t_3 \cdot \frac{1}{13}$	44 4199	7 668	3.5 · 10 ⁻⁷	3.4 · 10 ⁻⁵
$t_5 = \frac{69}{91} \cdot \frac{20}{17}$	<u>1380</u> 1547	<u>157</u> 176	3.7 - 10 ⁻⁶	4.1 · 10 ⁻⁶
$t_6 = t_5 \cdot \frac{77}{95}$	<u>3036</u> 4199	<u>449</u> 621	1.9 · 10 ⁻⁶	2.6 · 10 ⁻⁶
$t_7 = \frac{69}{91} \cdot \frac{7}{11}$	<u>69</u> 143	<u>69</u> 143	0	0
$\mathbf{r}_8 = \mathbf{r}_7 \cdot \frac{1}{13}$	<u>69</u> 1859	<u>17</u> 458	-1.17 · 10 ⁻⁶	3.2 · 10 ⁻⁵
$t_9 = \frac{11}{19} \cdot \frac{20}{17}$	220 323	220 323	0	0
$t_{10} = t_9 \cdot \frac{56}{65}$	2464 4199	<u>169</u> 288	8.2 • 10 ⁻⁷	1.4 · 10 ⁻⁶
$t_{11} = \frac{10}{13} \cdot \frac{4}{17}$	<u>40</u> 221	<u>40</u> 221	0	0
$t_{12} = t_{11} \cdot \frac{77}{95}$	<u>616</u> 4199	<u>109</u> 743	-9.6 - 10 ⁻⁷	6.6 · 10 ⁻⁶
$t_{13} = t_2 + t_4$	259524 600457	<u>51</u> 118	7.4 • 10 ^{-6 ·}	1.7 · 10 ⁻⁵
t ₁₄ = t ₁₃ + t ₆	<u>693672</u> 600457	<u>320</u> 277	5.4 • 10 ⁻⁶	$4.7 \cdot 10^{-6}$
$t_{15} = t_8 + t_{10}$	<u>374639</u> 600457	<u>73</u> 117	-8.5 • 10 ⁻⁶	1.3 • 10 ⁻⁵
$r_{16} = r_{15} + r_{12}$	<u>462727</u> 600457	<u>84</u> 109	$-1.7 \cdot 10^{-5}$	2.3 · 10 ⁻⁵
$D = c_{14} - c_{16}$	5 13	<u>5</u> 13	0	0

Table [. Illustration of the recovery of the exact simple fraction result in evaluating a particular determinant although intermediate approximate results were rounded to at most J digits in numerator and denominator.