

## On Errors and Error Bounds

The value of any arithmetic scheme is reflected in the quality of error bounds an analyst can construct for computations in that arithmetic. Examples that show that one system of arithmetic produces a smaller error in some specific computations are often irrelevant. In real life, we don't know the true answer, so we don't know the error. We may be able to get a bound on the error, either analytically or by interval arithmetic.

Since the bound depends on the style of arithmetic, we may prefer arithmetic that provides lower error bounds for equal analytic effort. Comparing error bounds is tricky, of course, since they reflect the skill of the analyst and the intended audience as well as properties of the arithmetic. That's why it's so gratifying when a single analysis like the one below graphically demonstrates the superiority of one system of arithmetic over another.

## The Last Example on Gradual Underflow?

W. Kahan suggested that Smith's algorithm for complex divide be analyzed under gradual underflow and flush to zero. My analysis surprised me in the strength of its argument for gradual underflow.

Hopefully we all agree that complex divide is a relevant useful operation. I hope that everyone can follow the error analysis below. We have two complex numbers p+qi and r+s. p, q, r, and s are all single precision EEC numbers, normalized or not, but r and s are not both zero. The usual formula for the quotient is

$$\left\{ \frac{p+q(\frac{s}{r})}{r+s} \right\} + \left\{ \frac{qr-sp}{r+s} \right\} i$$

but this familiar formula is prone to intermediate underflow and overflow and requires 3+, 3\*, and 2/ operations. Surprisingly, the computation can be rearranged with fewer operations, 3+, 3\*, and 3/, with less chance of intermediate over/underflow:

if  $|s| \leq |r|$  then compute

$$\tilde{z} := \left\{ \frac{p+q(\frac{s}{r})}{r+s} \right\} + \left\{ \frac{q-p(\frac{s}{r})}{r+s} \right\} i$$

else compute

$$\tilde{z} := \left\{ \frac{p(\frac{s}{r})+q}{r(\frac{s}{r})+s} \right\} + \left\{ \frac{q(\frac{s}{r})-p}{r(\frac{s}{r})+s} \right\} i$$

This formula dates back at least to E. Smith in 1963; you can find it in Knuth volume II, p. 195.

**CLAIM:** If no exception other than underflow or inexact is raised, the indicated formulas produce a computed complex result  $\tilde{z} \neq \tilde{z}'$  that differs from the correct result  $z$  by no more than a few units in the last place of  $|s|$ .

This claim is possible in EEC with gradual underflow but no comparably simple statement can be made in a system with flush to zero or UN symbols in place of gradual underflow. Note that the claim does not imply that both components of complex  $z$  are individually accurate to a few units in the last place.

The computation should be executed in normalizing mode. In Warning mode, Invalid Result may be raised on one of the final divides if  $p$  and  $q$  are tiny.

I won't analyze the entire computation; instead, let's just look at the denominator

$$D = r + s + \left( \frac{s}{r} \right)$$

De Bough

17 April 1980

in the case  $|s| \leq |r|$ .  $r$  must be normalized but  $s$  can be zero, denormalized or normalized. The model of arithmetic is

computed( $x \text{ op } y$ ) = ( $x \text{ op } y$ ) $(1 + \epsilon)$  +  $\delta$   
where  $|e| \leq \rho = 2^{-24}$  and  $|\delta| \leq \mu = 2^{-127}$  for flush to zero (FZ) and  $2^{-150}$  for gradual underflow (GU); further we may choose  $\epsilon$  and  $\delta$  so that  $\epsilon\delta = 0$ .

Denoting computed values by  $\sim$  we find

$$\frac{\tilde{s}}{\tilde{r}} = \frac{s}{r} (1 + \epsilon_1) + \delta_1.$$

This can't overflow; in fact  $s/r \leq 1$ .

$$\tilde{s} \cdot \left( \frac{\tilde{s}}{\tilde{r}} \right) = s \cdot \left( \frac{s}{r} \right) (1 + \epsilon_2) + \delta_2$$

This can't overflow either.

$$\tilde{D} = \tilde{r} + \tilde{s} \cdot \left( \frac{\tilde{s}}{\tilde{r}} \right) = \left\{ r + s \cdot \left( \frac{s}{r} \right) \right\} (1 + \epsilon_3).$$

A true add always occurs so no underflow is possible; overflow can happen if  $r$  and  $s$  are huge.

$$\tilde{D} = r + s \cdot \left( \frac{s}{r} \right) (1 + \epsilon_1)(1 + \epsilon_2)(1 + \epsilon_3) + s \cdot \delta_1 (1 + \epsilon_2)(1 + \epsilon_3) + \delta_2 (1 + \epsilon_3)$$

Thus the absolute error

$$|\tilde{D} - D| = \left| r\epsilon_3 + s \cdot \left( \frac{s}{r} \right) \left\{ (1 + \epsilon_1)(1 + \epsilon_2)(1 + \epsilon_3) - 1 \right\} + s \delta_1 (1 + \epsilon_2)(1 + \epsilon_3) + \delta_2 (1 + \epsilon_3) \right|$$

$$\leq |r|\rho + \left| \frac{s^2}{r} \right| \cdot 3\rho + |\delta_1|\rho + \mu$$

$$\leq \left( |r| + 3 \left| \frac{s^2}{r} \right| \right) \rho + \mu$$

and the relative error

$$\left| \frac{\tilde{D} - D}{D} \right| \leq \left( \frac{|r| + 3 \left| \frac{s^2}{r} \right|}{|r| + \left| \frac{s^2}{r} \right|} \right) \rho + \frac{\mu}{|r| + s \cdot \left( \frac{s}{r} \right)} \leq 2\rho + \frac{\mu}{|r|}$$

So the relative uncertainty in the denominator is bounded by:

	GU	FZ
$\rho$	$2^{-24}$	$2^{-24}$
minimum $ r $	$2^{-126}$	$2^{-126}$
$\mu$	$2^{-150}$	$2^{-127}$

Relative error bound,

$$\text{independent of } r \quad 3 \cdot 2^{-24}$$

In this analysis, with gradual underflow the uncertainty due to underflow is actually less than the uncertainty due to roundoff. With flush to zero, the uncertainty due to underflow overwhelms the uncertainty due to roundoff.

The error bound  $2\rho + \frac{\mu}{|r|}$  is realistic. For instance, take

	GU	FZ
$r$	$3 \cdot 2^{-126}$	$3 \cdot 2^{-126}$
$s$	$2^{-126}$	$2^{-126}$
$\frac{\tilde{s}}{r}$	$.AAAAAB_{16} 2^{-1}$	$.AAAAAB_{16} 2^{-1}$
$s \cdot \frac{\tilde{s}}{r}$	$.555556_{16} 2^{-126}$	0
$\tilde{s} + s \cdot \left( \frac{\tilde{s}}{r} \right) \approx \tilde{D}$	$3.555558_{16} 2^{-126}$	$3 \cdot 2^{-126}$
D correct	$(10/3) 2^{-126}$	$(10/3) 2^{-126}$
Relative error	$.8 \cdot 2^{-24}$	.1
Computed bound on error	$2.4 \cdot 2^{-24}$	.17