

## Underflow revisited

A. Cuyt\*, P. Kuterna\*\*, B. Verdonk\*\*\*, D. Verschaeren†

Department of Mathematics and Computer Science, Universiteit Antwerpen (UIA), Universiteitsplein 1, 2610 Wilrijk, Belgium  
 e-mail: {cuyt, kuterna, verdonk, vschaer}@uia.ua.ac.be

Received: January 2002 / Accepted: June 2002

**Abstract.** Underflow is a floating-point phenomenon. Although the use of gradual underflow as defended in [2] and [4] is now widespread, most numerical analysts may not be aware of the fact that several implementations of the same principle are in existence, leading to different behavior of code on different platforms, mainly with respect to exception signaling. We intend to thoroughly discuss the slight differences among these implementations. Examples will be taken from current hardware and from our own multiprecision software class library. Throughout the discussion the focus is on the analysis of the phenomenon and not on any implementation issues. Many programmers are also unaware of the fact that the IEEE 754 and 854 standards do not guarantee that a program will deliver identical results on all conforming systems. Of all the differences that can occur cross-platform, the underflow exception is just one.

### 1 Underflow and IEEE 754-854

We denote by  $\mathbb{F}(\beta, t, L, U)$  the set of normalized floating-point numbers  $\pm d_0.d_1 \dots d_{t-1} \times \beta^e$  where

$$0 \leq d_i \leq \beta - 1, \quad d_0 \neq 0, \quad L \leq e \leq U.$$

In order to be able to take care of overflow, the set is enlarged by the representation  $\pm 1.0 \dots 0 \times \beta^{U+1}$  for signed infinity. In order to take care of underflow

---

\* Research Director, FWO-Vlaanderen

\*\* Supported by an NOI-grant from the Universiteit Antwerpen (UIA)

\*\*\* Postdoctoral Fellow, FWO-Vlaanderen

† Research Fellow, IWT (Institute for Science and Technology)

the set is also enlarged by the denormalized numbers  $\pm 0.d_1 \dots d_{t-1} \times \beta^L$ , where  $\pm 0.0 \dots 0 \times \beta^L$  represents signed zero. We speak of denormalized numbers rather than unnormalized, because they cannot be normalized with the exponent range being bounded by  $[L, U]$ . The introduction of denormal numbers enables the implementation of the primary underflow mechanism of gradual underflow which, roughly speaking, rounds tiny floating-point numbers less than  $\beta^L$  in magnitude to denormal numbers, rather than to zero. The internal representation of these denormal numbers and signed zeroes is with exponent  $L - 1$ , flagging that the leading digit  $d_0 = 0$  while the actual exponent  $e = L$ .

Because the set of floating-point numbers is a discrete approximation of the real number set, each arithmetic operation introduces a rounding error. In round to nearest, the default rounding mode, the relative error made when approximating  $x * y$  by its nearest floating-point number  $\bigcirc(x * y)$ , where  $*$   $\in \{+, -, \times, /\}$ , is at most

$$\left| \frac{x * y - \bigcirc(x * y)}{x * y} \right| \leq \frac{\beta}{2} \beta^{-t} \quad x, y \in \mathbb{F}(\beta, t, L, U) \quad (1)$$

unless  $x * y$  overflows or underflows. Implementing  $x * y$  such that its machine version  $x \otimes y$  delivers  $\bigcirc(x * y)$  is called an exactly rounded implementation. When taking gradual underflow into account the error analysis has to be reformulated as

$$\begin{aligned} \bigcirc(x * y) &= (x * y)(1 + \epsilon) + \eta, \\ |\epsilon| &\leq \frac{\beta}{2} \beta^{-t}, \quad |\eta| \leq \left(\frac{\beta}{2} \beta^{-t}\right) \beta^L, \quad \epsilon \eta = 0 \end{aligned} \quad (2)$$

unless  $x * y$  overflows. Here  $\epsilon$  expresses the relative error that occurs when  $|\bigcirc(x * y)|$  is larger than or equal to the underflow threshold  $\beta^L$ . For denormalized results,  $\eta$  expresses the absolute error which is at most  $\frac{\beta}{2} \beta^{-t} \beta^L$  [4]. In round up, round down or trunc the relative bound of  $(\frac{\beta}{2} \beta^{-t})$  doubles to  $\beta^{-t+1}$  in (1), (2) and the bound on  $\eta$ . We also note that at most one of  $\epsilon$  and  $\eta$  is nonzero, in other words  $\epsilon \eta = 0$ , and that, moreover,  $\eta = 0$  whenever the arithmetic operation is either addition or subtraction [11]. It is easy to prove that the use of denormalized numbers allows us to represent all tiny results of the floating-point addition or subtraction operations without rounding error, in other words, exactly.

The underflow exception flag was introduced to signal that (1) is no longer valid and that  $\eta \neq 0$  has occurred in (2). The IEEE standards [6, 7] relax this condition in the sense that the underflow exception should be signaled “at least” when  $\eta \neq 0$ . A commentary [1] to the standard, however, encourages the stricter criterion for setting the underflow flag. That is, it should be set whenever the result delivered is different from what would be

delivered in a system with the same precision but with an unbounded (or large enough) exponent range. But it is the rule rather than the exception that many more underflow alarms go off.

The IEEE standards specify that underflow (in non-trapping mode) should be signaled at the occurrence of the following two correlated events. One is the creation of a tiny nonzero result in the interval  $] -\beta^L, +\beta^L[$ . The other is the loss of accuracy during the approximation of this tiny result, usually by a denormalized number. Tininess may be detected either after rounding or before rounding. Loss of accuracy may be detected as either a case of inexactness (the delivered result differs from the result computed with unbounded precision and unbounded exponent range) or a case of denormalization loss (the denormalized result differs from the result delivered with unbounded exponent range). Since the decision to denormalize is taken after rounding, only the following combinations of tininess and loss of accuracy can occur: tiny before rounding and inexact, tiny after rounding and inexact, tiny after rounding and denormalization loss. We now discuss these three possible ways to detect the underflow exception in accordance with the IEEE standards. Although several options are possible for the implementation of the underflow exception, the IEEE standard requires that underflow be detected in the same way for all operations in a single programming environment.

We introduce the notations `result_ext` for the exact result of the arithmetic operation  $x * y$  (unbounded precision and unbounded exponent range) and `result_tmp` for the normalized, rounded result (to  $t$   $\beta$ -digits) of  $x * y$  with unbounded exponent range. We let `result` denote the (possibly denormalized) floating-point result delivered. Then the following three slightly different implementations are consistent with the specifications of the IEEE standards for the underflow exception:

- (W) `result_ext` is tiny (this is before rounding) and cannot be delivered exactly to `result`, in other words,

$$\begin{aligned} |\text{result\_ext}| &< \beta^L \\ \text{result} &\neq \text{result\_ext} \end{aligned}$$

- (V) `result_tmp` is tiny (this is after rounding to  $t$   $\beta$ -digits) and different from `result_ext` or `result` or both, in other words,

$$\begin{aligned} |\text{result\_tmp}| &< \beta^L \\ \text{result} &\neq \text{result\_ext} \end{aligned}$$

- (U) `result_tmp` is tiny, possibly different from `result_ext` and has to be denormalized as in the previous case, but what is worse, in the process

of the denormalization, nonzero trailing  $\beta$ -digits are lost, implying the condition  $\eta \neq 0$  in (2),

$$\begin{aligned} |\text{result\_tmp}| &< \beta^L \\ \text{result} &\neq \text{result\_tmp}. \end{aligned}$$

It is clear that, if condition U is satisfied, condition V is also satisfied, because a denormalization loss implies inexactness. Moreover, if condition V is true, condition W is also true, because tininess after rounding implies tininess before rounding. The situation where case U does not occur while  $\text{result\_tmp}$  is tiny, is in fact not so alarming because here the tiny result can be denormalized without losing any nonzero  $\beta$ -digits. For instance, with  $t = 4$  and  $\beta = 2$ , the conditions for U-underflow are not satisfied in the case:

$$\begin{aligned} \text{result\_ext} &= 1.010111 \times 2^{L-1}, \\ \text{result\_tmp} &= 1.010 \times 2^{L-1}, \\ \text{result} &= 0.101 \times 2^L. \end{aligned} \tag{3}$$

whereas they are in the case:

$$\begin{aligned} \text{result\_ext} &= 1.00111 \times 2^{L-2}, \\ \text{result\_tmp} &= 1.010 \times 2^{L-2}, \\ \text{result} &= 0.010 \times 2^L. \end{aligned} \tag{4}$$

The above situations respectively occur in the multiplication  $x \otimes y$  with

$$\begin{aligned} x &= 1.100 \times \beta^{e(x)}, \\ y &= 1.101 \times \beta^{e(y)}, \\ e(x) + e(y) &= L - 1 \text{ or } L - 2. \end{aligned}$$

The difference between the two cases lies in the fact that in (3)  $\eta = 0$  whereas in (4)  $\eta \neq 0$ , the first case obeying (1) and the second case not. The difficulty in detecting pure U-underflow stems from the inexact flag which in most implementations does not allow one to distinguish between  $\text{result\_ext} \neq \text{result\_tmp}$  and  $\text{result\_tmp} \neq \text{result}$ . Observe that the inexact exception can be raised at several occasions during the computation of  $x \otimes y$ : the significand of the normalized  $\text{result\_ext}$  may contain more than  $t$  bits or, while both the significands of  $\text{result\_ext}$  and  $\text{result\_tmp}$  contain at most  $t$  bits, some trailing nonzero bits may be lost at the very end during denormalization.

In addition to signaling underflow at the occurrence of a U case, a V implementation also signals underflow when  $\text{result\_tmp}$  is tiny but does not suffer a denormalization loss: the inexact condition arises during rounding

but (1) is not violated. In addition to signaling underflow in the U and V cases, a W implementation even signals underflow when the delivered floating-point result, though inexact, is no longer tiny: the result delivered equals  $\pm\beta^L$  after rounding. Since all potentially dangerous underflow cases that complicate the error analysis should be flagged, the U implementation is the minimal implementation required. However, owing to the difficulty of implementing this scheme, the IEEE standard allows setting the underflow flag whenever the unrounded or final result is tiny and the infinitely precise result cannot be delivered exactly.

The IEEE 754-854 standards are currently up for revision and the underflow definition is one of the topics on the committee's revision list. There is a general consensus that a unique definition for underflow, rather than the current three possible definitions, is a good idea for implementation design and testing, but it does entail some kind of compromise. One proposal [9] is to choose criterion V as the unique definition for underflow: underflow shall be signaled when the result is tiny after rounding and inexact. However, it was subsequently observed by W. Kahan [10] that the V definition is inconvenient both to those looking for underflows signifying greater than normal rounding error (U underflow) and to those looking for subnormal results, exact or inexact, in order to get rid of them. In this respect, it is relevant to mention one unapproved proposal [8] which introduces, besides the underflow exception, a subnormal exception. This exception signals that the rounded result is tiny, irrespective of the inexactness of the result. A final consensus has not yet been reached at the time of writing.

## 2 Underflow signaling in some implementations

We discuss a W and a V hardware implementation, respectively, by SUN for their UltraSparc processors and by INTEL for their Pentium processors. At the end we also mention a multiprecision C++ class library developed at the University of Antwerp, which supports a U implementation. The consistency of each implementation was tested using a large set of test vectors which is a generalization of Coonen's [3] and Hough's [5] sets of test vectors and which is publicly available from [12, 13]. For each of the precisions (single, double, 64 bit extended on INTEL and 113 bit quadruple on SUN) our test set for multiplication contained 1152 cases of U underflow, an extra 176 cases of  $V \wedge \neg U$  underflow and an additional 64 cases of  $W \wedge \neg V$  underflow. Analogously the test set for division contained 286 cases of alarming U underflow and 51 cases of  $V \wedge \neg U$  underflow. One can show [12, 13] that  $W \wedge \neg V$  underflow cannot occur during division. It was already pointed out in the previous section that underflow does not occur in addition and subtraction when gradual underflow is supported.

## 2.1 SUN UltraSparc

Although SUN implemented V underflow signaling in their single and double precision SuperSparc processors, they have switched to a W implementation in their single and double precision UltraSparc. Both the V and the W implementations are easier and faster to handle than the detection of a pure U case. The implementation on the UltraSparc signals all W, V and U cases included in our test set, as is required for a proper W implementation, since U implies V, which in its turn implies W. We give an example of a test vector that signals underflow on the UltraSparc and not on the SuperSparc.

If we multiply  $x = 0.11 \dots 1 \times 2^L$  by  $y = 1.00 \dots 01 \times 2^0$  in any of the available precisions, then

$$\begin{aligned} \text{result\_ext} &= 0.11 \dots 1 \dots 111 \times 2^L && 2t - 1 \text{ bits} \\ &= 1.1 \dots 1 \dots 111 \times 2^{L-1}. \end{aligned}$$

When rounding the tiny `result_ext` to nearest or upward, one obtains `result_tmp = 1.0 × 2L` or the smallest normalized float. It is clear that the conditions for W underflow are satisfied while those for V underflow are not.

## 2.2 INTEL PC family

The INTEL processors are extended-based. The default working precision is  $t = 64$  and the default exponent range  $[L, U] = [-16382, 16383]$ . Single or double precision arithmetic can be mimicked by changing the precision control (often also referred to as rounding precision) to, respectively, 24 or 53. But this leads to a change in the underflow behavior (as described in Sect. 2.2.1) and to erroneous double rounding in some cases (as described in Sect. 2.2.2).

*2.2.1 Underflow strategy and precision control* In its extended precision the INTEL Pentium implements the V underflow strategy. However, it behaves like a U implementation in single and double precision as a result of the fact that the hardware is extended-based: only the precision of single ( $t = 24$ ) or double ( $t = 53$ ) is mimicked but the exponent range is still that of the extended precision (15 bits wide). Hence single or double precision tiny results are only recognized as tiny when stored from the extended precision register to memory. At that moment the single or double precision denormalization also takes place, resulting in U underflow detection.

We elaborate on this in more detail. Assume that the precision control is set to double ( $t = 53$ ) and consider for instance the two double precision

operands  $a = b = 1.0 \times 2^{-1022}$ , the smallest normal number in double precision. Now consider the C program statements

```
double a = b;
long double c = a * b;
double d = a * b;
```

Most numerical analysts will agree with the idea that the second statement consists of two operations, namely, a double precision multiplication and a conversion of the result to extended precision, and that the last statement is essentially one single operation, namely, the double precision product  $a \times b$  copied to the double precision variable  $d$ . So we expect to get

$$\begin{aligned} a \times b &= \text{double}(1.0 \times 2^{-2044}) = 0 && \text{underflow, inexact} \\ c &= 0 \\ d &= 0. \end{aligned}$$

This is not the case on an INTEL platform. Because the hardware is extended-based, the second statement actually behaves like most of us think the last statement does:  $a$  and  $b$  are copied to the extended precision registers of the INTEL FPU, the product is carried out in the extended precision register and rounded according to the precision of precision control (in this case  $t = 53$ ) but with extended precision exponent range. The result is then copied to the extended precision variable  $c$ . In the same way, the last statement is essentially a compound statement (something many do not realize):  $a$  and  $b$  are copied to the extended precision registers of the INTEL FPU, the product is again carried out in the extended precision register and rounded to the precision of precision control ( $t = 53$ ) with extended precision exponent range. The result is then converted to the double precision variable  $d$ . This conversion to double precision memory implies a reduction in the range of representable exponents (from  $[-16382, 16383]$  to  $[-1022, 1023]$ ) and hence we get

$$\begin{aligned} a \times b &= 1.0 \times 2^{-2044} \\ c &= 1.0 \times 2^{-2044} \\ d &= 0 && \text{underflow, inexact.} \end{aligned}$$

Even though the precision control is set to  $t = 53$ , the product  $a \times b$  does not underflow because the extended precision exponent range is large enough. Underflow and inexactness, which is due to denormalization loss, are only detected when the result of the product is stored in  $d$ .

Thus, for the evaluation of expressions with the precision control set to 24 or 53 bits, the INTEL actually uses a hybrid format with 15 bits for the

exponent (instead of the usual 8 for single or 11 for double) and 24 or 53 bits for the significand. While this, on the one hand, disturbs the semantics of arithmetic statements and neglects the formal model of single or double precision arithmetic, it can, on the other hand, deliver more accurate results as in the above example for  $c$ .

This also explains why, while the INTEL processors implement the V underflow strategy in extended precision, they behave like a U implementation when the precision control is set to  $t = 24$  or  $t = 53$ . To illustrate this, assume that the precision control is set to  $t = 53$  and consider the two double precision operands  $a = (1 + 2^{-52}) \times 2^{-1022}$  and  $b = 1.5 \times 2^{-1}$  of which the product is given by

$$a \times b = (1 + 2^{-1} + 2^{-52} + 2^{-53}) \times 2^{-1023}.$$

In pure double precision semantics, this is a  $V \wedge \neg U$  underflow case. When the INTEL processor executes the program statement

```
double d = a * b;
```

the value  $a \times b$  is first rounded to the precision of precision control (for our example  $t = 53$ ) with extended precision exponent range, yielding

$$(1 + 2^{-1} + 2^{-51}) \times 2^{-1023} \quad \text{inexact.}$$

This rounded result is then stored to the double precision variable  $d$  without unrounding it first: denormalization without denormalization loss takes place and

$$d = (2^{-1} + 2^{-2} + 2^{-52}) \times 2^{-1022}.$$

Since the store to double precision memory involves only tininess and no inexactness, the INTEL processors do not signal underflow for  $V \wedge \neg U$  cases when mimicking double precision. The same holds when the precision control is set to single.

**2.2.2 Erroneous double rounding** More remarkable is the way in which the floating-point unit deals with certain cases which, in pure single or double precision semantics, are U underflow cases. As an example, consider the multiplication of the single precision ( $t = 24$  and  $[L, U] = [-126, 127]$ ) operands  $x$  and  $y$  where

$$\begin{aligned} x &= 1.000\ 0000\ 0000\ 0000\ 0000\ 0001 \times 2^{-25}, \\ y &= 1.111\ 1111\ 1111\ 1111\ 1111\ 1111 \times 2^{-126}. \end{aligned}$$



For the computation of  $x \times y$  in single precision, the intermediate values `result_ext` and `result_tmp` defined above are given by

$$\begin{aligned} \text{result\_ext} &= 1.000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0111 \\ &\quad 1111\ 1111\ 1111\ 1111\ 1111 \times 2^{-126-24}, \quad (5) \\ \text{result\_tmp} &= 1.000\ 0000\ 0000\ 0000\ 0000\ 0000 \times 2^{-126-24}. \end{aligned}$$

Clearly, both `result_ext` and `result_tmp` are tiny. Moreover, the U conditions for underflow are satisfied since, in the process of denormalization, nonzero trailing bits will be lost. Double rounding will occur if erroneously `result_tmp` is denormalized rather than `result_ext`. Denormalizing `result_ext` means shifting the significand 24 bits to the right and adjusting the exponent accordingly, yielding  $L = -126$ . This unnormalized value then has to be rounded to single precision ( $t = 24$ ), resulting in

$$\text{result} = 0.000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0001 \times 2^{-126} \quad \text{underflow, inexact}$$

which is the correct floating-point approximation of  $x \times y$  given by (5). At the same time, both the inexact and underflow exceptions should be raised to signal denormalization loss.

When the computation of  $x \times y$  is carried out on an INTEL PC platform with precision control set to  $t = 24$ , erroneous double rounding occurs. Taking into account the hybrid format used by the INTEL when the precision control is set to  $t = 24$  (or  $t = 53$ ), the INTEL FPU first computes `result_tmp` in the extended precision register with  $t = 24$  and  $[L, U] = [-16382, 16383]$ , and apparently does not unround when storing the result to a single precision variable, delivering the doubly rounded value

$$\text{result} = 0.000\ 0000\ 0000\ 0000\ 0000\ 0000 \times 2^{-126} \quad \text{underflow, inexact.}$$

While the inexact and underflow exceptions are appropriately raised to signal denormalization loss, the returned value is not correct. It results from rounding to even the exact halfway value obtained by denormalizing `result_tmp` rather than `result_ext`. The IEEE standard, however, as it is formulated now, admits this deviation from the principle of exact rounding, for the sake of extended-based hardware platforms.

**2.2.3 The *MpIeee* class library** At <http://win-www.uia.ac.be/u/cant> a multiprecision, high performance and yet fully IEEE compliant class library can be found, with user definable precision  $t$  and base  $\beta = 2^k$  or  $10^j$ . Great care has been taken to support all IEEE features (basic operations including square root and IEEE specified remainder, exactly rounded decimal-to-binary and binary-to-decimal conversions, exactly rounded conversions between floating-point formats of different precisions including

the hardware precisions, round to integral value and conversions from and to integers) without performance penalty in comparison with other multi-precision libraries that offer arithmetical operators (instead of library calls) and a user-defined base (usually of the form  $2^k$  or  $10^j$ ). The library can be used for rather large precisions and supports an exponent range equal to the range of the C++ `long` integer type (which is currently often 32 bits), with the restriction that  $|L| < U$  in (1). In the library, U underflow signaling is implemented. The fact that, generally speaking, underflow will occur less often because of the wide exponent range, together with the fact that the unrounding of `result_tmp` is not so costly compared to the multiprecision basic operations, makes the cost of U underflow checking relatively small. The importance of the U underflow implementation in the context of a high precision library lies in the fact that a warning is only issued in case of a true potentially dangerous underflow, that is, in case of greater than normal rounding error.

## References

- [1] Cody, W.J., Coonen, J.T., Gay, D.M., Hanson, K., Hough, D., Kahan, W., Karpin-ski, R., Palmer, J., Ris, F.N., Stevenson, D.: A proposed radix- and word-length-independent standard for floating-point arithmetic. *IEEE Micro* **4**, 86–100 (1984)
- [2] Coonen, J.T.: Underflow and the denormalized numbers. *Comput. Math. Appl.* **14**, 75–87 (1981)
- [3] Coonen, J.T.: Contributions to a proposed standard for binary floating-point arithmetic. Ph.D. thesis. Berkeley: University of California 1984
- [4] Demmel, J.: Underflow and the reliability of numerical software. *SIAM J. Sci. Statist. Comput.* **5**, 887–919 (1984)
- [5] Hough, D.G. et al.: UCBTEST, a suite of programs for testing certain difficult cases of IEEE 754 floating-point arithmetic. Restricted public domain software from <http://netlib.bell-labs.com/netlib/fp/index.html>
- [6] IEEE: IEEE standard for binary floating-point arithmetic, ANSI/IEEE Standard 754-1985. New York: Institute of Electrical and Electronics Engineers 1985; reprinted in *SIGPLAN Notices* **22**, 9–25 (1987)
- [7] IEEE: IEEE standard for radix-independent floating-point arithmetic. ANSI/IEEE Standard 854-1987. New York: Institute of Electrical and Electronics Engineers 1987
- [8] IEEE 754R Revision Group: 754R change proposal: two-exception underflow. Available at <http://754r.ucbtest.org/proposals/underflow/underflow2.html>
- [9] IEEE 754R Revision Group: Minutes from 754R meeting 16 May 2001. Available at <http://grouper.ieee.org/groups/754/meeting-minutes/01-05-16.html>
- [10] IEEE 754R Revision Group: Minutes from 754R meeting 20 June 2001. Available at <http://grouper.ieee.org/groups/754/meeting-minutes/01-06-20.html>
- [11] Sun Microsystems Numerical computation guide. Revision A. Santa Clara, CA: Sun-Soft 1995
- [12] Verdonk, B., Cuyt, A., Verschaeren, D.: A precision and range independent tool for testing floating-point arithmetic. I: Basic operations, square root and remainder. *ACM Trans. Math. Software* **27**, 92–118 (2001)

- [13] Verdonk, B., Cuyt, A., Verschaeren, D.: A precision and range independent tool for testing floating-point arithmetic. II: Conversions. *ACM Trans. Math. Software* **27**, 119–140 (2001)