

**Submitter:** Fred Tydeman and Jim Thomas (USA)

**Submission Date:** 2009-05-04

**Previous version of paper:** N1303, N1317, N1352, N1378

**Related WG14 documents:** N1151, N1171

**Related WG21 documents:** N2798

**Subject:** xxx\_TRUE\_MIN macros for <float.h>

This is a rewrite based upon feedback from the committee.

### Changes to C1x

Add new bullets to **5.2.4.2.2 Characteristics of floating types <float.h>**

[bullet after FLT\_EVAL\_METHOD] the presence or absence of subnormal numbers is characterized by the values **FLT\_HAS\_SUBNORM**, **DBL\_HAS\_SUBNORM**, and **LDBL\_HAS\_SUBNORM**:

```
-1      indeterminable (footnote A)
0       absent (footnote B) (type does not allow subnormal numbers)
+1      present (type does allow subnormal numbers)
```

A) Characterization as indeterminable is intended if floating-point operations do not consistently interpret subnormal representations as zero, nor as nonzero.

B) Characterization as absent is intended if no floating-point operations produce subnormal results from non-subnormal inputs, even if the type format includes representations of subnormal numbers.

[bullet after FLT\_MIN] minimum positive floating-point number (footnote C)

```
FLT_TRUE_MIN      1E-37
DBL_TRUE_MIN      1E-37
LDBL_TRUE_MIN     1E-37
```

C) If the presence or absence of subnormal numbers is indeterminable, then the value is intended to be a positive number no greater than the minimum normalized positive number for the type.

[paragraph 15, example 2] Remove "normalized" from just before IEC 60559.

Add

```
FLT_TRUE_MIN      1.40129846E-45 // decimal constant
FLT_TRUE_MIN      0X1P-149F // hex constant
FLT_HAS_SUBNORM    +1
DBL_TRUE_MIN      4.9406564584124654E-324 // decimal constant
DBL_TRUE_MIN      0X1P-1074 // hex constant
DBL_HAS_SUBNORM    +1
```

after **FLT\_MIN** and **DBL\_MIN**.

### Words for Rationale:

[add to 5.2.4.2.2 section] Applications that depend on support for subnormal numbers (i.e., gradual underflow) can check, at translation time, if **xxx\_HAS\_SUBNORM** has the value +1, where

xxx is FLT, DBL, or LDBL. Applications that depend on subnormal numbers not arising in computations can check if xxx\_HAS\_SUBNORM has the value 0.

Today, most processors treat subnormal numbers according to IEC 60559, supporting xxx\_HAS\_SUBNORM = +1; or, they flush subnormal results (and perhaps also subnormal operands) to zero, supporting xxx\_HAS\_SUBNORM = 0. For processors that provide a control mode to do one or the other, the determination of xxx\_HAS\_SUBNORM depends on the startup setting of the control mode, because the standard does not provide means to change the control mode.

Implementations whose floating-point operations are not consistent in their treatment of subnormal representations should define xxx\_HAS\_SUBNORMAL to -1. Prior to IEEE 754-1985, processors treated subnormal representations in a great variety of ways, sometimes with bizarre and problematic results. See

<http://expertvoices.nsd.org/cornell-cs322/2008/03/26/the-old-man-of-floating-point-and-subnormal-numbers/> and <http://www.quadibloc.com/comp/cp0201.htm> for some examples.

The values of the smallest subnormal floating-point numbers (if supported) are typically, but not always, FLT\_MIN\*FLT\_EPSILON, DBL\_MIN\*DBL\_EPSILON, LDBL\_MIN\*LDBL\_EPSILON. The Motorola 68881 supports an 80-bit format whose minimum subnormal is a factor of 2 smaller than the formula would indicate, which was valid in IEEE 754-1985 (equivalent to IEC 60559), but not in IEEE 754-2008.