

浮動小数点演算プリミティブの実装基準（仮）

October 3, 2023

浮動小数点数のフォーマットは IEEE 単精度フォーマットに基づく。ただし、同等の機能が得られれば内部的に異なる表現を用いても構わない。IEEE 規格の +0 とノーマル数が表現でき、かつ、これらの浮動小数点値に対しては対応する実数値が定義されていなければならない。それ以外の浮動小数点値の定義／未定義および値の解釈は指定しないが、個々の浮動小数点値について対応する実数値が定義されているか否かを一意に定めねばならない。以下では、対応する実数値が定義されている浮動小数点数を「有効な浮動小数点数」と呼ぶ。

以下の説明では、 $\varepsilon = 2^{-126}$ とする。また、数学的関数と識別できるように、浮動小数点演算プリミティブはキャピタライズしたうえで“FADD”のようにタイプライタ体を用いて表記する。

fadd 有効な浮動小数点数 A, B の組が $-2^{127} < A, B, A+B < 2^{127}$ を満たすとき、FADD(A, B) は有効な浮動小数点数を返し、以下の条件を満たすこと。

$$|\text{FADD}(A, B) - (A + B)| < \max(|A| \cdot 2^{-23}, |B| \cdot 2^{-23}, |A + B| \cdot 2^{-23}, \varepsilon)$$

fsub 有効な浮動小数点数 A, B の組が $-2^{127} < A, B, A-B < 2^{127}$ を満たすとき、FSUB(A, B) は有効な浮動小数点数を返し、以下の条件を満たすこと。

$$|\text{FSUB}(A, B) - (A - B)| < \max(|A| \cdot 2^{-23}, |B| \cdot 2^{-23}, |A - B| \cdot 2^{-23}, \varepsilon)$$

fmul 有効な浮動小数点数 A, B の組が $-2^{127} < A, B, AB < 2^{127}$ を満たすとき、FMUL(A, B) は有効な浮動小数点数を返し、以下を満たすこと。

$$|\text{FMUL}(A, B) - (AB)| < \max(|AB| \cdot 2^{-22}, \varepsilon)$$

fdiv 有効な浮動小数点数 $A, B (B \neq 0)$ の組が $-2^{127} < A, B, \frac{A}{B} < 2^{127}$ を満たすとき、FDIV(A, B) は有効な浮動小数点数を返し、以下を満たすこと。

$$\left| \text{FDIV}(A, B) - \frac{A}{B} \right| < \max\left(\left| \frac{A}{B} \right| \cdot 2^{-20}, \varepsilon\right)$$

sqrt 有効な浮動小数点数 A が $0 \leq A < 2^{127}$ を満たすとき、SQRT(A) は有効な浮動小数点数を返し、以下の条件を満たすこと。

$$|\text{SQRT}(A) - \sqrt{A}| < \max(\sqrt{A} \cdot 2^{-20}, \varepsilon)$$

sin ある定数 $c (1 - 2^{-23} < c < 1 + 2^{-23})$ が存在して、 $-2^{127} < A < 2^{127}$ を満たす全ての有効な浮動小数点数 A に対して SIN(A) は有効な浮動小数点数を返し、以下の条件を満たすこと。

$$|\text{SIN}(A) - \sin(cA)| < \max(|\sin(cA)| \cdot 2^{-18}, \varepsilon)$$

cos ある定数 c ($1 - 2^{-23} < c < 1 + 2^{-23}$) が存在して、 $-2^{127} < A < 2^{127}$ を満たす全ての有効な浮動小数点数 A に対して $\text{COS}(A)$ は有効な浮動小数点数を返し、以下の条件を満たすこと。

$$|\text{COS}(A) - \cos(cA)| < \max(|\cos(cA)| \cdot 2^{-18}, \varepsilon)$$

atan $-2^{127} < A < 2^{127}$ を満たす全ての有効な浮動小数点数 A に対して $\text{ATAN}(A)$ は有効な浮動小数点数を返し、以下の条件を満たすこと。

$$|\text{ATAN}(A) - \arctan(A)| < \max(|\arctan(A)| \cdot 2^{-20}, \varepsilon)$$

fhalf $\text{FHALF}(A)$ は $\text{FMUL}(A, 0.5)$ が満たすべき基準を満たすこと。

fsqr $\text{FSQR}(A)$ は $\text{FMUL}(A, A)$ が満たすべき基準を満たすこと。

fabs $-2^{127} < A < 2^{127}$ を満たす全ての有効な浮動小数点数 A に対して $\text{FABS}(A)$ は有効な浮動小数点数を返し、 $\text{FABS}(A) = |A|$ を満たすこと。

fneg $-2^{127} < A < 2^{127}$ を満たす全ての有効な浮動小数点数 A に対して $\text{FNEG}(A)$ は有効な浮動小数点数を返し、 $\text{FNEG}(A) = -A$ を満たすこと。

fless 有効な浮動小数点数 A, B の組が $-2^{127} < A, B < 2^{127}$ を満たすとき、 $\text{FLESS}(A, B)$ は真偽値を返し、 $A < B$ と $\text{FLESS}(A) = \text{true}$ が同値であること。

fiszero 有効な浮動小数点数 A が $-2^{127} < A < 2^{127}$ を満たすとき、 $\text{FISZERO}(A)$ は真偽値を返し、 $A = 0$ と $\text{FISZERO}(A) = \text{true}$ が同値であること。

fispos $\text{FISPOS}(A)$ は $\text{FLESS}(0.0, A)$ が満たすべき基準を満たすこと。

fisneg $\text{FISNEG}(A)$ は $\text{FLESS}(A, 0.0)$ が満たすべき基準を満たすこと。

floor 有効な浮動小数点数 A が $-2^{127} < A < 2^{127}$ を満たすとき、 $\text{FLOOR}(A)$ は有効な浮動小数点数を返し、その値は整数で、かつ、 $\text{FLOOR}(A) \leq A < \text{FLOOR}(A) + 1$ が成立すること。

ftoi (int_of_float) 有効な浮動小数点数 A が $-2^{31} + 1 \leq A \leq 2^{31} - 1$ を満たすとき、 $\text{FTOI}(A)$ は 32bit 整数を返し、かつ、 $|I - A| < |\text{FTOI}(A) - A|$ を満たすような 32bit 整数 I が存在しないこと。

itof (float_of_int) 32bit 整数 I に対して $\text{ITOF}(A)$ は有効な浮動小数点数を返し、かつ、 $|A - I| < |\text{ITOF}(I) - I|$ を満たすような有効な浮動小数点数 A が存在しないこと。