

# 1 Integer

- Interpretation of binary representations

$$\begin{aligned} B2U_w(\vec{x}) &= \sum_{i=0}^{w-1} x_i 2^i & UMin_w &= 0 & UMax_w &= 2^w - 1 \\ B2T_w(\vec{x}) &= -x_{w-1} 2^{w-1} + \sum_{i=0}^{w-2} x_i 2^i & TMin_w &= -2^{w-1} & TMax_w &= 2^{w-1} - 1 \end{aligned}$$

- Transformation between unsigned and signed

$$\begin{aligned} T2U_w(x) &= \begin{cases} x & 0 \leq x \leq TMax_w \\ 2^w + x & TMin_w \leq x < 0 \end{cases} \\ U2T_w(u) &= \begin{cases} u & 0 \leq u \leq TMax_w \\ u - 2^w & TMax_w < u \leq UMax_w \end{cases} \end{aligned}$$

- Expansion

$$\begin{aligned} \vec{u} &= [u_{w-1}, \dots, u_0] \rightarrow \vec{u}' = [0, \dots, 0, u_{w-1}, \dots, u_0] & B2U_w(\vec{u}) &= B2U_{w'}(\vec{u}') \\ \vec{x} &= [x_{w-1}, \dots, x_0] \rightarrow \vec{x}' = [x_{w-1}, \dots, x_{w-1}, x_{w-1}, \dots, x_0] & B2T_w(\vec{x}) &= B2T_{w'}(\vec{x}') \end{aligned}$$

- Truncation

$$\begin{aligned} \vec{u} &= [u_{w-1}, \dots, u_0] \rightarrow \vec{u}' = [u_{k-1}, \dots, u_0] & B2U_k(\vec{u}') &= B2U_w(\vec{u}) \mod 2^k \\ \vec{x} &= [x_{w-1}, \dots, x_0] \rightarrow \vec{x}' = [x_{k-1}, \dots, x_0] & B2T_k(\vec{x}') &= U2T_k(B2U_w(\vec{x}) \mod 2^k) \end{aligned}$$

- Addition

$$\begin{aligned} x +_w^u y &= \begin{cases} x + y & 0 \leq x + y < 2^w \\ x + y - 2^w & 2^w \leq x + y \leq 2^{w+1} - 2 \end{cases} & 0 \leq x, y \leq 2^w - 1 \\ x +_w^t y &= \begin{cases} x + y + 2^w & -2^w \leq x + y < -2^{w-1} \\ x + y & -2^{w-1} \leq x + y \leq 2^{w-1} - 1 \\ x + y - 2^w & 2^{w-1} \leq x + y \leq 2^w - 2 \end{cases} & -2^{w-1} \leq x, y \leq 2^{w-1} - 1 \end{aligned}$$

- Negation

$$-^t_w x = \begin{cases} TMin_w & x = TMin_w \\ -x & TMin_w < x \leq TMax_w \end{cases}$$

- Multiplication

$$\begin{aligned} x *_w^u y &= (x \cdot y) \mod 2^w & u \ll k &= u *_w^u 2^k \\ x *_w^t y &= U2T_w((x \cdot y) \mod 2^w) & x \ll k &= x *_w^t 2^k \end{aligned}$$

- Division

$$\begin{aligned} x \gg k &= \lfloor x/2^k \rfloor \\ (x + (1 \ll k) - 1) \gg k &= \lceil x/2^k \rceil & x/y &\equiv \begin{cases} \lfloor x/y \rfloor & x \geq 0, y > 0 \\ \lceil x/y \rceil & x < 0, y > 0 \end{cases} \end{aligned}$$

## 2 Floating Point Number

- Definition  $V = (-1)^s \times M \times 2^E$ . 1-bit  $s$  encodes the sign  $s$ ,  $k$ -bit **exp** encodes the exponent  $E$ , and  $n$ -bit **frac** encodes the significand  $M$ . float: 1, 8, 23; double: 1, 11, 52.

**Normalized exp** is neither 0 or  $2^k - 1$  (all 1).

- $E = e - Bias$ , in which  $Bias = 2^{k-1} - 1$ .
- $M = 1.f_{n-1} \dots f_1 f_0$ .

**Denormalized exp** = 0

- $E = 1 - Bias$ .
- $M = 0.f_{n-1} \dots f_1 f_0$ .

**Infinity** **exp** =  $2^k - 1$  (all 1), **frac** = 0.

**NaN** **exp** =  $2^k - 1$  (all 1), **frac**  $\neq$  0.

- Use round to even.
- $x +^f y \equiv Round(x + y)$ . FP addition lacks associativity but features monotonicity:  $x + a \geq x + b$  if  $a \geq b$ .
- $x *^f y \equiv Round(x \times y)$ . FP Multiplication is not associative, and it does not distribute over addition. It features monotonicity:

$$a \geq b, c \geq 0 \Rightarrow a *^f c \geq b *^f c$$

$$a \geq b, c \leq 0 \Rightarrow a *^f c \leq b *^f c$$

$$a \neq NaN \Rightarrow a *^f a \geq 0$$

- Type conversions:
  - **int** to **float**: no overflow. Possible to be rounded.
  - **int/float** to **double**: precise conversion.
  - **double** to **float**: possible to overflow to infinity. Possible to be rounded.
  - **float/double** to **int**: round to 0. Possible to overflow.