

Latihan Soal Representasi Bilangan – Floating Point

- Sebuah representasi bilangan floating point menggunakan 5 bit, dengan pembagian 3 bit untuk eksponen dan 2 bit untuk fraction. Representasi ini tidak memiliki sign bit, sehingga hanya digunakan untuk merepresentasikan bilangan non negatif.

Hitunglah representasi bilangan untuk bilangan berikut ini:

- $9/32$
- $3/16$
- $15/2$

$K = 3$ Bias = 3. Tentukan batas max denormal = $2^{-bias+1} = 2^{-2} = 1/4$

- $9/32 = 0.01001 = 1.001 \times 2^{-2}$. $E = -2$. $exp = E + bias = -2 + 3 = 1$. $frac = 001$.

Pembulatan menjadi 00. Kode biner: **001 00 = $1.00 \times 2^{-2} = 1/4$**

- $3/16 = 0.0011 = 1.1 \times 2^{-3}$. $Exp = 0 \Rightarrow$ bentuk denormal, $E = -bias + 1 = -2$. $1.1 \times 2^{-3} = 0.11 \times 2^{-2}$. $Exp = 000$, $frac = 11$.

Kode biner: **000 11**

- $15/2 = 111.1 = 1.111 \times 2^2$.

$Exp = 2 + 3 = 5$. $Frac = 1.111$, dibulatkan ke bawah jadi 1.11, dibulatkan ke atas jadi 10.00, jadi akan dibulatkan ke genap terdekat, yaitu 10.00

pembulatan menjadi $10.00 \times 2^2 = 1.00 \times 2^3$. $Exp = 3 + 3 = 6 = 110$, $frac = 00$

Kode biner: **110 00**

- Jelaskan eksekusi kode C berikut (x adalah int, f adalah float, d adalah double):

- `x == (int)(float) x // nilai x mungkin mengalami pembulatan`

0100 0000 0000 1110 0100 1000 1001 0000

- `x == (int)(double) x // nilai x tetap`

- `f == (float)(double) f // nilai f tetap`

- `d == (float) d // nilai d mungkin mengalami pembulatan`

- diberikan a dan b adalah int (32 bit), dengan representasi two complements (signed).

MIN_INT adalah minimum integer, dan MAX_INT adalah maksimum integer.

Pasangkanlah bagian sebelah kiri dengan pasangan yang sesuai di sebelah kanan pada tabel berikut:

1. Komplemen dari a	a. $\sim(\sim a \mid (b \wedge (\text{MIN_INT} + \text{MAX_INT})))$
2. a	b. $((a \wedge b) \& \sim b) \mid (\sim(a \wedge b) \& b)$
3. a & b	c. $1 + (a \ll 3) + \sim a$
4. a * 7	d. $(a \ll 4) + (a \ll 2) + (a \ll 1)$
5. a/4	e. $((a < 0) ? (a + 3) : a) \gg 2$
6. $(a < 0) ? 1 : -1$	f. $a \wedge (\text{MIN_INT} + \text{MAX_INT})$
	g. $\sim((a \mid (\sim a + 1)) \gg 31) \& 1$
	h. $\sim((a \gg 31) \ll 1)$
	i. $a \gg 2$

- (a). $b^{(\text{MIN_INT}+\text{MAX_INT})} = \sim b$. $\sim(\sim a | \sim b) = \sim(\sim(a \& b)) = a \& b$
 $1010\ 1101 \wedge 1111\ 1111 = 0101\ 0010 = \sim(1010\ 1101)$
- (b). $(x \& \sim y) | (\sim x \& y) = x \wedge y$. $((a \wedge b) \& \sim b) | (\sim(a \wedge b) \& b) = (a \wedge b) \wedge b = a$
- (c). $-a = \sim a + 1$. $\sim a = -a + 1$. $1 + (a < 3) + \sim a = 1 + \sim a + (a < 3) = -a + (a < 3) = a * 8 - a = a * 7$.
- (f) $\text{MIN_INT} + \text{MAX_INT} = 10000000...00 + 01111..11 = 1111111111111...$
 $a \wedge 111111111111.. = \sim a$
- (h) $a \gg 31$ bernilai -1 (111...11) jika $a < 0$, dan 0 jika $a \geq 0$.
 $(a \gg 31) < 1$ bernilai -2 (11..110) jika $a < 0$, dan 0 jika $a \geq 0$.
 $\sim((a \gg 31) < 1)$ bernilai 1 (00..001) jika $a < 0$, dan -1 (11..111) jika $a \geq 0$.
- (i) $a \gg 2 = a/4$

1 == f

2 == b

3 == a

4 == c

5 == i

6 == h

4. Diberikan representasi bilangan floating point dengan 8 bit, dengan pembagian: 1 bit sign, 3 bit exponent dan 4 bit fractions, menggunakan standar floating point IEEE.

Lengkapilah tabel berikut:

Bias = $2^{3-1} - 1 = 3$. $E = \text{exp} - \text{bias}$. $E = 0 = \text{exp} - 3 \Rightarrow \text{exp} = 3$.

Deskripsi	Biner	Nilai
Minus zero	1 000 0000	-0.0
Smallest denormalized (negative) nearest to zero	1 000 0001	$-1 * 2^{-2} * 2^{-4} = -2^{-6}$
Largest normalized (positive)	0 110 1111	$1.1111 * 2^3 = 1\ 15/16$ $* 8 = 15.5$
1	0 011 0000	1
-	0 101 0110	$5.5 = 101.1 \Rightarrow 1.011 * 2^2$
Positive infinity	0 111 0000	-