

## Floating-point Instruction Set

Sl. no	Instructions	Short Description	Corresponding Fixed point instr	Ref. pg
<b>ALU</b>				
1	$F_n = F_x + F_y$	Adds the floating-point operands in $F_x$ and $F_y$ .	$R_n = R_x + R_y$	B-24
2	$F_n = F_x - F_y$	Subtracts the floating-point operand in $F_y$ from that in $F_x$ .	$R_n = R_x - R_y$	B-25
3	COMP( $F_x, F_y$ )	Compares the floating-point operand in $F_x$ with that in register $F_y$ .	COMP( $R_x, R_y$ )	B-29
4	$F_n = -F_x$	Complements the sign bit of the floating-point operand in $F_x$ .	$R_n = -R_x$	B-30
5	$F_n = \text{ABS } F_x$	Returns the absolute value of the floating-point operand in $F_x$ by setting the sign bit to 0.	$R_n = \text{ABS } R_x$	B-31
6	$F_n = \text{MIN}(F_x, F_y)$	Returns the smaller of the floating-point operands in $F_x$ and $F_y$ .	$R_n = \text{MIN}(R_x, R_y)$	B-42
7	$F_n = \text{MAX}(F_x, F_y)$	Returns the larger of the floating-point operands in $F_x$ and $F_y$ .	$R_n = \text{MAX}(R_x, R_y)$	B-43
8	$F_n = \text{PASS } F_x$	Passes the floating-point operand in $F_x$ through the ALU to the floating-point field in register $F_n$ .	$R_n = \text{PASS } R_x$	B-32
9	$F_n = (F_x + F_y)/2$	Adds the floating-point operands in $F_x$ and $F_y$ and divides the result by 2, by decrementing the exponent of the sum before rounding.	$R_n = (R_x + R_y)/2$	B-28
10	$F_n = \text{CLIP } F_x \text{ BY } F_y$	Returns the floating-point operand in $F_x$ if the absolute value of the operand in $F_x$ is less than the absolute value of the floating-point operand in $F_y$ . Else, returns $ F_y $ if $F_x$ is positive, and $- F_y $ if $F_x$ is negative.	$R_n = \text{CLIP } R_x \text{ BY } R_y$	B-44
11	$F_n = \text{ABS } (F_x + F_y)$	Adds the floating-point operands in $F_x$ and $F_y$ , and places the absolute value of the normalized result in register $F_n$ .	-	B-26
12	$F_n = \text{ABS } (F_x - F_y)$	Subtracts the floating-point operand in $F_y$ from that in $F_x$ and places the absolute value of the normalized result in register $F_n$ .	-	B-27
13	$F_n = \text{RND } F_x$	Rounds the floating-point operand in $F_x$ .	-	B-33

14	<b>Fn = SCALB Fx BY Ry</b>	Scales the exponent of the floating-point operand in Fx by adding to it the fixed-point twos-complement integer in Ry.	-	<b>B-34</b>
15	<b>Rn = MANT Fx</b>	Extracts the mantissa (fraction bits with the hidden bit, excluding the sign bit) from the floating-point operand in Fx.	-	<b>B-35</b>
16	<b>Rn = LOGB Fx</b>	Converts the exponent of the floating-point operand in Fx to an unbiased twos-complement fixed-point integer.	-	<b>B-36</b>
17	<b>Rn = FIX Fx BY Ry</b>	Converts the floating-point operand in Fx to a twos-complement 32-bit fixed-point integer result. For the scaling factor (Ry), the fixed-point twos-complement integer in Ry is added to the exponent of the floating-point operand in Fx before the conversion.	-	<b>B-37</b>
18	<b>Rn = FIX Fx</b>	Converts the floating-point operand in Fx to a twos-complement 32-bit fixed-point integer result.	-	<b>B-37</b>
19	<b>Rn = TRUNC Fx BY Ry</b>	Converts the floating-point operand in Fx to a twos-complement 32-bit fixed-point integer result. The TRUNC operation always truncates toward 0. For the scaling factor (Ry), the fixed-point twos-complement integer in Ry is added to the exponent of the floating-point operand in Fx before the conversion.	-	<b>B-37</b>
20	<b>Rn = TRUNC Fx</b>	Converts the floating-point operand in Fx to a twos-complement 32-bit fixed-point integer result. The TRUNC operation always truncates toward 0.	-	<b>B-37</b>
21	<b>Fn = FLOAT Rx BY Ry</b>	Converts the fixed-point operand in Rx to a floating-point result. For the scaling factor (Ry), the fixed-point twos-complement integer in Ry is added to the exponent of the floating-point result.	-	<b>B-38</b>
22	<b>Fn = FLOAT Rx</b>	Converts the fixed-point operand in Rx to a floating-point result.	-	<b>B-38</b>
23	<b>Fn = RECIPS Fx</b>	Creates an 8-bit accurate seed for 1/Fx, the reciprocal of Fx.	-	<b>B-39</b>
24	<b>Fn = RSQRTS Fx</b>	Creates a 4-bit accurate seed for 1/√Fx, the reciprocal square root of Fx.	-	<b>B-40</b>
25	<b>Fn = Fx COPYSIGN Fy</b>	Copies the sign of the floating-point operand in Fy to the floating-point operand from Fx without changing the exponent or the mantissa.	-	<b>B-41</b>

MULTIPLIER				
26	$F_n = F_x * F_y$	Multiplies the floating-point operands in registers $F_x$ and $F_y$ .	$R_n = R_x * R_y$	B-53
SHIFTER				
27	$R_n = \text{FPACK } F_x$	Converts the IEEE 32-bit floating-point value in $F_x$ to a 16-bit floating-point value stored in $R_n$ .	-	B-74
28	$F_n = \text{FUNPACK } R_x$	Converts the 16-bit floating-point value in $R_x$ to an IEEE 32-bit floating-point value stored in $F_x$ .	-	B-75

- ☐ - Floating-point *Specific* instructions
- ☐ - Floating-point instructions in *correspondence* with Fixed-point instructions in SHARC (*present* in MISA 1 instruction set)
- ☒ - Floating-point instructions in *correspondence* with Fixed-point instructions in SHARC (*Not present* in MISA 1 instruction set)