11713020 Due on April 9 23:55pm **3.27**, 3.29, 3.30. 3.27 [20] <\$3.5> IEEE 754-2008 contains a half precision that is only 16 bits wide. The leftmost bit is still the sign bit, the exponent is 5 bits wide and has a bias of 15, and the mantissa is 10 bits long. A hidden 1 is assumed. Write down the bit pattern to represent -1.5625×10^{-1} assuming a version of this format, which uses an excess-16 format to store the exponent. Comment on how the range and A: 1.10000 X24 B = 1.10/00/ X2-2 accuracy of this 16-bit floating point format compares to the single precision IEEE 754 standard. 1, Align binary points 5 4 3 4 3 4 4 5 5 5 5 5 5 5 5 1.10/0/00/X2-2 = 0.00000/10/0/00/ X24 signed exponent bit bits mantissa (fraction) 2. Add significands 1.1010001 x 24 + 0.00000110101001 X 24 = 1.10101000101001 X24 (bias 15) (10 bits) Decinal: -1.5625 X10-1 = -0.15625 (-1)'=-) the signed bit is | 3. Normalize result & check for over / underflow · 1.10101000101001 X 24, with no overflow/underflow. $(0.15625)_{10} = (0.00|0|)_{2}$ Binam: 0,00/0 = 1.01 x 23 4. Round and renormalize if necessary The exponent is-3, as the bias is 15, the exponent bits is 0/100 Grand 15 0, R5 15 10, |.|0|0|000|0|00|
Guard 11 = Striky Bat
Round Bit And the mantissa is 0/00 0000 00 after sound to meanest even => The reguired bit pattern is 101100 0100000000 finally get . 1.1010100010 and it's mimalized. As the exponent bits 00000 and 11111 reserved. So the result is 1.1010100010 X24 0/00///0/0/000/0 1° smallest value Exponent bits: 00001. Actual exponent is (1-15) = -14 Montissa: 0000000000 => Significand: 1.0 Decimal: -8.0546875X10° = -8.0546875. => ±1.0×2-14 2=16.1×10-5 (-8.0546875), = (-1000,0000 111)2 Binary: -1000,0000|| = -1.0000000||| X 13 2° largest value Exponent bit: 11110. Actual exponent is (30-15) = 15 Decimal: -1.7981640625×10-1 = -0.179931640625 Martissa: 111111111 => significand & 2.0 (-0.17913/640625), = (-0.00/011/00001). => ± 2.0 X 2 = 46.5536 X 104 Binary: -0.000|0|110000| =-1.0|110000| X2-3 1.01110000 0 => Accuracy: $\frac{\Delta A}{|A|} = \frac{2^{-h} \times 2^{eq}}{1 \times 2^{eq}} = 2^{-h}$ X 1.0000000111 Step 1. Adding the exponents without bias: 10111000010 when equivalent to 10 lg 2 & 3 decimal digits of precision. 3+(-3)=0 10111000010 Step 2. Multiplying the significands: 10111000010 1.01110011000101001110 16-bit floating point format ±1.0X2-4 & t6.1 X10-5 ±1.0x2726 × ±1.2X/0-38 The product is 1.01110011000101001110 & ± 2,0x2 127 € ± 3.4 X/038 x ±2.0x2 = ± 65536 We need to keep to 12 bits (10 bits + 1 round bit + 1 stilly bit) 2-10/3 decimal digits 1.01110011000101001110 GALS Me result is 1.011100110001 X2 26 decimal digits Acchacy **3.29** [20] <\$3.5> Calculate the sum of 2.6125×10^{1} and $4.150390625 \times 10^{-1}$ Step 3. The result is normalized and the exponent, which is o. by hand, assuming A and B are stored in the 16-bit half precision described in Exercise 3.27. Assume 1 guard, 1 round bit, and 1 sticky bit, and round to the as not overflow or underflow. nearest even. Show all the steps. Step 4. Rounding the product Decimal: 2.6125 X/0 = 26.125 >0, the signed bit is 0 RS=01, Trucate the result by discarding RS (26.125) = (11010.001)2 The result is 1.0111001100 X20 the mantissa is 10/000/000 Binay: 100000 = 1.100001 X24 Step 5. Since the signs of the original operands are some. The exponent is 4, as the bias is 15, the exponent bit is 10011 make the sign of the product positive. The 16-bit had precion of A = 0 poll 1010001000 Hove, the product is 1.0111001100 X12° Decimal: 4.150390625 X107 = 0.4150390625 >0 signed bit is 0 Converting to decimal to check our result: (04150390625) = (0,0110101001)2 $(1.01|00|00)_{2} \times 2^{0} = (1.01|00|1)_{2}$ Binary: 0.01/0/0100 = 1.10/0/00/X2-2 the montissa is 10/0/00/00 $= 1 + 115/2^{8} = 1.44921875$ The exponent is -2, as the bias is 15, the exponent bit is 0/10/ The according of my result is 2-10, The 16-bit but precision of B = 00/10/10/00/00 which is about 3 decimal objects. Compared swith the result contributed by calculator, the difference is less than 0.0001.