

This is an open-book examination.  
Please write clearly any assumptions you make.

---

1. The IEEE 754-2008 standard specifies the result of  $\text{maxnum}(x,y)$  as the number  $y$  if  $x < y$ ,  $x$  if  $y < x$ , the number if one operand is a number and the other a quiet NaN. Otherwise it is either  $x$  or  $y$  (this means results might differ among implementations). When either  $x$  or  $y$  is a signaling NaN, then the result is a quiet NaN and the operation signals invalid (under default exception handling).

- (a) Explain the effect of the rounding direction on the result. (2 points)

**Answer:** The result is always exact. The rounding direction has no effect.

- (b) Prove that the only exception which may be signaled by  $\text{maxnum}(x,y)$  is the invalid exception. (2 points)

**Answer:** Since the result is always exact and is equal to one of the two operands or a NaN then there is no overflow, no underflow, and no inexact results. These exceptions are not signaled. The operation does not involve division and thus the division by zero is not signaled. The only exception which may be signaled is invalid in case of an operand sNaN.

- (c) Explain by an example the comment “this means results might differ among implementations”. (2 points)

**Answer:** This comment is for cases when we do not have  $x < y$ , we do not have  $y < x$ , and we do not have one operand a number while the other is a quiet NaN. The possible cases are  $x = y$  or two unordered operands. The unordered case for example may have  $\text{maxnum}(qNaN1,qNaN2)=qNaN1$  on one implementation while it is  $\text{maxnum}(qNaN1,qNaN2)=qNaN2$  on another implementation. The case of  $x = y$  may arise in  $\text{maxnum}(+0,-0)=+0$  on one implementation and  $-0$  on another.

- (d) For  $\text{maxnum}(\text{maxum}(1,2), \text{maxnum}(3,sNaN))$ , what are the result and exception signaled? (2 points)

**Answer:** The final result is  $\text{maxnum}(2,qNaN)=2$  and the invalid exception is signaled because of the sNaN.

- (e) Is  $\text{maxnum}$  as defined by the standard commutative? Is it associative? (2 points)

**Answer:** As we just proved above, the result of  $\text{maxnum}(qNaN1,qNaN2)$  may differ from  $\text{maxnum}(qNaN2,qNaN1)$  hence it is not commutative. Also the result of  $\text{maxnum}(\text{maxum}(1,2), \text{maxnum}(3,sNaN))$  differs from  $\text{maxnum}(\text{maxnum}(\text{maxum}(1,2),3),sNaN)$  hence it is not associative.

2. Consider a redundant number system with radix  $\beta$  and digits  $d_i$  in the digit set  $\mathcal{D} = \{-\alpha, -\alpha + 1, \dots, \alpha - 1, \alpha\}$ . Addition for the inputs  $x, y$ , and result  $s$  where  $x_i, y_i, s_i \in \mathcal{D}$  with a guaranteed limited carry propagation is possible if we implement the following rules given an appropriate choice of the threshold  $\tau$ :

- At each position  $i$ , form the primary sum  $p_i = x_i + y_i$  of the two operands  $x$  and  $y$ .
- If  $p_i > \tau$  generate a carry  $c_{i+1} = 1$ . If  $p_i < -\tau$  generate a carry  $c_{i+1} = -1$ . Otherwise,  $c_{i+1} = 0$ .
- The intermediate sum at position  $i$  is  $w_i = p_i - \beta c_{i+1}$ .
- The final sum at position  $i$  is  $s_i = w_i + c_i$ .

(a) What is the condition on  $\alpha$  to have a unique representation for zero? (2 points)

**Answer:** We must have  $\alpha < \beta$ , otherwise we might get  $1\bar{\beta} = 00$ .

(b) In order to have no carry propagation (i.e. to absorb  $c_{i+1}$  into  $s_{i+1}$ ) and given that  $c_{i+1} \in \{-1, 0, 1\}$  prove that the possible range for  $\tau$  is  $1 \leq \beta - \alpha \leq \tau \leq \alpha - 1$ . (4 points)

**Answer:** In order to always be able to absorb any incoming carry we must have  $|w_i| < \alpha$  so that  $|w_i + c_i| \leq \alpha$ . If  $\tau \geq \alpha$  and  $p_i = \alpha$  we get  $w_i = \alpha$  which contradicts the above condition, hence  $\tau \leq \alpha - 1$ . Similarly, if  $\tau \leq \beta - \alpha - 1$  while  $p_i = \beta - \alpha$  then  $w_i = \beta - \alpha - 1 \times \beta = -\alpha$  which also contradicts the condition.