# Default Rules for Rounding
## ... Precision Floating Point Arithmetic
### Ignoring Over/underflow.

... these rules would have to be amended only slightly to ... over/underflow, which is a nearly independent and much ... complicated topic. For simplicity here we consider the ... "representable numbers" to be an infinite discrete subset of the ... of real numbers.

... representable numbers must include 0, 1 and, if $x$ then $-x$ too.

... representable number must be represented uniquely by a ... symbol string that represents nothing else.

... arithmetic operation* which, when executed without roundoff ... would produce a representable number, must actually be ... ed without error.

... discard information unnecessarily.

... arithmetic operation which would be ... off error must result in a representable number ... would have been produced in the absence of roundoff ...

... preceding rule is ambiguous when two representable ... ers are nearest the unrounded result. This ambiguity ... is resolved in a systematic way which preserves sign ... etry ( e.g. $x - y = -(y - x)$ ) and is "unbiased" in ... sense that "drift" cannot occur; e.g. the sequence ... $_{1, 2, ...}$ defined for arbitrary $x_0$ and $y$ by $x_{1} := (x_0 + y) - y$
$$= x_2 = x_3 = \cdots$$

... metic operations include $+, -, \times, /, |\cdot|,$ and conversion; ... be extended to include $\sqrt{\phantom{x}}$ and other FORTRAN functions ... above were slightly relaxed.

W. Kahan
Univ. of Calif. @ Berkeley
Sept. 26, 1987