

# A Formal Model of IEEE Floating Point Arithmetic

Lei Yu

December 14, 2021

## Abstract

This development provides a formal model of IEEE-754 floating-point arithmetic. This formalization, including formal specification of the standard and proofs of important properties of floating-point arithmetic, forms the foundation for verifying programs with floating-point computation. There is also a code generation setup for floats so that we can execute programs using this formalization in functional programming languages. The definitions of the IEEE standard in Isabelle is ported from HOL Light [1].

## Contents

<b>1</b>	<b>Specification of the IEEE standard</b>	<b>2</b>
1.1	Derived parameters for floating point formats . . . . .	2
1.2	Predicates for the four IEEE formats . . . . .	2
1.3	Extractors for fields . . . . .	3
1.4	Partition of numbers into disjoint classes . . . . .	3
1.5	Special values . . . . .	4
1.6	Negation operation on floating point values . . . . .	4
1.7	Real number valuations . . . . .	4
1.8	Rounding . . . . .	5
1.9	Definitions of the arithmetic operations . . . . .	6
1.10	Comparison operations . . . . .	8
<b>2</b>	<b>Specify float to be double precision and round to even</b>	<b>9</b>
<b>3</b>	<b>Proofs of Properties about Floating Point Arithmetic</b>	<b>11</b>
3.1	Theorems derived from definitions . . . . .	11
3.2	Properties about ordering and bounding . . . . .	14
3.3	Algebraic properties about basic arithmetic . . . . .	19
3.4	Properties about Rounding Errors . . . . .	20

<b>4</b>	<b>Concrete encodings</b>	<b>22</b>
<b>5</b>	<b>Code Generation Setup for Floats</b>	<b>27</b>
5.1	Conversion from int . . . . .	29
5.2	Conversion to and from software floats, extracting information	30

## 1 Specification of the IEEE standard

```

theory IEEE
  imports
    HOL-Library.Float
    Word-Lib.Word-Lemmas
  begin

  typedef (overloaded) ('e::len, 'f::len) float = UNIV::(1 word × 'e word × 'f
  word) set
  <proof>

  setup-lifting type-definition-float

  syntax -float :: type ⇒ type ⇒ type (('(-, -') float)

  parse ('a, 'b) float as ('a::len, 'b::len) float.

  <ML>

```

### 1.1 Derived parameters for floating point formats

```

definition wordlength :: ('e, 'f) float itself ⇒ nat
  where wordlength x = LENGTH('e) + LENGTH('f) + 1

```

```

definition bias :: ('e, 'f) float itself ⇒ nat
  where bias x = 2LENGTH('e) - 1 - 1

```

```

definition emax :: ('e, 'f) float itself ⇒ nat
  where emax x = unat (- 1::'e word)

```

```

abbreviation fracwidth::('e, 'f) float itself ⇒ nat where
  fracwidth - ≡ LENGTH('f)

```

### 1.2 Predicates for the four IEEE formats

```

definition is-single :: ('e, 'f) float itself ⇒ bool
  where is-single x ⟷ LENGTH('e) = 8 ∧ wordlength x = 32

```

```

definition is-double :: ('e, 'f) float itself ⇒ bool
  where is-double x ⟷ LENGTH('e) = 11 ∧ wordlength x = 64

```

```

definition is-single-extended :: ('e, 'f) float itself ⇒ bool

```

**where** *is-single-extended*  $x \longleftrightarrow \text{LENGTH}('e) \geq 11 \wedge \text{wordlength } x \geq 43$

**definition** *is-double-extended*  $:: ('e, 'f) \text{ float} \Rightarrow \text{bool}$

**where** *is-double-extended*  $x \longleftrightarrow \text{LENGTH}('e) \geq 15 \wedge \text{wordlength } x \geq 79$

### 1.3 Extractors for fields

**lift-definition** *sign*  $:: ('e, 'f) \text{ float} \Rightarrow \text{nat is}$

$\lambda(s::1 \text{ word}, ::'e \text{ word}, ::'f \text{ word}). \text{unat } s \langle \text{proof} \rangle$

**lift-definition** *exponent*  $:: ('e, 'f) \text{ float} \Rightarrow \text{nat is}$

$\lambda(-, e::'e \text{ word}, -). \text{unat } e \langle \text{proof} \rangle$

**lift-definition** *fraction*  $:: ('e, 'f) \text{ float} \Rightarrow \text{nat is}$

$\lambda(-, -, f::'f \text{ word}). \text{unat } f \langle \text{proof} \rangle$

**abbreviation** *real-of-word*  $x \equiv \text{real } (\text{unat } x)$

**lift-definition** *valof*  $:: ('e, 'f) \text{ float} \Rightarrow \text{real}$

**is**  $\lambda(s, e, f).$

*let*  $x = (\text{TYPE} (('e, 'f) \text{ float}))$  *in*

*(if*  $e = 0$

*then*  $(-1::\text{real})^{\wedge}(\text{unat } s) * (2 / (2^{\wedge} \text{bias } x)) * (\text{real-of-word } f / 2^{\wedge}(\text{LENGTH}('f)))$

*else*  $(-1::\text{real})^{\wedge}(\text{unat } s) * ((2^{\wedge}(\text{unat } e)) / (2^{\wedge} \text{bias } x)) * (1 + \text{real-of-word } f / 2^{\wedge}(\text{LENGTH}('f)))$

*proof*)

### 1.4 Partition of numbers into disjoint classes

**definition** *is-nan*  $:: ('e, 'f) \text{ float} \Rightarrow \text{bool}$

**where** *is-nan*  $a \longleftrightarrow \text{exponent } a = \text{emax } \text{TYPE} (('e, 'f) \text{ float}) \wedge \text{fraction } a \neq 0$

**definition** *is-infinity*  $:: ('e, 'f) \text{ float} \Rightarrow \text{bool}$

**where** *is-infinity*  $a \longleftrightarrow \text{exponent } a = \text{emax } \text{TYPE} (('e, 'f) \text{ float}) \wedge \text{fraction } a = 0$

**definition** *is-normal*  $:: ('e, 'f) \text{ float} \Rightarrow \text{bool}$

**where** *is-normal*  $a \longleftrightarrow 0 < \text{exponent } a \wedge \text{exponent } a < \text{emax } \text{TYPE} (('e, 'f) \text{ float})$

**definition** *is-denormal*  $:: ('e, 'f) \text{ float} \Rightarrow \text{bool}$

**where** *is-denormal*  $a \longleftrightarrow \text{exponent } a = 0 \wedge \text{fraction } a \neq 0$

**definition** *is-zero*  $:: ('e, 'f) \text{ float} \Rightarrow \text{bool}$

**where** *is-zero*  $a \longleftrightarrow \text{exponent } a = 0 \wedge \text{fraction } a = 0$

**definition** *is-finite*  $:: ('e, 'f) \text{ float} \Rightarrow \text{bool}$

**where** *is-finite*  $a \longleftrightarrow (\text{is-normal } a \vee \text{is-denormal } a \vee \text{is-zero } a)$

## 1.5 Special values

**lift-definition** *plus-infinity* :: ('e, 'f) float ( $\infty$ ) **is** (0, - 1, 0) *<proof>*

**lift-definition** *topfloat* :: ('e, 'f) float **is** (0, - 2,  $2^{\text{LENGTH('f)} - 1}$ ) *<proof>*

**instantiation** *float::(len, len) zero begin*

**lift-definition** *zero-float* :: ('e, 'f) float **is** (0, 0, 0) *<proof>*

**instance** *<proof>*

**end**

## 1.6 Negation operation on floating point values

**instantiation** *float::(len, len) uminus begin*

**lift-definition** *uminus-float* :: ('e, 'f) float  $\Rightarrow$  ('e, 'f) float **is**  $\lambda(s, e, f). (1 - s, e, f)$  *<proof>*

**instance** *<proof>*

**end**

**abbreviation** (*input*) *minus-zero*  $\equiv - (0::('e, 'f)float)$

**abbreviation** (*input*) *minus-infinity*  $\equiv - \infty$

**abbreviation** (*input*) *bottomfloat*  $\equiv - \text{topfloat}$

## 1.7 Real number valuations

The largest value that can be represented in floating point format.

**definition** *largest* :: ('e, 'f) float *itself*  $\Rightarrow$  real

**where** *largest*  $x = (2^{\text{emax } x - 1} / 2^{\text{bias } x}) * (2 - 1 / (2^{\text{fracwidth } x}))$

Threshold, used for checking overflow.

**definition** *threshold* :: ('e, 'f) float *itself*  $\Rightarrow$  real

**where** *threshold*  $x = (2^{\text{emax } x - 1} / 2^{\text{bias } x}) * (2 - 1 / (2^{\text{Suc}(\text{fracwidth } x)}))$

Unit of least precision.

**lift-definition** *one-lp::('e, 'f) float*  $\Rightarrow$  ('e, 'f) float **is**  $\lambda(s, e, f). (0, e::'e \text{ word}, 1)$  *<proof>*

**lift-definition** *zero-lp::('e, 'f) float*  $\Rightarrow$  ('e, 'f) float **is**  $\lambda(s, e, f). (0, e::'e \text{ word}, 0)$  *<proof>*

**definition** *ulp* :: ('e, 'f) float  $\Rightarrow$  real **where** *ulp*  $a = \text{valof } (\text{one-lp } a) - \text{valof } (\text{zero-lp } a)$

Enumerated type for rounding modes.

**datatype** *roundmode* = *To-nearest* | *float-To-zero* | *To-pinfinity* | *To-ninfinity*

Characterization of best approximation from a set of abstract values.

**definition** *is-closest*  $v\ s\ x\ a \longleftrightarrow a \in s \wedge (\forall b. b \in s \longrightarrow |v\ a - x| \leq |v\ b - x|)$

Best approximation with a deciding preference for multiple possibilities.

**definition** *closest*  $v\ p\ s\ x =$

(*SOME*  $a. \text{is-closest } v\ s\ x\ a \wedge ((\exists b. \text{is-closest } v\ s\ x\ b \wedge p\ b) \longrightarrow p\ a)$ )

## 1.8 Rounding

**fun** *round* :: *roundmode*  $\Rightarrow$  *real*  $\Rightarrow$  ('e, 'f) *float*

**where**

*round To-nearest*  $y =$

(if  $y < -\text{threshold } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *minus-infinity*  
 else if  $y \geq \text{threshold } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *plus-infinity*  
 else *closest (valof)* ( $\lambda a. \text{even (fraction } a)$ ) {*a. is-finite*  $a$ }  $y$ )

| *round float-To-zero*  $y =$

(if  $y < -\text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *bottomfloat*  
 else if  $y > \text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *topfloat*  
 else *closest (valof)* ( $\lambda a. \text{True}$ ) {*a. is-finite*  $a \wedge |\text{valof } a| \leq |y|$ }  $y$ )

| *round To-pinfinity*  $y =$

(if  $y < -\text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *bottomfloat*  
 else if  $y > \text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *plus-infinity*  
 else *closest (valof)* ( $\lambda a. \text{True}$ ) {*a. is-finite*  $a \wedge \text{valof } a \geq y$ }  $y$ )

| *round To-ninfinity*  $y =$

(if  $y < -\text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *minus-infinity*  
 else if  $y > \text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *topfloat*  
 else *closest (valof)* ( $\lambda a. \text{True}$ ) {*a. is-finite*  $a \wedge \text{valof } a \leq y$ }  $y$ )

Rounding to integer values in floating point format.

**definition** *is-integral* :: ('e, 'f) *float*  $\Rightarrow$  *bool*

**where** *is-integral*  $a \longleftrightarrow \text{is-finite } a \wedge (\exists n::\text{nat}. |\text{valof } a| = \text{real } n)$

**fun** *intround* :: *roundmode*  $\Rightarrow$  *real*  $\Rightarrow$  ('e, 'f) *float*

**where**

*intround To-nearest*  $y =$

(if  $y \leq -\text{threshold } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *minus-infinity*  
 else if  $y \geq \text{threshold } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *plus-infinity*  
 else *closest (valof)* ( $\lambda a. (\exists n::\text{nat}. \text{even } n \wedge |\text{valof } a| = \text{real } n)$ ) {*a. is-integral*  
 $a$ }  $y$ )

| *intround float-To-zero*  $y =$

(if  $y < -\text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *bottomfloat*  
 else if  $y > \text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *topfloat*  
 else *closest (valof)* ( $\lambda x. \text{True}$ ) {*a. is-integral*  $a \wedge |\text{valof } a| \leq |y|$ }  $y$ )

| *intround To-pinfinity*  $y =$

(if  $y < -\text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *bottomfloat*  
 else if  $y > \text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *plus-infinity*  
 else *closest (valof)* ( $\lambda x. \text{True}$ ) {*a. is-integral*  $a \wedge \text{valof } a \geq y$ }  $y$ )

| *intround To-ninfinity*  $y =$

(if  $y < -\text{largest } \text{TYPE}((\text{'e}, \text{'f}) \text{float})$  then *minus-infinity*

else if  $y > \text{largest TYPE}('e, 'f) \text{ float}$  then *topfloat*  
 else *closest (valof) ( $\lambda x. \text{True}$ ) \{a. is-integral a  $\wedge$  valof a  $\geq$  y\} y*)

Round, choosing between -0.0 or +0.0

**definition** *float-round::roundmode  $\Rightarrow$  bool  $\Rightarrow$  real  $\Rightarrow$  ('e, 'f) float*  
**where** *float-round mode toneg r =*  
 (*let x = round mode r in*  
*if is-zero x*  
*then if toneg*  
*then minus-zero*  
*else 0*  
*else x*)

Non-standard of NaN.

**definition** *some-nan :: ('e, 'f) float*  
**where** *some-nan = (SOME a. is-nan a)*

Coercion for signs of zero results.

**definition** *zerosign :: nat  $\Rightarrow$  ('e, 'f) float  $\Rightarrow$  ('e, 'f) float*  
**where** *zerosign s a =*  
 (*if is-zero a then (if s = 0 then 0 else - 0) else a*)

Remainder operation.

**definition** *rem :: real  $\Rightarrow$  real  $\Rightarrow$  real*  
**where** *rem x y =*  
 (*let n = closest id ( $\lambda x. \exists n::nat. \text{even } n \wedge |x| = \text{real } n$ ) \{x.  $\exists n::nat. |x| = \text{real } n$ \}*  
*in x / y*)  
*in x - n \* y*)

**definition** *frem :: roundmode  $\Rightarrow$  ('e, 'f) float  $\Rightarrow$  ('e, 'f) float  $\Rightarrow$  ('e, 'f) float*  
**where** *frem m a b =*  
 (*if is-nan a  $\vee$  is-nan b  $\vee$  is-infinity a  $\vee$  is-zero b then some-nan*  
*else zerosign (sign a) (round m (rem (valof a) (valof b)))*)

## 1.9 Definitions of the arithmetic operations

**definition** *fintrnd :: roundmode  $\Rightarrow$  ('e, 'f) float  $\Rightarrow$  ('e, 'f) float*  
**where** *fintrnd m a =*  
 (*if is-nan a then (some-nan)*  
*else if is-infinity a then a*  
*else zerosign (sign a) (intround m (valof a))*)

**definition** *fadd :: roundmode  $\Rightarrow$  ('e, 'f) float  $\Rightarrow$  ('e, 'f) float  $\Rightarrow$  ('e, 'f) float*  
**where** *fadd m a b =*  
 (*if is-nan a  $\vee$  is-nan b  $\vee$  (is-infinity a  $\wedge$  is-infinity b  $\wedge$  sign a  $\neq$  sign b)*  
*then some-nan*  
*else if (is-infinity a) then a*  
*else if (is-infinity b) then b*  
*else*

*zerosign*  
 (if *is-zero* *a*  $\wedge$  *is-zero* *b*  $\wedge$  *sign* *a* = *sign* *b* then *sign* *a*  
 else if *m* = *To-ninfinity* then 1 else 0)  
 (round *m* (valof *a* + valof *b*)))

**definition** *fsub* :: *roundmode*  $\Rightarrow$  ('e , 'f) float  $\Rightarrow$  ('e , 'f) float  $\Rightarrow$  ('e , 'f) float  
**where** *fsub* *m* *a* *b* =  
 (if *is-nan* *a*  $\vee$  *is-nan* *b*  $\vee$  (*is-infinity* *a*  $\wedge$  *is-infinity* *b*  $\wedge$  *sign* *a* = *sign* *b*)  
 then *some-nan*  
 else if *is-infinity* *a* then *a*  
 else if *is-infinity* *b* then - *b*  
 else  
*zerosign*  
 (if *is-zero* *a*  $\wedge$  *is-zero* *b*  $\wedge$  *sign* *a*  $\neq$  *sign* *b* then *sign* *a*  
 else if *m* = *To-ninfinity* then 1 else 0)  
 (round *m* (valof *a* - valof *b*)))

**definition** *fmul* :: *roundmode*  $\Rightarrow$  ('e , 'f) float  $\Rightarrow$  ('e , 'f) float  $\Rightarrow$  ('e , 'f) float  
**where** *fmul* *m* *a* *b* =  
 (if *is-nan* *a*  $\vee$  *is-nan* *b*  $\vee$  (*is-zero* *a*  $\wedge$  *is-infinity* *b*)  $\vee$  (*is-infinity* *a*  $\wedge$  *is-zero* *b*)  
 then *some-nan*  
 else if *is-infinity* *a*  $\vee$  *is-infinity* *b*  
 then (if *sign* *a* = *sign* *b* then *plus-infinity* else *minus-infinity*)  
 else *zerosign* (if *sign* *a* = *sign* *b* then 0 else 1) (round *m* (valof *a* \* valof *b*)))

**definition** *fdiv* :: *roundmode*  $\Rightarrow$  ('e , 'f) float  $\Rightarrow$  ('e , 'f) float  $\Rightarrow$  ('e , 'f) float  
**where** *fdiv* *m* *a* *b* =  
 (if *is-nan* *a*  $\vee$  *is-nan* *b*  $\vee$  (*is-zero* *a*  $\wedge$  *is-zero* *b*)  $\vee$  (*is-infinity* *a*  $\wedge$  *is-infinity* *b*)  
 then *some-nan*  
 else if *is-infinity* *a*  $\vee$  *is-zero* *b*  
 then (if *sign* *a* = *sign* *b* then *plus-infinity* else *minus-infinity*)  
 else if *is-infinity* *b*  
 then (if *sign* *a* = *sign* *b* then 0 else - 0)  
 else *zerosign* (if *sign* *a* = *sign* *b* then 0 else 1) (round *m* (valof *a* / valof *b*)))

**definition** *fsqrt* :: *roundmode*  $\Rightarrow$  ('e , 'f) float  $\Rightarrow$  ('e , 'f) float  
**where** *fsqrt* *m* *a* =  
 (if *is-nan* *a* then *some-nan*  
 else if *is-zero* *a*  $\vee$  *is-infinity* *a*  $\wedge$  *sign* *a* = 0 then *a*  
 else if *sign* *a* = 1 then *some-nan*  
 else *zerosign* (*sign* *a*) (round *m* (sqrt (valof *a*))))

**definition** *fmul-add* :: *roundmode*  $\Rightarrow$  ('t , 'w) float  $\Rightarrow$  ('t , 'w) float  $\Rightarrow$  ('t , 'w) float  
 $\Rightarrow$  ('t , 'w) float  
**where** *fmul-add* *mode* *x* *y* *z* =  
 (let *signP* = if *sign* *x* = *sign* *y* then 0 else 1 in  
 let *infP* = *is-infinity* *x*  $\vee$  *is-infinity* *y*  
 in  
 if *is-nan* *x*  $\vee$  *is-nan* *y*  $\vee$  *is-nan* *z* then *some-nan*

```

else if is-infinity x ∧ is-zero y ∨
      is-zero x ∧ is-infinity y ∨
      is-infinity z ∧ infP ∧ signP ≠ sign z then
  some-nan
else if is-infinity z ∧ (sign z = 0) ∨ infP ∧ (signP = 0)
  then plus-infinity
else if is-infinity z ∧ (sign z = 1) ∨ infP ∧ (signP = 1)
  then minus-infinity
else
  let r1 = valof x * valof y;
      r2 = valof z
  in
    float-round mode
      (if (r1 = 0) ∧ (r2 = 0) ∧ (signP = sign z) then
        signP = 1
      else mode = To-ninfinity) (r1 + r2))

```

## 1.10 Comparison operations

**datatype** *ccode* = *Gt* | *Lt* | *Eq* | *Und*

**definition** *fcompare* :: ('e, 'f) float ⇒ ('e, 'f) float ⇒ *ccode*

**where** *fcompare* *a b* =

```

(if is-nan a ∨ is-nan b then Und
 else if is-infinity a ∧ sign a = 1
  then (if is-infinity b ∧ sign b = 1 then Eq else Lt)
 else if is-infinity a ∧ sign a = 0
  then (if is-infinity b ∧ sign b = 0 then Eq else Gt)
 else if is-infinity b ∧ sign b = 1 then Gt
 else if is-infinity b ∧ sign b = 0 then Lt
 else if valof a < valof b then Lt
 else if valof a = valof b then Eq
 else Gt)

```

**definition** *flt* :: ('e, 'f) float ⇒ ('e, 'f) float ⇒ bool

**where** *flt* *a b* ⇔ *fcompare* *a b* = *Lt*

**definition** *fle* :: ('e, 'f) float ⇒ ('e, 'f) float ⇒ bool

**where** *fle* *a b* ⇔ *fcompare* *a b* = *Lt* ∨ *fcompare* *a b* = *Eq*

**definition** *fgt* :: ('e, 'f) float ⇒ ('e, 'f) float ⇒ bool

**where** *fgt* *a b* ⇔ *fcompare* *a b* = *Gt*

**definition** *fge* :: ('e, 'f) float ⇒ ('e, 'f) float ⇒ bool

**where** *fge* *a b* ⇔ *fcompare* *a b* = *Gt* ∨ *fcompare* *a b* = *Eq*

**definition** *feq* :: ('e, 'f) float ⇒ ('e, 'f) float ⇒ bool

**where** *feq* *a b* ⇔ *fcompare* *a b* = *Eq*

## 2 Specify float to be double precision and round to even

**instantiation** *float* :: (*len*, *len*) *plus*  
**begin**

**definition** *plus-float* :: ('*a*', '*b*') *float* ⇒ ('*a*', '*b*') *float* ⇒ ('*a*', '*b*') *float*  
  **where**  $a + b = fadd$  *To-nearest* *a b*

**instance** ⟨*proof*⟩

**end**

**instantiation** *float* :: (*len*, *len*) *minus*  
**begin**

**definition** *minus-float* :: ('*a*', '*b*') *float* ⇒ ('*a*', '*b*') *float* ⇒ ('*a*', '*b*') *float*  
  **where**  $a - b = fsub$  *To-nearest* *a b*

**instance** ⟨*proof*⟩

**end**

**instantiation** *float* :: (*len*, *len*) *times*  
**begin**

**definition** *times-float* :: ('*a*', '*b*') *float* ⇒ ('*a*', '*b*') *float* ⇒ ('*a*', '*b*') *float*  
  **where**  $a * b = fmul$  *To-nearest* *a b*

**instance** ⟨*proof*⟩

**end**

**instantiation** *float* :: (*len*, *len*) *one*  
**begin**

**lift-definition** *one-float* :: ('*a*', '*b*') *float* **is**  $(0, 2^{\wedge}(LENGTH('a) - 1) - 1, 0)$   
  ⟨*proof*⟩

**instance** ⟨*proof*⟩

**end**

**instantiation** *float* :: (*len*, *len*) *inverse*  
**begin**

**definition** *divide-float* :: ('*a*', '*b*') *float* ⇒ ('*a*', '*b*') *float* ⇒ ('*a*', '*b*') *float*  
  **where**  $a \div b = fdiv$  *To-nearest* *a b*

**definition** *inverse-float* :: ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  
**where** *inverse-float* a = fdiv To-nearest 1 a

**instance** <proof>

**end**

**definition** *float-rem* :: ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  
**where** *float-rem* a b = frem To-nearest a b

**definition** *float-sqrt* :: ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  
**where** *float-sqrt* a = fsqrt To-nearest a

**definition** *ROUNDFLOAT* :: ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  
**where** *ROUNDFLOAT* a = fintrnd To-nearest a

**instantiation** *float* :: (len, len) ord  
**begin**

**definition** *less-float* :: ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  $\Rightarrow$  bool  
**where**  $a < b \longleftrightarrow$  flt a b

**definition** *less-eq-float* :: ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  $\Rightarrow$  bool  
**where**  $a \leq b \longleftrightarrow$  fle a b

**instance** <proof>

**end**

**definition** *float-eq* :: ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  $\Rightarrow$  bool (**infixl**  $\doteq$  70)  
**where** *float-eq* a b = feq a b

**instantiation** *float* :: (len, len) abs  
**begin**

**definition** *abs-float* :: ('a, 'b) float  $\Rightarrow$  ('a, 'b) float  
**where** *abs-float* a = (if sign a = 0 then a else - a)

**instance** <proof>

**end**

The  $1 + \varepsilon$  property.

**definition** *normalizes* :: - itself  $\Rightarrow$  real  $\Rightarrow$  bool  
**where** *normalizes* float-format x =  
 $(1 / (2 :: \text{real}) \sim (\text{bias float-format} - 1) \leq |x| \wedge |x| < \text{threshold float-format})$

**end**

### 3 Proofs of Properties about Floating Point Arithmetic

```
theory IEEE-Properties
imports IEEE
begin
```

#### 3.1 Theorems derived from definitions

**lemma** *valof-eq*:

```
valof x =
  (if exponent x = 0
   then (- 1) ^ sign x * (2 / 2 ^ bias TYPE(('a, 'b) float)) *
    (real (fraction x) / 2 ^ LENGTH('b))
   else (- 1) ^ sign x * (2 ^ exponent x / 2 ^ bias TYPE(('a, 'b) float)) *
    (1 + real (fraction x) / 2 ^ LENGTH('b)))
for x::('a, 'b) float
⟨proof⟩
```

**lemma** *exponent-le* [*simp*]:

```
⟨exponent a ≤ mask LENGTH('a)⟩ for a :: ⟨('a, -) float⟩
⟨proof⟩
```

**lemma** *exponent-not-less* [*simp*]:

```
⟨¬ mask LENGTH('a) < IEEE.exponent a⟩ for a :: ⟨('a, -) float⟩
⟨proof⟩
```

**lemma** *infinity-simps*:

```
sign (plus-infinity::('e, 'f)float) = 0
sign (minus-infinity::('e, 'f)float) = 1
exponent (plus-infinity::('e, 'f)float) = emax TYPE(('e, 'f)float)
exponent (minus-infinity::('e, 'f)float) = emax TYPE(('e, 'f)float)
fraction (plus-infinity::('e, 'f)float) = 0
fraction (minus-infinity::('e, 'f)float) = 0
⟨proof⟩
```

**lemma** *zero-simps*:

```
sign (0::('e, 'f)float) = 0
sign (- 0::('e, 'f)float) = 1
exponent (0::('e, 'f)float) = 0
exponent (- 0::('e, 'f)float) = 0
fraction (0::('e, 'f)float) = 0
fraction (- 0::('e, 'f)float) = 0
⟨proof⟩
```

**lemma** *emax-eq*:  $emax\ x = 2 \wedge LENGTH('e) - 1$

```
for x::('e, 'f)float itself
⟨proof⟩
```

**lemma** *topfloat-simps*:

$$\text{sign } (\text{topfloat}::('e, 'f)\text{float}) = 0$$

$$\text{exponent } (\text{topfloat}::('e, 'f)\text{float}) = \text{emax } \text{TYPE} (('e, 'f)\text{float}) - 1$$

$$\text{fraction } (\text{topfloat}::('e, 'f)\text{float}) = 2 \wedge (\text{fracwidth } \text{TYPE} (('e, 'f)\text{float})) - 1$$

**and** *bottomfloat-simps*:

$$\text{sign } (\text{bottomfloat}::('e, 'f)\text{float}) = 1$$

$$\text{exponent } (\text{bottomfloat}::('e, 'f)\text{float}) = \text{emax } \text{TYPE} (('e, 'f)\text{float}) - 1$$

$$\text{fraction } (\text{bottomfloat}::('e, 'f)\text{float}) = 2 \wedge (\text{fracwidth } \text{TYPE} (('e, 'f)\text{float})) - 1$$

*<proof>*

**lemmas** *float-defs* =

*is-finite-def is-infinity-def is-zero-def is-nan-def*

*is-normal-def is-denormal-def valof-eq*

*less-eq-float-def less-float-def*

*flt-def fgt-def fle-def fge-def feq-def*

*fcompare-def*

*infinity-simps*

*zero-simps*

*topfloat-simps*

*bottomfloat-simps*

*float-eq-def*

**lemma** *float-cases*:  $\text{is-nan } a \vee \text{is-infinity } a \vee \text{is-normal } a \vee \text{is-denormal } a \vee \text{is-zero } a$

*<proof>*

**lemma** *float-cases-finite*:  $\text{is-nan } a \vee \text{is-infinity } a \vee \text{is-finite } a$

*<proof>*

**lemma** *float-zero1[simp]*:  $\text{is-zero } 0$

*<proof>*

**lemma** *float-zero2[simp]*:  $\text{is-zero } (- x) \longleftrightarrow \text{is-zero } x$

*<proof>*

**lemma** *emax-pos[simp]*:  $0 < \text{emax } x \text{ emax } x \neq 0$

*<proof>*

The types of floating-point numbers are mutually distinct.

**lemma** *float-distinct*:

$\neg (\text{is-nan } a \wedge \text{is-infinity } a)$

$\neg (\text{is-nan } a \wedge \text{is-normal } a)$

$\neg (\text{is-nan } a \wedge \text{is-denormal } a)$

$\neg (\text{is-nan } a \wedge \text{is-zero } a)$

$\neg (\text{is-infinity } a \wedge \text{is-normal } a)$

$\neg (\text{is-infinity } a \wedge \text{is-denormal } a)$

$\neg (\text{is-infinity } a \wedge \text{is-zero } a)$

$\neg (\text{is-normal } a \wedge \text{is-denormal } a)$

$\neg (\text{is-denormal } a \wedge \text{is-zero } a)$

*<proof>*

**lemma** *denormal-imp-not-zero*: *is-denormal f*  $\implies$   $\neg$ *is-zero f*  
*<proof>*

**lemma** *normal-imp-not-zero*: *is-normal f*  $\implies$   $\neg$ *is-zero f*  
*<proof>*

**lemma** *normal-imp-not-denormal*: *is-normal f*  $\implies$   $\neg$ *is-denormal f*  
*<proof>*

**lemma** *denormal-zero[simp]*:  $\neg$ *is-denormal 0*  $\neg$ *is-denormal minus-zero*  
*<proof>*

**lemma** *normal-zero[simp]*:  $\neg$ *is-normal 0*  $\neg$ *is-normal minus-zero*  
*<proof>*

**lemma** *float-distinct-finite*:  $\neg$  (*is-nan a*  $\wedge$  *is-finite a*)  $\neg$ (*is-infinity a*  $\wedge$  *is-finite a*)  
*<proof>*

**lemma** *finite-infinity*: *is-finite a*  $\implies$   $\neg$  *is-infinity a*  
*<proof>*

**lemma** *finite-nan*: *is-finite a*  $\implies$   $\neg$  *is-nan a*  
*<proof>*

For every real number, the floating-point numbers closest to it always exist.

**lemma** *is-closest-exists*:  
**fixes** *v* :: ('e, 'f)float  $\Rightarrow$  *real*  
**and** *s* :: ('e, 'f)float set  
**assumes** *finite*: *finite s*  
**and** *non-empty*: *s*  $\neq$  {}  
**shows**  $\exists a.$  *is-closest v s x a*  
*<proof>*

**lemma** *closest-is-everything*:  
**fixes** *v* :: ('e, 'f)float  $\Rightarrow$  *real*  
**and** *s* :: ('e, 'f)float set  
**assumes** *finite*: *finite s*  
**and** *non-empty*: *s*  $\neq$  {}  
**shows** *is-closest v s x* (*closest v p s x*)  $\wedge$   
(( $\exists b.$  *is-closest v s x b*  $\wedge$  *p b*)  $\longrightarrow$  *p* (*closest v p s x*)  
*<proof>*

**lemma** *closest-in-set*:  
**fixes** *v* :: ('e, 'f)float  $\Rightarrow$  *real*  
**assumes** *finite s* **and** *s*  $\neq$  {}  
**shows** *closest v p s x*  $\in$  *s*  
*<proof>*

**lemma** *closest-is-closest-finite*:  
**fixes**  $v :: ('e, 'f)\text{float} \Rightarrow \text{real}$   
**assumes** *finite s and  $s \neq \{\}$*   
**shows** *is-closest v s x (closest v p s x)*  
 $\langle \text{proof} \rangle$

**instance** *float::(len, len) finite*  $\langle \text{proof} \rangle$

**lemma** *is-finite-nonempty*:  $\{a. \text{is-finite } a\} \neq \{\}$   
 $\langle \text{proof} \rangle$

**lemma** *closest-is-closest*:  
**fixes**  $v :: ('e, 'f)\text{float} \Rightarrow \text{real}$   
**assumes**  $s \neq \{\}$   
**shows** *is-closest v s x (closest v p s x)*  
 $\langle \text{proof} \rangle$

### 3.2 Properties about ordering and bounding

Lifting of non-exceptional comparisons.

**lemma** *float-lt* [*simp*]:  
**assumes** *is-finite a is-finite b*  
**shows**  $a < b \iff \text{valof } a < \text{valof } b$   
 $\langle \text{proof} \rangle$

**lemma** *float-eq* [*simp*]:  
**assumes** *is-finite a is-finite b*  
**shows**  $a \doteq b \iff \text{valof } a = \text{valof } b$   
 $\langle \text{proof} \rangle$

**lemma** *float-le* [*simp*]:  
**assumes** *is-finite a is-finite b*  
**shows**  $a \leq b \iff \text{valof } a \leq \text{valof } b$   
 $\langle \text{proof} \rangle$

Reflexivity of equality for non-NaN's.

**lemma** *float-eq-refl* [*simp*]:  $a \doteq a \iff \neg \text{is-nan } a$   
 $\langle \text{proof} \rangle$

Properties about Ordering.

**lemma** *float-lt-trans*:  $\text{is-finite } a \implies \text{is-finite } b \implies \text{is-finite } c \implies a < b \implies b < c \implies a < c$   
 $\langle \text{proof} \rangle$

**lemma** *float-le-less-trans*:  $\text{is-finite } a \implies \text{is-finite } b \implies \text{is-finite } c \implies a \leq b \implies b < c \implies a < c$   
 $\langle \text{proof} \rangle$

**lemma** *float-le-trans*:  $is\_finite\ a \implies is\_finite\ b \implies is\_finite\ c \implies a \leq b \implies b \leq c \implies a \leq c$   
 ⟨proof⟩

**lemma** *float-le-neg*:  $is\_finite\ a \implies is\_finite\ b \implies \neg a < b \iff b \leq a$   
 ⟨proof⟩

Properties about bounding.

**lemma** *float-le-infinity* [*simp*]:  $\neg is\_nan\ a \implies a \leq plus\_infinity$   
 ⟨proof⟩

**lemma** *zero-le-topfloat*[*simp*]:  $0 \leq topfloat - 0 \leq topfloat$   
 ⟨proof⟩

**lemma** *LENGTH-contr*:  
 $Suc\ 0 < LENGTH('e) \implies 2 \wedge LENGTH('e::len) \leq Suc\ (Suc\ 0) \implies False$   
 ⟨proof⟩

**lemma** *valof-topfloat*:  $valof\ (topfloat::('e, 'f)float) = largest\ TYPE(('e, 'f)float)$   
**if**  $LENGTH('e) > 1$   
 ⟨proof⟩

**lemma** *float-frac-le*:  $fraction\ a \leq 2^{LENGTH('f)} - 1$   
**for**  $a::('e, 'f)float$   
 ⟨proof⟩

**lemma** *float-exp-le*:  $is\_finite\ a \implies exponent\ a \leq emax\ TYPE(('e, 'f)float) - 1$   
**for**  $a::('e, 'f)float$   
 ⟨proof⟩

**lemma** *float-sign-le*:  $(-1::real) \wedge (sign\ a) = 1 \vee (-1::real) \wedge (sign\ a) = -1$   
 ⟨proof⟩

**lemma** *exp-less*:  $a \leq b \implies (2::real) \wedge a \leq 2 \wedge b$  **for**  $a\ b::nat$   
 ⟨proof⟩

**lemma** *div-less*:  $a \leq b \wedge c > 0 \implies a/c \leq b/c$  **for**  $a\ b\ c::'a::linordered\_field$   
 ⟨proof⟩

**lemma** *finite-topfloat*:  $is\_finite\ topfloat$   
 ⟨proof⟩

**lemmas** *float-leI* = *float-le*[*THEN iffD2*]

**lemma** *factor-minus*:  $x * a - x = x * (a - 1)$   
**for**  $x\ a::'a::comm\_semiring-1-cancel$   
 ⟨proof⟩

**lemma** *real-le-power-numeral-diff*:  $\text{real } a \leq \text{numeral } b \wedge n - 1 \iff a \leq \text{numeral } b \wedge n - 1$

*<proof>*

**definition** *denormal-exponent*::('e, 'f)float itself  $\Rightarrow$  int **where**

*denormal-exponent*  $x = 1 - (\text{int } (\text{LENGTH } ('f)) + \text{int } (\text{bias TYPE } (('e, 'f)\text{float})))$

**definition** *normal-exponent*::('e, 'f)float  $\Rightarrow$  int **where**

*normal-exponent*  $x = \text{int } (\text{exponent } x) - \text{int } (\text{bias TYPE } (('e, 'f)\text{float})) - \text{int } (\text{LENGTH } ('f))$

**definition** *denormal-mantissa*::('e, 'f)float  $\Rightarrow$  int **where**

*denormal-mantissa*  $x = (-1::\text{int}) \wedge \text{sign } x * \text{int } (\text{fraction } x)$

**definition** *normal-mantissa*::('e, 'f)float  $\Rightarrow$  int **where**

*normal-mantissa*  $x = (-1::\text{int}) \wedge \text{sign } x * (2^{\text{LENGTH } ('f)} + \text{int } (\text{fraction } x))$

**lemma** *unat-one-word-le*:  $\text{unat } a \leq \text{Suc } 0$  **for**  $a::1$  word

*<proof>*

**lemma** *one-word-le*:  $a \leq 1$  **for**  $a::1$  word

*<proof>*

**lemma** *sign-cases*[*case-names pos neg*]:

**obtains**  $\text{sign } x = 0 \mid \text{sign } x = 1$

*<proof>*

**lemma** *is-infinity-cases*:

**assumes** *is-infinity*  $x$

**obtains**  $x = \text{plus-infinity} \mid x = \text{minus-infinity}$

*<proof>*

**lemma** *is-zero-cases*:

**assumes** *is-zero*  $x$

**obtains**  $x = 0 \mid x = -0$

*<proof>*

**lemma** *minus-minus-float*[*simp*]:  $-(-f) = f$  **for**  $f::('e, 'f)\text{float}$

*<proof>*

**lemma** *sign-minus-float*:  $\text{sign } (-f) = (1 - \text{sign } f)$  **for**  $f::('e, 'f)\text{float}$

*<proof>*

**lemma** *exponent-uminus*[*simp*]:  $\text{exponent } (-f) = \text{exponent } f$  *<proof>*

**lemma** *fraction-uminus*[*simp*]:  $\text{fraction } (-f) = \text{fraction } f$  *<proof>*

**lemma** *is-normal-minus-float*[*simp*]:  $\text{is-normal } (-f) = \text{is-normal } f$  **for**  $f::('e, 'f)\text{float}$

*<proof>*

**lemma** *is-denormal-minus-float*[simp]: *is-denormal*  $(-f) = \text{is-denormal } f$  **for**  $f::('e, 'f)\text{float}$   
 ⟨proof⟩

**lemma** *bitlen-normal-mantissa*:  
*bitlen*  $(\text{abs } (\text{normal-mantissa } x)) = \text{Suc } \text{LENGTH}(f)$  **for**  $x::('e, 'f)\text{float}$   
 ⟨proof⟩

**lemma** *less-int-natI*:  $x < y$  **if**  $0 \leq x$  *nat*  $x < \text{nat } y$   
 ⟨proof⟩

**lemma** *normal-exponent-bounds-int*:  
 $2 - 2^{\wedge}(\text{LENGTH}(e) - 1) - \text{int } \text{LENGTH}(f) \leq \text{normal-exponent } x$   
 $\text{normal-exponent } x \leq 2^{\wedge}(\text{LENGTH}(e) - 1) - \text{int } \text{LENGTH}(f) - 1$   
**if** *is-normal*  $x$   
**for**  $x::('e, 'f)\text{float}$   
 ⟨proof⟩

**lemmas** *of-int-leI = of-int-le-iff*[THEN *iffD2*]

**lemma** *normal-exponent-bounds-real*:  
 $2 - 2^{\wedge}(\text{LENGTH}(e) - 1) - \text{real } \text{LENGTH}(f) \leq \text{normal-exponent } x$   
 $\text{normal-exponent } x \leq 2^{\wedge}(\text{LENGTH}(e) - 1) - \text{real } \text{LENGTH}(f) - 1$   
**if** *is-normal*  $x$   
**for**  $x::('e, 'f)\text{float}$   
 ⟨proof⟩

**lemma** *float-eqI*:  
 $x = y$  **if** *sign*  $x = \text{sign } y$  *fraction*  $x = \text{fraction } y$  *exponent*  $x = \text{exponent } y$   
 ⟨proof⟩

**lemma** *float-induct*[*induct type:float, case-names normal denormal neg zero infinity nan*]:  
**fixes**  $a::('e, 'f)\text{float}$   
**assumes** *normal*:  
 $\bigwedge x. \text{is-normal } x \implies \text{valof } x = \text{normal-mantissa } x * 2^{\text{powr } \text{normal-exponent } x}$   
 $\implies P x$   
**assumes** *denormal*:  
 $\bigwedge x. \text{is-denormal } x \implies$   
 $\text{valof } x = \text{denormal-mantissa } x * 2^{\text{powr } \text{denormal-exponent } \text{TYPE}((('e, 'f)\text{float}))}$   
 $\implies$   
 $P x$   
**assumes** *zero*:  $P 0$  *minus-zero*  
**assumes** *infty*:  $P \text{ plus-infinity } P \text{ minus-infinity}$   
**assumes** *nan*:  $\bigwedge x. \text{is-nan } x \implies P x$   
**shows**  $P a$   
 ⟨proof⟩

**lemma** *infinite-infinity*[simp]:  $\neg \text{is-finite plus-infinity} \neg \text{is-finite minus-infinity}$

*<proof>*

**lemma** *nan-not-finite[simp]*: *is-nan x*  $\implies \neg$  *is-finite x*  
*<proof>*

**lemma** *valof-nonneg*:  
*valof x*  $\geq 0$  **if** *sign x = 0* **for** *x::('e, 'f)float*  
*<proof>*

**lemma** *valof-nonpos*:  
*valof x*  $\leq 0$  **if** *sign x = 1* **for** *x::('e, 'f)float*  
*<proof>*

**lemma** *real-le-intI*: *x*  $\leq$  *y* **if** *floor x*  $\leq$  *floor y* *x*  $\in$   $\mathbb{Z}$  **for** *x y::real*  
*<proof>*

**lemma** *real-of-int-le-2-powr-bitlenI*:  
*real-of-int x*  $\leq 2$  *powr n - 1* **if** *bitlen (abs x)*  $\leq$  *m* *m*  $\leq$  *n*  
*<proof>*

**lemma** *largest-eq*:  
*largest TYPE (('e, 'f)float)* =  
 $(2^{\wedge}(\text{LENGTH('f)} + 1) - 1) * 2$  *powr real-of-int*  $(2^{\wedge}(\text{LENGTH('e)} - 1) -$   
*int LENGTH('f) - 1)  
*<proof>**

**lemma** *bitlen-denormal-mantissa*:  
*bitlen (abs (denormal-mantissa x))*  $\leq$  *LENGTH('f)* **for** *x::('e, 'f)float*  
*<proof>*

**lemma** *float-le-topfloat*:  
**fixes** *a::('e, 'f)float*  
**assumes** *is-finite a* *LENGTH('e) > 1*  
**shows** *a*  $\leq$  *topfloat*  
*<proof>*

**lemma** *float-val-le-largest*:  
*valof a*  $\leq$  *largest TYPE (('e, 'f)float)*  
**if** *is-finite a* *LENGTH('e) > 1*  
**for** *a::('e, 'f)float*  
*<proof>*

**lemma** *float-val-lt-threshold*:  
*valof a*  $<$  *threshold TYPE (('e, 'f)float)*  
**if** *is-finite a* *LENGTH('e) > 1*  
**for** *a::('e, 'f)float*  
*<proof>*

### 3.3 Algebraic properties about basic arithmetic

Commutativity of addition.

**lemma**

**assumes** *is-finite a is-finite b*

**shows** *float-plus-comm-eq:  $a + b = b + a$*

**and** *float-plus-comm:  $is-finite (a + b) \implies (a + b) \doteq (b + a)$*

*<proof>*

The floating-point number  $a$  falls into the same category as the negation of  $a$ .

**lemma** *is-zero-uminus[simp]:  $is-zero (- a) \longleftrightarrow is-zero a$*

*<proof>*

**lemma** *is-infinity-uminus [simp]:  $is-infinity (- a) = is-infinity a$*

*<proof>*

**lemma** *is-finite-uminus[simp]:  $is-finite (- a) \longleftrightarrow is-finite a$*

*<proof>*

**lemma** *is-nan-uminus[simp]:  $is-nan (- a) \longleftrightarrow is-nan a$*

*<proof>*

The sign of  $a$  and the sign of  $a$ 's negation is different.

**lemma** *float-neg-sign:  $(sign a) \neq (sign (- a))$*

*<proof>*

**lemma** *float-neg-sign1:  $sign a = sign (- b) \longleftrightarrow sign a \neq sign b$*

*<proof>*

**lemma** *valof-uminus:*

**assumes** *is-finite a*

**shows** *valof (- a) = - valof a (is ?L = ?R)*

*<proof>*

Showing  $a + (- b) = a - b$ .

**lemma** *float-neg-add:*

*is-finite a  $\implies is-finite b \implies is-finite (a - b) \implies valof a + valof (- b) = valof a - valof b$*

*<proof>*

**lemma** *float-plus-minus:*

**assumes** *is-finite a is-finite b is-finite (a - b)*

**shows**  *$(a + - b) \doteq (a - b)$*

*<proof>*

**lemma** *finite-bottomfloat: is-finite bottomfloat*

*<proof>*

**lemma** *bottomfloat-eq-m-largest*:  $\text{valof } (\text{bottomfloat}::('e, 'f)\text{float}) = - \text{largest } \text{TYPE}((('e, 'f)\text{float}))$   
**if**  $\text{LENGTH}('e) > 1$   
 $\langle \text{proof} \rangle$

**lemma** *float-val-ge-bottomfloat*:  $\text{valof } a \geq \text{valof } (\text{bottomfloat}::('e, 'f)\text{float})$   
**if**  $\text{LENGTH}('e) > 1$  *is-finite*  $a$   
**for**  $a::('e, 'f)\text{float}$   
 $\langle \text{proof} \rangle$

**lemma** *float-ge-bottomfloat*:  $\text{is-finite } a \implies a \geq \text{bottomfloat}$   
**if**  $\text{LENGTH}('e) > 1$  *is-finite*  $a$   
**for**  $a::('e, 'f)\text{float}$   
 $\langle \text{proof} \rangle$

**lemma** *float-val-ge-largest*:  
**fixes**  $a::('e, 'f)\text{float}$   
**assumes**  $\text{LENGTH}('e) > 1$  *is-finite*  $a$   
**shows**  $\text{valof } a \geq - \text{largest } \text{TYPE}((('e, 'f)\text{float}))$   
 $\langle \text{proof} \rangle$

**lemma** *float-val-gt-threshold*:  
**fixes**  $a::('e, 'f)\text{float}$   
**assumes**  $\text{LENGTH}('e) > 1$  *is-finite*  $a$   
**shows**  $\text{valof } a > - \text{threshold } \text{TYPE}((('e, 'f)\text{float}))$   
 $\langle \text{proof} \rangle$

Showing  $\text{abs } (- a) = \text{abs } a$ .

**lemma** *float-abs [simp]*:  $\neg \text{is-nan } a \implies \text{abs } (- a) = \text{abs } a$   
 $\langle \text{proof} \rangle$

**lemma** *neg-zerosign*:  $- (\text{zerosign } s a) = \text{zerosign } (1 - s) (- a)$   
 $\langle \text{proof} \rangle$

### 3.4 Properties about Rounding Errors

**definition** *error* ::  $('e, 'f)\text{float}$  *itself*  $\Rightarrow \text{real} \Rightarrow \text{real}$   
**where**  $\text{error } - x = \text{valof } (\text{round } \text{To-nearest } x::('e, 'f)\text{float}) - x$

**lemma** *bound-at-worst-lemma*:  
**fixes**  $a::('e, 'f)\text{float}$   
**assumes** *threshold*:  $|x| < \text{threshold } \text{TYPE}((('e, 'f)\text{float}))$   
**assumes** *finite*: *is-finite*  $a$   
**shows**  $|\text{valof } (\text{round } \text{To-nearest } x::('e, 'f)\text{float}) - x| \leq |\text{valof } a - x|$   
 $\langle \text{proof} \rangle$

**lemma** *error-at-worst-lemma*:  
**fixes**  $a::('e, 'f)\text{float}$

**assumes** *threshold*:  $|x| < \text{threshold } \text{TYPE}('e, 'f)\text{float}$   
**and** *is-finite* *a*  
**shows**  $|\text{error } \text{TYPE}('e, 'f)\text{float } x| \leq |\text{valof } a - x|$   
 $\langle \text{proof} \rangle$

**lemma** *error-is-zero* [*simp*]:  
**fixes** *a*::('e, 'f)float  
**assumes** *is-finite* *a*  $1 < \text{LENGTH}('e)$   
**shows**  $\text{error } \text{TYPE}('e, 'f)\text{float } (\text{valof } a) = 0$   
 $\langle \text{proof} \rangle$

**lemma** *is-finite-zerosign*[*simp*]:  $\text{is-finite } (\text{zerosign } s \ a) \longleftrightarrow \text{is-finite } a$   
 $\langle \text{proof} \rangle$

**lemma** *is-finite-closest*:  $\text{is-finite } (\text{closest } (v::\Rightarrow\text{real}) \ p \ \{a. \ \text{is-finite } a\} \ x)$   
 $\langle \text{proof} \rangle$

**lemma** *defloat-float-zerosign-round-finite*:  
**assumes** *threshold*:  $|x| < \text{threshold } \text{TYPE}('e, 'f)\text{float}$   
**shows**  $\text{is-finite } (\text{zerosign } s \ (\text{round } \text{To-nearest } x::('e, 'f)\text{float}))$   
 $\langle \text{proof} \rangle$

**lemma** *valof-zero*[*simp*]:  $\text{valof } 0 = 0 \ \text{valof } \text{minus-zero} = 0$   
 $\langle \text{proof} \rangle$

**lemma** *signzero-zero*:  
 $\text{is-zero } a \implies \text{valof } (\text{zerosign } s \ a) = 0$   
 $\langle \text{proof} \rangle$

**lemma** *val-zero*:  $\text{is-zero } a \implies \text{valof } a = 0$   
 $\langle \text{proof} \rangle$

**lemma** *float-add*:  
**fixes** *a b*::('e, 'f)float  
**assumes** *is-finite* *a*  
**and** *is-finite* *b*  
**and** *threshold*:  $|\text{valof } a + \text{valof } b| < \text{threshold } \text{TYPE}('e, 'f)\text{float}$   
**shows** *finite-float-add*:  $\text{is-finite } (a + b)$   
**and** *error-float-add*:  $\text{valof } (a + b) = \text{valof } a + \text{valof } b + \text{error } \text{TYPE}('e, 'f)\text{float } (\text{valof } a + \text{valof } b)$   
 $\langle \text{proof} \rangle$

**lemma** *float-sub*:  
**fixes** *a b*::('e, 'f)float  
**assumes** *is-finite* *a*  
**and** *is-finite* *b*  
**and** *threshold*:  $|\text{valof } a - \text{valof } b| < \text{threshold } \text{TYPE}('e, 'f)\text{float}$   
**shows** *finite-float-sub*:  $\text{is-finite } (a - b)$   
**and** *error-float-sub*:  $\text{valof } (a - b) = \text{valof } a - \text{valof } b + \text{error } \text{TYPE}('e, 'f)\text{float } (\text{valof } a - \text{valof } b)$

'f)float) (valof a - valof b)  
 ⟨proof⟩

**lemma float-mul:**

**fixes** a b::('e, 'f)float  
**assumes** is-finite a  
**and** is-finite b  
**and** threshold: |valof a \* valof b| < threshold TYPE(('e, 'f)float)  
**shows** finite-float-mul: is-finite (a \* b)  
**and** error-float-mul: valof (a \* b) = valof a \* valof b + error TYPE(('e, 'f)float)  
 (valof a \* valof b)  
 ⟨proof⟩

**lemma float-div:**

**fixes** a b::('e, 'f)float  
**assumes** is-finite a  
**and** is-finite b  
**and** not-zero: ¬ is-zero b  
**and** threshold: |valof a / valof b| < threshold TYPE(('e, 'f)float)  
**shows** finite-float-div: is-finite (a / b)  
**and** error-float-div: valof (a / b) = valof a / valof b + error TYPE(('e, 'f)float)  
 (valof a / valof b)  
 ⟨proof⟩

**lemma valof-one[simp]:** valof (1 :: ('e, 'f) float) = of-bool (LENGTH('e) > 1)  
 ⟨proof⟩

**end**

**theory FP64**

**imports**

IEEE

Word-Lib.Word-64

**begin**

## 4 Concrete encodings

Floating point operations defined as operations on words. Called "fixed precision" (fp) word in HOL4.

**type-synonym** float64 = (11,52)float

**type-synonym** fp64 = 64 word

**lift-definition** fp64-of-float :: float64 ⇒ fp64 **is**

λ(s::1 word, e::11 word, f::52 word). word-cat s (word-cat e f::63 word) ⟨proof⟩

**lift-definition** float-of-fp64 :: fp64 ⇒ float64 **is**

λx. apsnd word-split (word-split x::1 word \* 63 word) ⟨proof⟩

**definition** rel-fp64 ≡ (λx (y::word64). x = float-of-fp64 y)

**definition**  $eq\_fp64::float64 \Rightarrow float64 \Rightarrow bool$  **where**  $[simp]: eq\_fp64 \equiv (=)$

**lemma**  $float\_of\_fp64\_inverse[simp]: fp64\_of\_float (float\_of\_fp64 a) = a$   
 $\langle proof \rangle$

**lemma**  $float\_of\_fp64\_inj\_iff[simp]: fp64\_of\_float r = fp64\_of\_float s \iff r = s$   
 $\langle proof \rangle$

**lemma**  $fp64\_of\_float\_inverse[simp]: float\_of\_fp64 (fp64\_of\_float a) = a$   
 $\langle proof \rangle$

**lemma**  $Quotientfp: Quotient eq\_fp64 fp64\_of\_float float\_of\_fp64 rel\_fp64$   
—  $eq\_fp64$  is a workaround to prevent a (failing – TODO: why?) code setup in  $setup\_lifting$ .  
 $\langle proof \rangle$

**setup-lifting**  $Quotientfp$

**lift-definition**  $fp64\_lessThan::fp64 \Rightarrow fp64 \Rightarrow bool$  **is**  
 $flt::float64 \Rightarrow float64 \Rightarrow bool$   $\langle proof \rangle$

**lift-definition**  $fp64\_lessEqual::fp64 \Rightarrow fp64 \Rightarrow bool$  **is**  
 $fle::float64 \Rightarrow float64 \Rightarrow bool$   $\langle proof \rangle$

**lift-definition**  $fp64\_greaterThan::fp64 \Rightarrow fp64 \Rightarrow bool$  **is**  
 $fgt::float64 \Rightarrow float64 \Rightarrow bool$   $\langle proof \rangle$

**lift-definition**  $fp64\_greaterEqual::fp64 \Rightarrow fp64 \Rightarrow bool$  **is**  
 $fge::float64 \Rightarrow float64 \Rightarrow bool$   $\langle proof \rangle$

**lift-definition**  $fp64\_equal::fp64 \Rightarrow fp64 \Rightarrow bool$  **is**  
 $feq::float64 \Rightarrow float64 \Rightarrow bool$   $\langle proof \rangle$

**lift-definition**  $fp64\_abs::fp64 \Rightarrow fp64$  **is**  
 $abs::float64 \Rightarrow float64$   $\langle proof \rangle$

**lift-definition**  $fp64\_negate::fp64 \Rightarrow fp64$  **is**  
 $uminus::float64 \Rightarrow float64$   $\langle proof \rangle$

**lift-definition**  $fp64\_sqrt::roundmode \Rightarrow fp64 \Rightarrow fp64$  **is**  
 $fsqrt::roundmode \Rightarrow float64 \Rightarrow float64$   $\langle proof \rangle$

**lift-definition**  $fp64\_add::roundmode \Rightarrow fp64 \Rightarrow fp64 \Rightarrow fp64$  **is**  
 $fadd::roundmode \Rightarrow float64 \Rightarrow float64 \Rightarrow float64$   $\langle proof \rangle$

**lift-definition**  $fp64\_sub::roundmode \Rightarrow fp64 \Rightarrow fp64 \Rightarrow fp64$  **is**  
 $fsub::roundmode \Rightarrow float64 \Rightarrow float64 \Rightarrow float64$   $\langle proof \rangle$

**lift-definition**  $fp64\text{-mul}::\text{roundmode} \Rightarrow fp64 \Rightarrow fp64 \Rightarrow fp64$  is  
 $fmul::\text{roundmode} \Rightarrow \text{float64} \Rightarrow \text{float64} \Rightarrow \text{float64}$   $\langle \text{proof} \rangle$

**lift-definition**  $fp64\text{-div}::\text{roundmode} \Rightarrow fp64 \Rightarrow fp64 \Rightarrow fp64$  is  
 $fdiv::\text{roundmode} \Rightarrow \text{float64} \Rightarrow \text{float64} \Rightarrow \text{float64}$   $\langle \text{proof} \rangle$

**lift-definition**  $fp64\text{-mul-add}::\text{roundmode} \Rightarrow fp64 \Rightarrow fp64 \Rightarrow fp64 \Rightarrow fp64$  is  
 $fmul\text{-add}::\text{roundmode} \Rightarrow \text{float64} \Rightarrow \text{float64} \Rightarrow \text{float64} \Rightarrow \text{float64}$   $\langle \text{proof} \rangle$

**end**

**theory** *Conversion-IEEE-Float*

**imports**

*HOL-Library.Float*

*IEEE-Properties*

*HOL-Library.Code-Target-Numeral*

**begin**

**definition**  $of\text{-finite} (x::('e, 'f)\text{float}) =$   
*(if is-normal x then (Float (normal-mantissa x) (normal-exponent x))*  
*else if is-denormal x then (Float (denormal-mantissa x) (denormal-exponent*  
 $TYPE(('e, 'f)\text{float}))$   
*else 0)*

**lemma**  $float\text{-val-of-finite}: is\text{-finite} x \Longrightarrow of\text{-finite} x = \text{valof } x$   
 $\langle \text{proof} \rangle$

**definition**  $is\text{-normal-Float}::('e, 'f)\text{float} \text{ itself} \Rightarrow \text{Float.float} \Rightarrow \text{bool}$  **where**  
 $is\text{-normal-Float} x f \longleftrightarrow$   
 $\text{mantissa } f \neq 0 \wedge$   
 $\text{bitlen } |\text{mantissa } f| \leq \text{fracwidth } x + 1 \wedge$   
 $-\text{int } (\text{bias } x) - \text{bitlen } |\text{mantissa } f| + 1 < \text{Float.exponent } f \wedge$   
 $\text{Float.exponent } f < 2^{\wedge}(\text{LENGTH}('e)) - \text{bitlen } |\text{mantissa } f| - \text{bias } x$

**definition**  $is\text{-denormal-Float}::('e, 'f)\text{float} \text{ itself} \Rightarrow \text{Float.float} \Rightarrow \text{bool}$  **where**  
 $is\text{-denormal-Float} x f \longleftrightarrow$   
 $\text{mantissa } f \neq 0 \wedge$   
 $\text{bitlen } |\text{mantissa } f| \leq 1 - \text{Float.exponent } f - \text{int } (\text{bias } x) \wedge$   
 $1 - 2^{\wedge}(\text{LENGTH}('e) - 1) - \text{int } \text{LENGTH}('f) < \text{Float.exponent } f$

**lemmas**  $is\text{-denormal-FloatD} =$

$is\text{-denormal-Float-def}[\text{THEN } iffD1, \text{THEN } conjunct1]$

$is\text{-denormal-Float-def}[\text{THEN } iffD1, \text{THEN } conjunct2]$

**definition**  $is\text{-finite-Float}::('e, 'f)\text{float} \text{ itself} \Rightarrow \text{Float.float} \Rightarrow \text{bool}$  **where**  
 $is\text{-finite-Float} x f \longleftrightarrow is\text{-normal-Float} x f \vee is\text{-denormal-Float} x f \vee f = 0$

**lemma**  $is\text{-finite-Float-eq}:$

$is\text{-finite-Float } TYPE(('e, 'f)\text{float}) f \longleftrightarrow$

$(let\ e = Float.exponent\ f;\ bm = bitlen\ (abs\ (mantissa\ f))$   
 $in\ bm \leq Suc\ LENGTH('f) \wedge$   
 $bm \leq 2^{\wedge}(LENGTH('e) - 1) - e \wedge$   
 $1 - 2^{\wedge}(LENGTH('e) - 1) - int\ LENGTH('f) < e)$   
 $\langle proof \rangle$

**lift-definition**  $normal-of-Float :: Float.float \Rightarrow ('e, 'f)float$   
**is**  $\lambda x.$   $let\ m = mantissa\ x;\ e = Float.exponent\ x\ in$   
 $(if\ m > 0\ then\ 0\ else\ 1,$   
 $word-of-int\ (e + int\ (bias\ TYPE(('e, 'f)float)) + bitlen\ |m| - 1),$   
 $word-of-int\ (|m| * 2^{\wedge}(Suc\ LENGTH('f) - nat\ (bitlen\ |m|)) - 2^{\wedge}(LENGTH('f))))$   
 $\langle proof \rangle$

**lemma**  $sign-normal-of-Float:sign\ (normal-of-Float\ x) = (if\ x > 0\ then\ 0\ else\ 1)$   
 $\langle proof \rangle$

**lemma**  $uint-word-of-int-bitlen-eq:$   
 $uint\ (word-of-int\ x::'a::len\ word) = x\ if\ bitlen\ x \leq LENGTH('a)\ x \geq 0$   
 $\langle proof \rangle$

**lemma**  $fraction-normal-of-Float:fraction\ (normal-of-Float\ x::('e, 'f)float) =$   
 $(nat\ |mantissa\ x| * 2^{\wedge}(Suc\ LENGTH('f) - nat\ (bitlen\ |mantissa\ x|)) - 2^{\wedge}$   
 $LENGTH('f))$   
**if**  $is-normal-Float\ TYPE(('e, 'f)float)\ x$   
 $\langle proof \rangle$

**lemma**  $exponent-normal-of-Float:exponent\ (normal-of-Float\ x::('e, 'f)float) =$   
 $nat\ (Float.exponent\ x + (bias\ TYPE(('e, 'f)float)) + bitlen\ |mantissa\ x| - 1)$   
**if**  $is-normal-Float\ TYPE(('e, 'f)float)\ x$   
 $\langle proof \rangle$

**lift-definition**  $denormal-of-Float :: Float.float \Rightarrow ('e, 'f)float$   
**is**  $\lambda x.$   $let\ m = mantissa\ x;\ e = Float.exponent\ x\ in$   
 $(if\ m \geq 0\ then\ 0\ else\ 1,\ 0,$   
 $word-of-int\ (|m| * 2^{\wedge}nat\ (e + bias\ TYPE(('e, 'f)float) + fracwidth\ TYPE(('e,$   
 $'f)float) - 1)))$   
 $\langle proof \rangle$

**lemma**  $sign-denormal-of-Float:sign\ (denormal-of-Float\ x) = (if\ x \geq 0\ then\ 0\ else\ 1)$   
 $\langle proof \rangle$

**lemma**  $exponent-denormal-of-Float:exponent\ (denormal-of-Float\ x::('e, 'f)float) =$   
 $0$   
 $\langle proof \rangle$

**lemma**  $fraction-denormal-of-Float:fraction\ (denormal-of-Float\ x::('e, 'f)float) =$   
 $(nat\ |mantissa\ x| * 2^{\wedge}nat\ (Float.exponent\ x + bias\ TYPE(('e, 'f)float) +$   
 $LENGTH('f) - 1))$

**if** *is-denormal-Float* *TYPE*((*'e*, *'f*)*float*) *x*  
⟨*proof*⟩

**definition** *of-finite-Float* :: *Float.float*  $\Rightarrow$  (*'e*, *'f*) *float* **where**  
*of-finite-Float* *x* = (if *is-normal-Float* *TYPE*((*'e*, *'f*)*float*) *x* then *normal-of-Float* *x*  
else if *is-denormal-Float* *TYPE*((*'e*, *'f*)*float*) *x* then *denormal-of-Float* *x*  
else 0)

**lemma** *valof-normal-of-Float*: *valof* (*normal-of-Float* *x*::(*'e*, *'f*)*float*) = *x*  
**if** *is-normal-Float* *TYPE*((*'e*, *'f*)*float*) *x*  
⟨*proof*⟩

**lemma** *valof-denormal-of-Float*: *valof* (*denormal-of-Float* *x*::(*'e*, *'f*)*float*) = *x*  
**if** *is-denormal-Float* *TYPE*((*'e*, *'f*)*float*) *x*  
⟨*proof*⟩

**lemma** *valof-of-finite-Float*:  
*is-finite-Float* (*TYPE*((*'e*, *'f*) *IEEE.float*)) *x*  $\implies$  *valof* (*of-finite-Float* *x*::(*'e*,  
*'f*)*float*) = *x*  
⟨*proof*⟩

**lemma** *is-normal-normal-of-Float*:  
*is-normal* (*normal-of-Float* *x*::(*'e*, *'f*)*float*) **if** *is-normal-Float* *TYPE*((*'e*, *'f*)*float*)  
*x*  
⟨*proof*⟩

**lemma** *is-denormal-denormal-of-Float*: *is-denormal* (*denormal-of-Float* *x*::(*'e*, *'f*)*float*)  
**if** *is-denormal-Float* *TYPE*((*'e*, *'f*)*float*) *x*  
⟨*proof*⟩

**lemma** *is-finite-of-finite-Float*: *is-finite* (*of-finite-Float* *x*)  
⟨*proof*⟩

**lemma** *Float-eq-zero-iff*: *Float* *m e* = 0  $\longleftrightarrow$  *m* = 0  
⟨*proof*⟩

**lemma** *bitlen-mantissa-Float*:  
**shows** *bitlen* |*mantissa* (*Float* *m e*)| = (if *m* = 0 then 0 else *bitlen* |*m*| + *e*) -  
*Float.exponent* (*Float* *m e*)  
⟨*proof*⟩

**lemma** *exponent-Float*:  
**shows** *Float.exponent* (*Float* *m e*) = (if *m* = 0 then 0 else *bitlen* |*m*| + *e*) -  
*bitlen* |*mantissa* (*Float* *m e*)|  
⟨*proof*⟩

**lemma** *is-normal-Float-normal*:  
*is-normal-Float* *TYPE*((*'e*, *'f*)*float*) (*Float* (*normal-mantissa* *x*) (*normal-exponent*

```

x))
  if is-normal x for x::('e, 'f)float
  <proof>

```

```

lemma is-denormal-Float-denormal:
  is-denormal-Float TYPE(('e, 'f)float)
  (Float (denormal-mantissa x) (denormal-exponent TYPE(('e, 'f)float)))
  if is-denormal x for x::('e, 'f)float
  <proof>

```

```

lemma is-finite-Float-of-finite: is-finite-Float TYPE(('e, 'f)float) (of-finite x) for
x::('e, 'f)float
  <proof>

```

```

end

```

## 5 Code Generation Setup for Floats

```

theory Double
  imports
    Conversion-IEEE-Float
    HOL-Library.Monad-Syntax
    HOL-Library.Code-Target-Numeral
begin

```

A type for code generation to SML/OCaml

```

typedef double = UNIV::(11, 52) float set <proof>

```

```

setup-lifting type-definition-double

```

```

instantiation double :: {uminus, plus, times, minus, zero, one, abs, ord, inverse} begin
lift-definition uminus-double::double  $\Rightarrow$  double is uminus <proof>
lift-definition plus-double::double  $\Rightarrow$  double  $\Rightarrow$  double is plus <proof>
lift-definition times-double::double  $\Rightarrow$  double  $\Rightarrow$  double is times <proof>
lift-definition divide-double::double  $\Rightarrow$  double  $\Rightarrow$  double is divide <proof>
lift-definition inverse-double::double  $\Rightarrow$  double is inverse <proof>
lift-definition minus-double::double  $\Rightarrow$  double  $\Rightarrow$  double is minus <proof>
lift-definition zero-double::double is 0 <proof>
lift-definition one-double::double is 1 <proof>
lift-definition less-eq-double::double  $\Rightarrow$  double  $\Rightarrow$  bool is (≤) <proof>
lift-definition less-double::double  $\Rightarrow$  double  $\Rightarrow$  bool is (<) <proof>
instance <proof>
end

```

```

lift-definition eq-double::double $\Rightarrow$ double $\Rightarrow$ bool is float-eq <proof>

```

```

lift-definition sqrt-double::double $\Rightarrow$ double is float-sqrt <proof>

```

**no-notation** *plus-infinity* ( $\infty$ )

**lift-definition** *infinity-double::double* ( $\infty$ ) **is** *plus-infinity*  $\langle$ *proof* $\rangle$

**lift-definition** *is-normal::double*  $\Rightarrow$  *bool* **is** *IEEE.is-normal*  $\langle$ *proof* $\rangle$

**lift-definition** *is-zero::double*  $\Rightarrow$  *bool* **is** *IEEE.is-zero*  $\langle$ *proof* $\rangle$

**lift-definition** *is-finite::double*  $\Rightarrow$  *bool* **is** *IEEE.is-finite*  $\langle$ *proof* $\rangle$

**lift-definition** *is-nan::double*  $\Rightarrow$  *bool* **is** *IEEE.is-nan*  $\langle$ *proof* $\rangle$

**code-printing type-constructor** *double*  $\rightarrow$   
*(SML) real* **and** *(OCaml) float*

**code-printing constant** *0* :: *double*  $\rightarrow$   
*(SML) 0.0* **and** *(OCaml) 0.0*

**declare** *zero-double.abs-eq*[*code del*]

**code-printing constant** *1* :: *double*  $\rightarrow$   
*(SML) 1.0* **and** *(OCaml) 1.0*

**declare** *one-double.abs-eq*[*code del*]

**code-printing constant** *eq-double* :: *double*  $\Rightarrow$  *double*  $\Rightarrow$  *bool*  $\rightarrow$   
*(SML) Real.==* (*(-:real)*, *(-)*) **and** *(OCaml) Pervasives.(=)*

**declare** *eq-double.abs-eq*[*code del*]

**code-printing constant** *Orderings.less-eq* :: *double*  $\Rightarrow$  *double*  $\Rightarrow$  *bool*  $\rightarrow$   
*(SML) Real.<=* (*(-)*, *(-)*) **and** *(OCaml) Pervasives.(<=)*

**declare** *less-double-def* [*code del*]

**code-printing constant** *Orderings.less* :: *double*  $\Rightarrow$  *double*  $\Rightarrow$  *bool*  $\rightarrow$   
*(SML) Real.<* (*(-)*, *(-)*) **and** *(OCaml) Pervasives.<*

**declare** *less-eq-double-def*[*code del*]

**code-printing constant** *(+)* :: *double*  $\Rightarrow$  *double*  $\Rightarrow$  *double*  $\rightarrow$   
*(SML) Real.+* (*(-)*, *(-)*) **and** *(OCaml) Pervasives.(+)*

**declare** *plus-double-def*[*code del*]

**code-printing constant** *(\*)* :: *double*  $\Rightarrow$  *double*  $\Rightarrow$  *double*  $\rightarrow$   
*(SML) Real.\** (*(-)*, *(-)*) **and** *(OCaml) Pervasives.(\*)*

**declare** *times-double-def* [*code del*]

**code-printing constant** *(-)* :: *double*  $\Rightarrow$  *double*  $\Rightarrow$  *double*  $\rightarrow$   
*(SML) Real.-* (*(-)*, *(-)*) **and** *(OCaml) Pervasives.(-)*

**declare** *minus-double-def* [*code del*]

**code-printing constant** *uminus* :: *double*  $\Rightarrow$  *double*  $\rightarrow$   
*(SML) Real.~* **and** *(OCaml) Pervasives.(~)*

**code-printing constant** *(/)* :: *double*  $\Rightarrow$  *double*  $\Rightarrow$  *double*  $\rightarrow$   
*(SML) Real./* (*(-)*, *(-)*) **and** *(OCaml) Pervasives.(/)*

```

declare divide-double-def [code del]

code-printing constant sqrt-double :: double ⇒ double ↯
  (SML) Math.sqrt and (OCaml) Pervasives.sqrt
declare sqrt-def[code del]

code-printing constant infinity-double :: double ↯
  (SML) Real.posInf
declare infinity-double.abs-eq[code del]

code-printing constant is-normal :: double ⇒ bool ↯
  (SML) Real.isNormal
declare [[code drop: is-normal]]

code-printing constant is-finite :: double ⇒ bool ↯
  (SML) Real.isFinite
declare [[code drop: is-finite]]

code-printing constant is-nan :: double ⇒ bool ↯
  (SML) Real.isNaN
declare [[code drop: is-nan]]

```

Mapping natural numbers to doubles.

```

fun float-of :: nat ⇒ double
where
  float-of 0 = 0
  | float-of (Suc n) = float-of n + 1

lemma float-of 2 < float-of 3 + float-of 4
  ⟨proof⟩

export-code float-of in SML

```

## 5.1 Conversion from int

```

lift-definition double-of-int::int ⇒ double is  $\lambda i.$  round To-nearest (real-of-int i)
  ⟨proof⟩

context includes integer.lifting begin
lift-definition double-of-integer::integer ⇒ double is double-of-int ⟨proof⟩
end

lemma float-of-int[code]:
  double-of-int i = double-of-integer (integer-of-int i)
  ⟨proof⟩

code-printing
  constant double-of-integer :: integer ⇒ double ↯ (SML) Real.fromLargeInt
declare [[code drop: double-of-integer]]

```

## 5.2 Conversion to and from software floats, extracting information

Need to trust a lot of code here...

**lemma** *is-finite-double-eq*:

```
is-finite-Float TYPE((11, 52)float) f  $\longleftrightarrow$ 
  (let e = Float.exponent f; bm = bitlen (abs (mantissa f))
   in (bm  $\leq$  53  $\wedge$  e + bm < 1025  $\wedge$  - 1075 < e))
<proof>
```

**code-printing**

```
code-module IEEE-Mantissa-Exponent  $\rightarrow$  (SML)
<
structure IEEE-Mantissa-Exponent =
struct
fun to-man-exp-double x =
  if Real.isFinite x
  then case Real.toManExp x of {man = m, exp = e} =>
    SOME (Real.floor (Real.* (m, Math.pow (2.0, 53.0))), IntInf.- (e, 53))
  else NONE
fun normfloat (m, e) =
  (if m mod 2 = 0 andalso m <> 0 then normfloat (m div 2, e + 1)
   else if m = 0 then (0, 0) else (m, e))
fun bitlen x = (if 0 < x then bitlen (x div 2) + 1 else 0)
fun is-finite-double-eq m e =
  let
    val (m, e) = normfloat (m, e)
    val bm = bitlen (abs m)
  in bm <= 53 andalso e + bm < 1025 andalso e > ~1075 end
fun from-man-exp-double m e =
  if is-finite-double-eq m e
  then SOME (Real.fromManExp {man = Real.fromLargeInt m, exp = e})
  else NONE
end
>
```

**lift-definition** *of-finite::double*  $\Rightarrow$  *Float.float* **is** *Conversion-IEEE-Float.of-finite*  
 <proof>

**definition** *man-exp-of-double::double*  $\Rightarrow$  (*integer \* integer*)*option* **where**  
*man-exp-of-double* d = (if *is-finite* d then let f = *of-finite* d in  
 Some (*integer-of-int* (mantissa f), *integer-of-int* (Float.exponent f)) else None)

**lift-definition** *of-finite-Float::Float.float*  $\Rightarrow$  *double* **is** *Conversion-IEEE-Float.of-finite-Float*  
 <proof>

**definition** *double-of-man-exp::integer*  $\Rightarrow$  *integer*  $\Rightarrow$  *double option* **where**  
*double-of-man-exp* m e = (let f = *Float* (*int-of-integer* m) (*int-of-integer* e) in  
 if *is-finite-Float* TYPE((11, 52)float) f

then Some (of-finite-Float f)  
 else None)

**code-printing**

**constant** man-exp-of-double :: double  $\Rightarrow$  (integer \* integer) option  $\rightarrow$   
 (SML) IEEE'-Mantissa'-Exponent.to'-man'-exp'-double |

**constant** double-of-man-exp :: integer  $\Rightarrow$  integer  $\Rightarrow$  double option  $\rightarrow$   
 (SML) IEEE'-Mantissa'-Exponent.from'-man'-exp'-double

**declare** [[code drop: man-exp-of-double]]

**declare** [[code drop: double-of-man-exp]]

**lift-definition** Float-of-double::double  $\Rightarrow$  Float.float option **is**

$\lambda x$ . if is-finite x then Some (of-finite x) else None <proof>

**lift-definition** double-of-Float::Float.float  $\Rightarrow$  double option **is**

$\lambda x$ . if is-finite-Float TYPE((11, 52)float) x then Some (of-finite-Float x) else None <proof>

**lemma** compute-Float-of-double[code]:

Float-of-double x =  
 map-option ( $\lambda(m, e)$ . Float (int-of-integer m) (int-of-integer e)) (man-exp-of-double x)  
 <proof>

**lemma** compute-double-of-Float[code]:

double-of-Float f = double-of-man-exp (integer-of-int (mantissa f)) (integer-of-int (Float.exponent f))  
 <proof>

**definition** check-conversion m e =

(let f = Float (int-of-integer m) (int-of-integer e) in  
 do {  
 d  $\leftarrow$  double-of-Float f;  
 Float-of-double d  
 } = (if is-finite-Float TYPE((11, 52)float) f then Some f else None))

**primrec** check-all::nat  $\Rightarrow$  (nat  $\Rightarrow$  bool)  $\Rightarrow$  bool **where**

check-all 0 P  $\longleftrightarrow$  True  
 | check-all (Suc i) P  $\longleftrightarrow$  P i  $\wedge$  check-all i P

**definition** check-conversions dm de =

check-all (nat (2 \* de)) ( $\lambda e$ . check-all (nat (2 \* dm)) ( $\lambda m$ .  
 check-conversion (integer-of-int (int m - dm)) (integer-of-int (int e - de))))

**lemma** check-conversions 100 100

<proof>

**end**

## References

- [1] J. Harrison. *Floating point verification in HOL light: the exponential function*. Springer, 1997.