

An Automatable Formal Semantics for IEEE-754 Floating-Point Arithmetic

Martin Brain, Cesare Tinelli, Philipp Rümmer, Thomas Wahl
(and the rest of the SMT community)

University of Oxford



June 24, 2015

Hasn't this been done before?

- Isabelle *A formal model of IEEE floating point arithmetic*
- HOL *Interpretation of IEEE-854 floating-point standard and definition in the HOL system.*
- HOL Light *Floating point verification in HOL light: The exponential function (Intel)*
- ACL2 *A mechanically checked proof of the AMD5K86TM floating-point division program (AMD and Centaur)*
- PVS *Defining the IEEE-854 floating-point standard in PVS*
- Coq *A generic library for floating-point numbers and its application to exact computing*
- Coq *Floating-point arithmetic in the Coq system*
- Coq *Flocq: A Unified Library for Proving Floating-point Algorithms in Coq*

Hasn't this been done before?

Isabelle *A formal model of IEEE floating point arithmetic*

HOL *Interpretation of IEEE-854 floating-point standard and definition in the HOL system.*

HOL Light *Floating point verification in HOL light: The exponential function (Intel)*

ACL2 *A mechanically checked proof of the AMD5K86TM floating-point division program (AMD and Centaur)*

PVS *Defining the IEEE-854 floating-point standard in PVS*

Coq *A generic library for floating-point numbers and its application to exact computing*

Coq *Floating-point arithmetic in the Coq system*

Coq *Flocq: A Unified Library for Proving Floating-point Algorithms in Coq*

... there is another way to think about
theorem-proving ...

Is there an x and y such that ...

$$0 < x$$

$$0 < y$$

$$x + y < x$$

Is there an x and y such that ...

$$\begin{array}{ccccc} & 0 & \clubsuit & & x \\ & & & & \\ & 0 & \clubsuit & & y \\ x & \spadesuit & y & \clubsuit & x \end{array}$$

Is there an x and y such that ...

$$\begin{array}{rcl} 0 & \clubsuit & x \\ 0 & \clubsuit & y \\ x \spadesuit y & \clubsuit & x \end{array}$$

It depends on the interpretation (of \clubsuit and \spadesuit)!

$$\begin{array}{rcl} D & = & \mathbb{Z} \\ [\clubsuit] & = & <\mathbb{Z} \\ [\spadesuit] & = & +\mathbb{Z} \end{array}$$

NO!

Is there an x and y such that ...

$$\begin{array}{rcl} 0 & \clubsuit & x \\ 0 & \clubsuit & y \\ x & \spadesuit & y \quad \clubsuit \quad x \end{array}$$

It depends on the interpretation (of \clubsuit and \spadesuit)!

$$\begin{aligned} D &= \{00, 01, 10, 11\} \\ [\clubsuit] &= \text{bvult} \\ [\spadesuit] &= \text{bvplus} \end{aligned}$$

Yes ($x = 01, y = 11$)

First Order Logic

Syntax

Fix a *signature* Σ
(i.e. $\Sigma = \{\clubsuit, \spadesuit\}$)

Semantics

An *interpretation* is
 $M = (D, \llbracket \cdot \rrbracket : \Sigma \rightarrow (2^{D^n}))$

Satisfiability

An interpretation M *satisfies* a formula ϕ :

$$M \models \phi$$

If ϕ evaluated over D (using $\llbracket \cdot \rrbracket$) is true.

How Do We Fix The *Meaning* of Symbols?

Option 1 – Axiomatic

$$M \models \text{Axioms} \Rightarrow \phi$$

$$\text{Axioms} = \forall a, b, c . a \clubsuit b \wedge b \clubsuit c \Rightarrow a \clubsuit c$$

$$= \forall a . \neg a \clubsuit a$$

...

Formalisation is solver *INPUT*.

Pros

- + Easy to implement
- + Flexible
- + Can add theorems

Cons

- All formulae quantified
- Axioms not always simple
- Hard to solve

How Do We Fix The *Meaning* of Symbols?

Option 2 – Algebraic

Fix signature Σ' and its interpretation $M' = (D, \llbracket \cdot \rrbracket : \Sigma' \rightarrow (2^{D^n}))$.

$$D = \mathbb{Z} \quad \llbracket \clubsuit \rrbracket = <_{\mathbb{Z}} \quad \llbracket \spadesuit \rrbracket = +_{\mathbb{Z}}$$

Is there M extension of M' such that:

$$M \models \phi$$

Formalisation is solver *SPECIFICATION*.

Pros

- + Fast decision procedures
- + Counter-examples
- + Few quantifiers

Cons

- Theory has to be built into solver
- Implementation harder

How Do We Fix The *Meaning* of Symbols?

Option 2 – Algebraic

Fix signature Σ' and its interpretation $M' = (D, [\cdot]) : \Sigma' \rightarrow (2^{D^n})$.

$$D = \mathbb{Z} \quad [\clubsuit] = <_{\mathbb{Z}} \quad [\spadesuit] = +_{\mathbb{Z}}$$

Is there M extension of M' such that:

$$M \models \phi$$

Formalisation is solver *SPECIFICATION*.

Pros

- + Fast decision procedures
- + Counter-examples
- + Few quantifiers

Cons

- Theory has to be built into solver
- Implementation harder

SAT Modulo Theory (SMT)

- The major school of *algebraic* solvers.
- Theories = specifications of (sets of) interpretations.
- SMT-LIB : international standard for SMT solvers.
- Mature implementations :
CVC4, Z3, MathSAT, Yices, STP, Boolector, OpenSMT, ...
- Near ubiquitous in software verification.

- 1 SMT
- 2 SMT-LIB Theory of Floating-Point
- 3 Conclusions

Requirements

Principles

Bit-Exact Must do *exactly* what the hardware does

Precise Gives SAT / UNSAT (ideally with model / proof)

Automated Ideally fast and “out of the box”

Flexible Support different decision procedures

Target Applications

- Path feasibility / test-case generation
- Generation of special values
- Numerical instability
- Undefined behaviour
- Hardware verification
- Functional correctness
- Automated numerical analysis

IEEE-754 2008

IEEE Std 754-2008
IEEE Standard for Floating-Point Arithmetic

3.2 Specification levels

Floating-point arithmetic is a systematic approximation of real arithmetic, as illustrated in Table 3.1. Floating-point arithmetic can only represent a finite subset of the continuum of real numbers. Consequently certain properties of real arithmetic, such as associativity of addition, do not always hold for floating-point arithmetic.

Table 3.1—Relationships between different specification levels for a particular format

Level 1	$\{-\infty \dots 0 \dots +\infty\}$	Extended real numbers.
many-to-one ↓	<i>rounding</i>	↑ projection (except for NaN)
Level 2	$\{-\infty \dots -0\} \cup \{+0 \dots +\infty\} \cup \text{NaN}$	Floating-point data—an algebraically closed system.
one-to-many ↓	<i>representation specification</i>	↑ many-to-one
Level 3	$(\textit{sign}, \textit{exponent}, \textit{significand}) \cup \{-\infty, +\infty\} \cup \text{qNaN} \cup \text{sNaN}$	Representations of floating-point data.
one-to-many ↓	<i>encoding for representations of floating-point data</i>	↑ many-to-one
Level 4	0111000...	Bit strings.

The mathematical structure underpinning the arithmetic in this standard is the extended reals, that is, the set of real numbers together with positive and negative infinity. For a given format, the process of *rounding* (see 4) maps an extended real number to a *floating-point number* included in that format. A *floating-point datum*, which can be a signed zero, finite non-zero number, signed infinity, or a NaN (not-a-number), can be mapped to one or more *representations of floating-point data* in a format.

Level 1 : Extended Reals

$$\mathbb{R}^* = \mathbb{R} \cup \{+\infty, -\infty, \text{NaN}\}$$

(partially ordered, additive and multiplicative commutative monoid with the distributivity property)

	$+\infty \leq w \Leftrightarrow w = +\infty$
$u + \text{NaN} = \text{NaN} + u = \text{NaN}$	$w \leq -\infty \Leftrightarrow w = -\infty$
$-\text{NaN} = \text{NaN}$	$-(+\infty) = -\infty$
$u \cdot \text{NaN} = \text{NaN} \cdot u = \text{NaN}$	$-(-\infty) = +\infty$
$\text{NaN}^{-1} = \text{NaN}$	$+\infty^{-1} = 0$
$\text{NaN} \leq u \Leftrightarrow u = \text{NaN}$	$-\infty^{-1} = 0$
$u \leq \text{NaN} \Leftrightarrow u = \text{NaN}$	$0^{-1} = +\infty$

...

Level 2(ish) : Domain

$$\mathbb{F}_{\varepsilon,\sigma} = \mathbb{F}_{\varepsilon,\sigma} \cup \{\text{NaN}\}$$

$$\mathbb{F}_{\varepsilon,\sigma} = \text{FZ}_{\varepsilon,\sigma} \cup \text{FS}_{\varepsilon,\sigma} \cup \text{FN}_{\varepsilon,\sigma} \cup \text{FI}_{\varepsilon,\sigma}$$

$$\text{FZ}_{\varepsilon,\sigma} = \{(s, e, m) \in B_{\varepsilon,\sigma} \mid e = \mathbf{0}_{\varepsilon}, m = \mathbf{0}_{\sigma-1}\}$$

$$\text{FS}_{\varepsilon,\sigma} = \{(s, e, m) \in B_{\varepsilon,\sigma} \mid e = \mathbf{0}_{\varepsilon}, m \neq \mathbf{0}_{\sigma-1}\}$$

$$\text{FN}_{\varepsilon,\sigma} = \{(s, e, m) \in B_{\varepsilon,\sigma} \mid e \neq \mathbf{1}_{\varepsilon}, e \neq \mathbf{0}_{\varepsilon}\}$$

$$\text{FI}_{\varepsilon,\sigma} = \{(s, e, m) \in B_{\varepsilon,\sigma} \mid e = \mathbf{1}_{\varepsilon}, m = \mathbf{0}_{\sigma-1}\}$$

$$v_{\varepsilon,\sigma} : \mathbb{F}_{\varepsilon,\sigma} \rightarrow \mathbb{R}^*$$

IEEE-754 2008 again

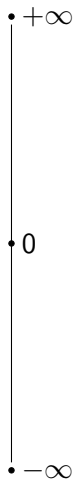
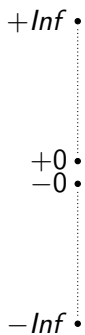
IEEE Std 754-2008
IEEE Standard for Floating-Point Arithmetic

5. Operations

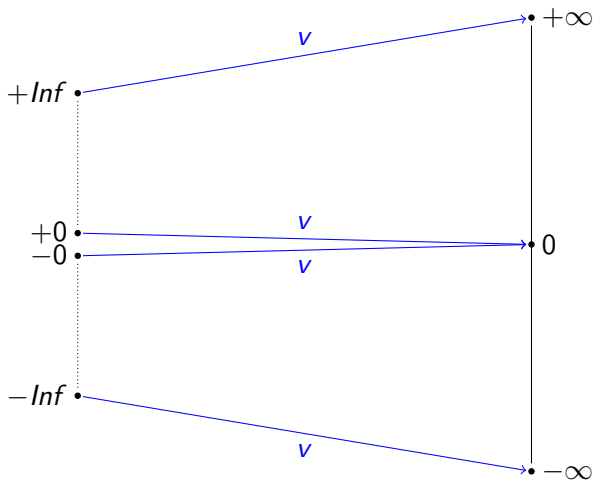
5.1 Overview

All conforming implementations of this standard shall provide the operations listed in this clause for all supported arithmetic formats, except as stated below. Each of the computational operations that return a numeric result specified by this standard shall be performed as if it first produced an intermediate result correct to infinite precision and with unbounded range, and then rounded that intermediate result, if necessary, to fit in the destination's format (see 4 and 7). Clause 6 augments the following specifications to cover ± 0 , $\pm\infty$, and NaN. Clause 7 describes default exception handling.

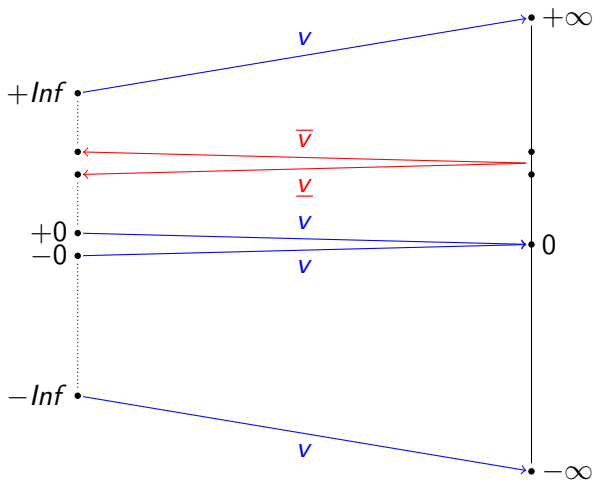
Upper and Lower Adjoints



Upper and Lower Adjoints



Upper and Lower Adjoints



Rounding is (Just) Selecting Between Adjoints!

$$\text{rnd}(v, \textit{mode}, \textit{sz}, r) = \bar{v}(r) \text{ or } \underline{v}(r)$$

This allows us to round to any format of float, bit-vectors, \mathbb{Z} , integer valued floats, ...

Operations

$$\text{add}_{\varepsilon,\sigma}(rm, f, g) = \text{rnd}(v, rm, \text{addSign}(rm, f, g), v(f) + v(g))$$

$$\text{sub}_{\varepsilon,\sigma}(rm, f, g) = \text{rnd}(v, rm, \text{subSign}(rm, f, g), v(f) - v(g))$$

$$\text{mul}_{\varepsilon,\sigma}(rm, f, g) = \text{rnd}(v, rm, \text{xorSign}(f, g), v(f) * v(g))$$

$$\text{div}_{\varepsilon,\sigma}(rm, f, g) = \begin{cases} \text{neg}_{\varepsilon,\sigma}(\text{rnd}(v, rm, \top, -(v(f)/v(g)))) & \text{xorSign}(f, g) \\ \text{rnd}(v, rm, \perp, v(f)/v(g)) & \neg\text{xorSign}(f, g) \end{cases}$$

$$\text{fma}_{\varepsilon,\sigma}(rm, f, g, h) = \text{rnd}(v, rm, \text{fmaSign}(rm, f, g, h), (v(f) * v(g)) + v(h))$$

Limitations and Omissions

- No decimal floats
- Only one NaN (no signaling / quiet, no payload)
- No exceptions
- No attributes
- No trigonometric functions

Implementations

	Bit-blast	ACDL	Axiomatic
CVC4	(✓)	(✓)	
Z3	✓		
MathSAT	✓	✓	
Sonolar	✓		
Alt-Ergo			(✓)
CBMC	✓		

- 1 SMT
- 2 SMT-LIB Theory of Floating-Point
- 3 Conclusions**

Help Needed!

- Correctness
- Examples
(edge cases, tests, challenge problems)
- “Diamond free” circuits
(multiply, divide, shift, float add, normalise)
- Elementary functions
- Floating-point remainder

Conclusions

- 1 Formalisation as input (axiomatic) vs. formalisation as specification (algebraic)
- 2 Rounding as choice of adjoints.
- 3 Have a specification (and implementations) of an SMT-LIB standard of floating-point.

Conclusions

- 1 Formalisation as input (axiomatic) vs. formalisation as specification (algebraic)
- 2 Rounding as choice of adjoints.
- 3 Have a specification (and implementations) of an SMT-LIB standard of floating-point.

Thank you for your time and attention.

Made using only Free Software