

Floating point verification in HOL Light: the exponential function

John Harrison

University of Cambridge

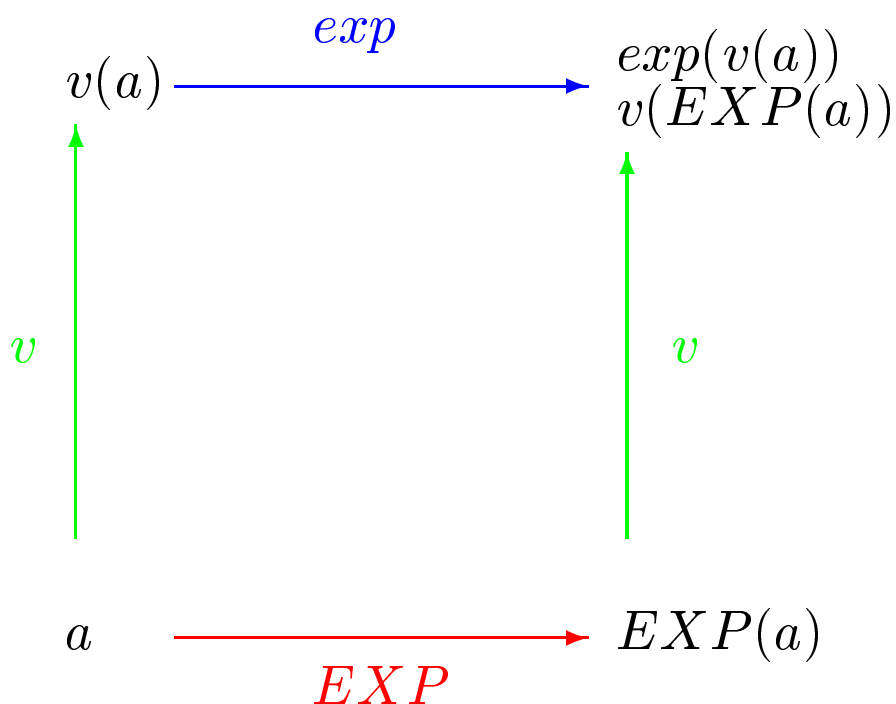
- Introduction
- Floating point correctness
- Our implementation language
- The algorithm
- Outline of the HOL proof
- General conclusions

Introduction

- Floating point algorithms are fairly small, but often complicated mathematically.
- There have been errors in commercial systems, e.g. the Pentium FDIV bug in 1994.
- In the case of transcendental functions it's difficult even to say what correctness *means*.
- Verification using model checkers is difficult because of the need for mathematical apparatus.
- It can even be difficult using theorem provers since not many of them have good theories of real numbers etc.

Floating point correctness

We want to specify the correctness according to the following diagram:



We measure the difference between $v(\overline{EXP}(a))$ and $\text{exp}(v(a))$ in ‘units in the last place’ of $\overline{EXP}(a)$.

Our implementation language

This includes the following constructs:

```
command = variable := expression
          | command ; command
          | if expression then command
            else command
          | if expression then command
          | while expression do command
          | do command while expression
          | skip
          | { expression }
```

We define a simple relational semantics in HOL, and derive weakest preconditions and total correctness rules. We then prove total correctness via VC generation.

The idea is that this language can be formally linked to C, Verilog, Handel, ...

The algorithm

The algorithm we verify is taken from a paper by Tang in *ACM Transactions on Mathematical Software*, 1989.

Similar techniques are widely used for floating point libraries, and, probably, for hardware implementations.

The algorithm relies on a table of precomputed constants. Tang's paper gives actual values as hex representations of IEEE numbers.

The algorithm works in three phases:

- Perform range reduction
- Use polynomial approximation
- Reconstruct answer using tables

The correctness proof reflects this.

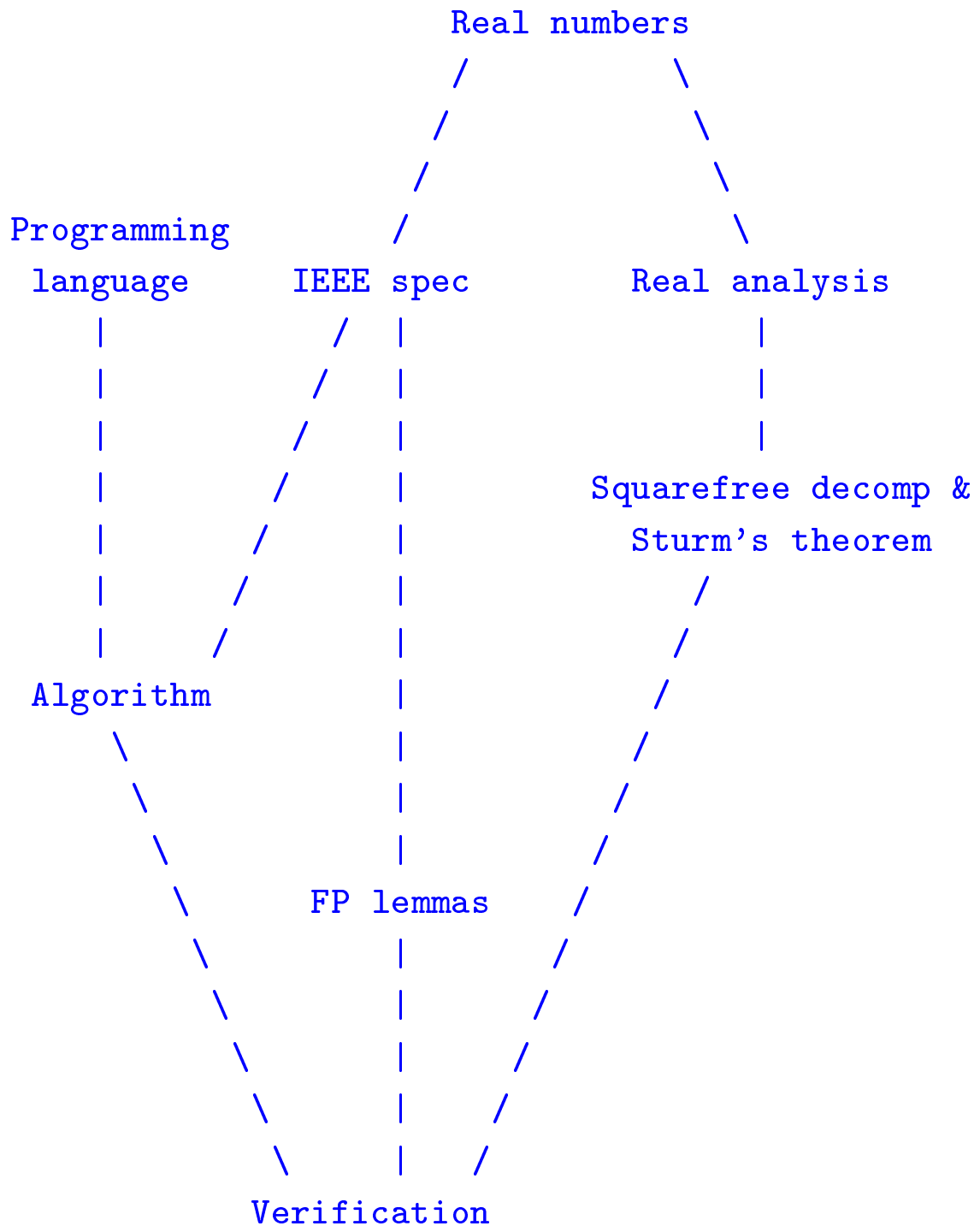
Code for the algorithm

```

if Isnan(X) then E := X
else if X == Plus_infinity then E := Plus_infinity
else if X == Minus_infinity then E := Plus_zero
else if abs(X) > THRESHOLD_1 then
  if X > Plus_zero then E := Plus_infinity
  else E := Plus_zero
else if abs(X) < THRESHOLD_2 then E := Plus_one + X
else
  (N := INTRND(X * Inv_L);
  N2 := N % Int_32;
  N1 := N - N2;
  if abs(N) >= Int_2e9 then
    R1 := (X - Tofloat(N1) * L1) - Tofloat(N2) * L1
  else
    R1 := X - Tofloat(N) * L1;
  R2 := Tofloat(--N) * L2;
  M := N1 / Int_32;
  J := N2;
  R := R1 + R2;
  Q := R * R * (A1 + R * A2);
  P := R1 + (R2 + Q);
  S := S_Lead(J) + S_Trail(J);
  E1 := S_Lead(J) + (S_Trail(J) + S * P);
  E := Scalb(E1,M)
)

```

Structure of the HOL proof



Floating point lemmas (1)

We define the error $\text{error}(x)$ resulting from rounding a real number x to a floating point value.

Because of the regular way in which the operations are defined, all the operations then relate to their abstract mathematical counterparts according to the same pattern:

$$\begin{aligned} &|- \text{Finite}(a) \wedge \text{Finite}(b) \wedge \\ &\quad \text{abs}(\text{Val}(a) + \text{Val}(b)) < \text{threshold}(\text{float_format}) \\ &\implies \text{Finite}(a + b) \wedge \\ &\quad (\text{Val}(a + b) = (\text{Val}(a) + \text{Val}(b)) + \\ &\quad \text{error}(\text{Val}(a) + \text{Val}(b))) \end{aligned}$$

The comparisons are even more straightforward:

$$\begin{aligned} &|- \text{Finite}(a) \wedge \text{Finite}(b) \\ &\implies (a < b = \text{Val}(a) < \text{Val}(b)) \end{aligned}$$

Floating point lemmas (2)

We have several lemmas quantifying the error, e.g.

$$\begin{aligned} &|- \text{abs}(x) < \text{threshold}(\text{float_format}) \wedge \\ &\quad \text{abs}(x) < (2^j / 2^{125}) \\ &\implies \text{abs}(\text{error}(x)) \leq 2^j / 2^{150} \end{aligned}$$

There are many important situations, however, where the operations are exact, because the result is exactly representable, e.g. subtraction of nearby values with the same sign:

$$\begin{aligned} &|- \text{Finite}(a) \wedge \text{Finite}(b) \wedge \\ &\quad 2 * \text{abs}(\text{Val}(a) - \text{Val}(b)) \leq \text{abs}(\text{Val}(a)) \\ &\implies \text{Finite}(a - b) \wedge \\ &\quad (\text{Val}(a - b) = \text{Val}(a) - \text{Val}(b)) \end{aligned}$$

This is a classic result in floating point error analysis.

Informal error analysis

Tang's error analysis translates quite directly into HOL. One needs to:

1. Prove that clever implementation tricks ensure certain remainder terms are calculated exactly. This relies on cancellation, and the fact that pre-stored constants have trailing zeroes.
2. Prove that the polynomial approximation obeys the appropriate error bounds.
3. Prove that the rounding errors when reconstructing the final answer do not get too large.

In Tang's paper, 1 is quite brief, 2 is dismissed in a few lines, while 3 is given a long and detailed proof.

HOL error analysis

In the HOL version, this order of difficulty is reversed!

1. The first part is not fundamentally difficult, but quite tricky because it involves a lot of special cases and low-level proofs.
2. The second part involves numerical approximation, which needs a lot of work to translate into a formal proof (e.g. Taylor series, Sturm's theorem ...). In fact Tang makes a small mistake here, though it doesn't affect the final result.
3. The last part is quite routine, and we can program HOL to compose the rounding errors automatically. Actually, we derive better bounds than Tang does since we avoid making simplifying assumptions to cut down the work.

The final result

Under the various ‘definitional’ assumptions, we confirm Tang’s bottom-line result:

$$\begin{aligned}
 & (\text{Isnan}(X) \implies \text{Isnan}(E)) \wedge \\
 & (X == \text{Plus_infinity} \vee \\
 & \quad \text{Finite}(X) \wedge \\
 & \quad \text{exp}(\text{Val } X) \geq \text{threshold}(\text{float_format}) \\
 & \implies E == \text{Plus_infinity}) \wedge \\
 & (X == \text{Minus_infinity} \implies E == \text{Plus_zero}) \wedge \\
 & (\text{Finite}(X) \wedge \text{exp}(\text{Val } X) < \text{threshold}(\text{float_format}) \\
 & \implies \text{Isnormal}(E) \wedge \\
 & \quad \text{abs}(\text{Val}(E) - \text{exp}(\text{Val } X)) < (\&54 / \&100) * \text{Ulp}(E) \\
 & \quad \vee (\text{Isdenormal}(E) \vee \text{Iszero}(E)) \wedge \\
 & \quad \text{abs}(\text{Val}(E) - \text{exp}(\text{Val } X)) < (\&77 / \&100) * \text{Ulp}(E))
 \end{aligned}$$

In fact, this specification is a bit more precise than Tang’s, e.g. we are explicit about the overflow threshold.

Conclusions

- We confirm (and strengthen) the main results of the hand proof. But we detect a few slips and uncover subtle issues. This class of proofs is a good target for verification.
- The proof was very long (over 3 months of work), but most of this was devoted to general results that could be re-used.
- It's a mistake to believe that only 'trivial' mathematics is needed for verification applications. HOL Light's mathematical theories are essential.
- Automation of linear arithmetic is practically indispensable. Better tools for nonlinear reasoning are needed.
- The proof runtimes are very long owing to the extensive use of arithmetic done by inference.