Numeric Programs

Guillaume Melquiond

MPRI 2-36-1 "Preuve de Programme"

January 23rd, 2013

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

ML-like programs:

- ML-like programs:
 - mutable variables,

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへで

- ML-like programs:
 - mutable variables,
 - recursive functions,

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで

arrays,

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

- arrays,
- algebraic datatypes.

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,
 - arrays,
 - algebraic datatypes.
- Program specification and verification:

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,
 - arrays,
 - algebraic datatypes.
- Program specification and verification:
 - Hoare logic: safety, validity, termination,

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,
 - arrays,
 - algebraic datatypes.
- Program specification and verification:
 - Hoare logic: safety, validity, termination,

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

weakest precondition computations,

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,
 - arrays,
 - algebraic datatypes.
- Program specification and verification:
 - Hoare logic: safety, validity, termination,
 - weakest precondition computations,
 - modular verification: function contract,

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,
 - arrays,
 - algebraic datatypes.
- Program specification and verification:
 - Hoare logic: safety, validity, termination,
 - weakest precondition computations,
 - modular verification: function contract,

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

abstract types and axioms,

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,
 - arrays,
 - algebraic datatypes.
- Program specification and verification:
 - Hoare logic: safety, validity, termination,
 - weakest precondition computations,
 - modular verification: function contract,

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

- abstract types and axioms,
- logic functions and predicates,

- ML-like programs:
 - mutable variables,
 - recursive functions,
 - basic types: integers, reals,
 - arrays,
 - algebraic datatypes.
- Program specification and verification:
 - Hoare logic: safety, validity, termination,
 - weakest precondition computations,
 - modular verification: function contract,

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

- abstract types and axioms,
- logic functions and predicates,
- recursive and inductive definitions.

► 32-bit signed integers in two-complement:

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

► 32-bit signed integers in two-complement:



► 32-bit signed integers in two-complement:

- ▶ $1 + 1 \rightarrow 2$,
- ▶ $2147483647 + 1 \rightarrow -2147483648$,

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

► 32-bit signed integers in two-complement:

- ► 1 + 1 → 2,
- ▶ $2147483647 + 1 \rightarrow -2147483648$,

(ロ) (同) (三) (三) (三) (○) (○)

▶ $100000^2 \rightarrow 1410065408$,

► 32-bit signed integers in two-complement:

- ▶ $1 + 1 \rightarrow 2$,
- ▶ $2147483647 + 1 \rightarrow -2147483648$,
- ▶ $100000^2 \rightarrow 1410065408$,
- ▶ $-2147483648 \mod -1 \rightarrow \text{boom}$ (floating-point exn?!).

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

32-bit signed integers in two-complement:

- ▶ $1 + 1 \rightarrow 2$,
- ▶ $2147483647 + 1 \rightarrow -2147483648$,
- ▶ $100000^2 \rightarrow 1410065408$,
- ▶ $-2147483648 \mod -1 \rightarrow \text{boom}$ (floating-point exn?!).

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

IEEE-754 binary64 floating-point numbers:

32-bit signed integers in two-complement:

- ▶ $1 + 1 \rightarrow 2$,
- ▶ $2147483647 + 1 \rightarrow -2147483648$,
- ▶ $100000^2 \rightarrow 1410065408$,
- ▶ $-2147483648 \mod -1 \rightarrow \text{boom}$ (floating-point exn?!).

(日) (日) (日) (日) (日) (日) (日)

- IEEE-754 binary64 floating-point numbers:
 - ▶ $2 \times 2 \times \cdots \times 2 \rightarrow +\infty$,

32-bit signed integers in two-complement:

- ▶ $1 + 1 \rightarrow 2$,
- ▶ $2147483647 + 1 \rightarrow -2147483648$,
- ▶ $100000^2 \rightarrow 1410065408$,
- ▶ $-2147483648 \mod -1 \rightarrow \text{boom}$ (floating-point exn?!).

(日) (日) (日) (日) (日) (日) (日)

IEEE-754 binary64 floating-point numbers:

•
$$2 \times 2 \times \cdots \times 2 \rightarrow +\infty$$
,

▶ $1 \div 0 \rightarrow +\infty$,

32-bit signed integers in two-complement:

- ▶ $1 + 1 \rightarrow 2$,
- ▶ $2147483647 + 1 \rightarrow -2147483648$,
- ▶ $100000^2 \rightarrow 1410065408$,
- ▶ $-2147483648 \mod -1 \rightarrow \text{boom}$ (floating-point exn?!).

(日) (日) (日) (日) (日) (日) (日)

- IEEE-754 binary64 floating-point numbers:
 - $2 \times 2 \times \cdots \times 2 \rightarrow +\infty$,
 - ▶ $1 \div 0 \rightarrow +\infty$,
 - ▶ $1 \div -0 \rightarrow -\infty$,

32-bit signed integers in two-complement:

- ▶ $1 + 1 \rightarrow 2$,
- ▶ $2147483647 + 1 \rightarrow -2147483648$,
- ▶ $100000^2 \rightarrow 1410065408$,
- ▶ $-2147483648 \mod -1 \rightarrow \text{boom}$ (floating-point exn?!).

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

- IEEE-754 binary64 floating-point numbers:
 - $2 \times 2 \times \cdots \times 2 \rightarrow +\infty$,
 - ▶ $1 \div 0 \rightarrow +\infty$,
 - ▶ 1 \div -0 \rightarrow -∞,
 - ▶ $0 \div 0 \rightarrow NaN.$

1983, truncation when computing an index of Vancouver Stock Exchange drops it to half its value after 6 months.

- 1983, truncation when computing an index of Vancouver Stock Exchange drops it to half its value after 6 months.
- 1987, the inflation in UK is computed with a rounding error: pensions are off by £100M for 21 months.

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

- 1983, truncation when computing an index of Vancouver Stock Exchange drops it to half its value after 6 months.
- 1987, the inflation in UK is computed with a rounding error: pensions are off by £100M for 21 months.
- 1991, during Gulf War 1, a Patriot system fails to intercept a Scud missile: 28 casualties.

(ロ) (同) (三) (三) (三) (○) (○)

- 1983, truncation when computing an index of Vancouver Stock Exchange drops it to half its value after 6 months.
- 1987, the inflation in UK is computed with a rounding error: pensions are off by £100M for 21 months.
- 1991, during Gulf War 1, a Patriot system fails to intercept a Scud missile: 28 casualties.

(ロ) (同) (三) (三) (三) (○) (○)

 1992, Green Party of Schleswig-Holstein seats in Parliament for a few hours, until a rounding error is discovered.

- 1983, truncation when computing an index of Vancouver Stock Exchange drops it to half its value after 6 months.
- 1987, the inflation in UK is computed with a rounding error: pensions are off by £100M for 21 months.
- 1991, during Gulf War 1, a Patriot system fails to intercept a Scud missile: 28 casualties.
- 1992, Green Party of Schleswig-Holstein seats in Parliament for a few hours, until a rounding error is discovered.
- 1995, Ariane 5 explodes during its maiden flight due to an overflow: insurance cost is \$500M.

(ロ) (同) (三) (三) (三) (○) (○)

- 1983, truncation when computing an index of Vancouver Stock Exchange drops it to half its value after 6 months.
- 1987, the inflation in UK is computed with a rounding error: pensions are off by £100M for 21 months.
- 1991, during Gulf War 1, a Patriot system fails to intercept a Scud missile: 28 casualties.
- 1992, Green Party of Schleswig-Holstein seats in Parliament for a few hours, until a rounding error is discovered.
- 1995, Ariane 5 explodes during its maiden flight due to an overflow: insurance cost is \$500M.
- ▶ 2007, Excel displays 77.1 × 850 as 100000.

 1991, during Gulf War 1, a Patriot system fails to intercept a Scud missile: 28 casualties.

Internal clock ticks every 0.1 second. Time is tracked by fixed-point arith.: $0.1 \simeq 209715 \cdot 2^{-24}$. Cumulated skew after 24h: -0.08s, distance: 160m. System was supposed to be rebooted periodically.

(ロ) (同) (三) (三) (三) (○) (○)

 1991, during Gulf War 1, a Patriot system fails to intercept a Scud missile: 28 casualties.

Internal clock ticks every 0.1 second. Time is tracked by fixed-point arith.: $0.1 \simeq 209715 \cdot 2^{-24}$. Cumulated skew after 24h: -0.08s, distance: 160m. System was supposed to be rebooted periodically.

2007, Excel displays 77.1 × 850 as 100000.

Bug in binary/decimal conversion. Failing inputs: 12 FP numbers. Probability to uncover them by random testing: 10⁻¹⁸.

Outline

Handling Machine Integers

Floating-Point Computations

Numerical Analysis

Automation

Numerical Algorithms



Binary Search

Exercise: Find appropriate precondition, postcondition, loop invariant, and variant, for this program:

```
function binary search(a:map int, n v:int): int
body
  try
    let ref 1 = 0 in
    let ref u = n - 1 in
    while 1 < u do
      let m = div (1 + u) 2 in
      if a[m] < v then</pre>
        1 := m + 1
      else if a[m] > v then
       u := m - 1
      else
        raise (Break m)
    done;
    raise Not found
  with Break i \rightarrow i
```

Target Type: int32

S2-bit signed integers in two-complement representation: integers between −2³¹ and 2³¹ − 1.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Target Type: int32

- S2-bit signed integers in two-complement representation: integers between −2³¹ and 2³¹ − 1.
- If the mathematical result of an operation fits in that range, that is the computed result.

(ロ) (同) (三) (三) (三) (○) (○)
Target Type: int32

- S2-bit signed integers in two-complement representation: integers between −2³¹ and 2³¹ − 1.
- If the mathematical result of an operation fits in that range, that is the computed result.

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

 Otherwise, an overflow occurs.
 Behavior depends on language and environment: modulo arith, saturated arith, abrupt termination, etc.

Target Type: int32

- S2-bit signed integers in two-complement representation: integers between −2³¹ and 2³¹ − 1.
- If the mathematical result of an operation fits in that range, that is the computed result.

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

 Otherwise, an overflow occurs.
 Behavior depends on language and environment: modulo arith, saturated arith, abrupt termination, etc.

A program is safe if no overflow occurs.

Idea: replace all arithmetic operations by abstract functions with preconditions. x - y becomes int32_sub(x, y).

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

function int32_sub(x: int, y: int): int requires $-2^{31} \le x - y < 2^{31}$ ensures result = x - y

Safety Checking, Try 2

Idea: replace

type int with an abstract type coercible to it,

all operations by abstract functions with preconditions, and add an axiom about the range of *int*32.

```
type int32
function of_int32(x: int32): int
axiom bounded_int32:
  forall x: int32. -2^31 ≤ of_int32(x) < 2^31</pre>
```

```
function int32_sub(x: int32, y: int32): int32
requires -2^31 ≤ of_int32(x) + of_int32(y) < 2^31
ensures of_int32(result) = of_int32(x) - of_int32(y)</pre>
```



1. How to handle int32 constants in programs?

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

2. How to specify saturating arithmetic?

Outline

Handling Machine Integers

Floating-Point Computations

Numerical Analysis

Automation

Numerical Algorithms



Floating-Point Arithmetic

- Limited range \Rightarrow exceptional behaviors.
- Limited precision \Rightarrow inaccurate results.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

IEEE-754 Binary Floating-Point Arithmetic. Width: $1 + w_e + w_m = 32$, or 64, or 128. Bias: $2^{w_e-1} - 1$. Precision: $p = w_m + 1$.

(ロ) (同) (三) (三) (三) (○) (○)

IEEE-754 Binary Floating-Point Arithmetic. Width: $1 + w_e + w_m = 32$, or 64, or 128. Bias: $2^{w_e-1} - 1$. Precision: $p = w_m + 1$.

A floating-point datum

sign *s* biased exponent $e'(w_e \text{ bits})$ mantissa $m(w_m \text{ bits})$ represents

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

IEEE-754 Binary Floating-Point Arithmetic. Width: $1 + w_e + w_m = 32$, or 64, or 128. Bias: $2^{w_e-1} - 1$. Precision: $p = w_m + 1$.

A floating-point datum

sign *s* biased exponent *e*' (w_e bits) mantissa *m* (w_m bits) represents

• if $0 < e' < 2^{w_e} - 1$, the real $(-1)^s \cdot \overline{1.m'} \cdot 2^{e'-bias}$, normal

IEEE-754 Binary Floating-Point Arithmetic. Width: $1 + w_e + w_m = 32$, or 64, or 128. Bias: $2^{w_e-1} - 1$. Precision: $p = w_m + 1$.

A floating-point datum

sign s | biased exponent e' (w_e bits) | mantissa m (w_m bits) represents

• if
$$0 < e' < 2^{w_e} - 1$$
, the real $(-1)^s \cdot \overline{1.m'} \cdot 2^{e'-bias}$, normal

- ▶ ± 0 if m' = 0, zeros • the real $(-1)^s \cdot \overline{0.m'} \cdot 2^{-bias+1}$ otherwise.
- subnormal

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

IEEE-754 Binary Floating-Point Arithmetic. Width: $1 + w_e + w_m = 32$, or 64, or 128. Bias: $2^{w_e-1} - 1$. Precision: $p = w_m + 1$.

A floating-point datum

sign *s* biased exponent $e'(w_e \text{ bits})$ mantissa $m(w_m \text{ bits})$ represents

• if
$$0 < e' < 2^{w_e} - 1$$
, the real $(-1)^s \cdot \overline{1.m'} \cdot 2^{e'-bias}$, normal

- ▶ ±0 if m' = 0, zeros ▶ the real $(-1)^s \cdot \overline{0.m'} \cdot 2^{-bias+1}$ otherwise. subnormal
- Fine real $(-1)^{n} \cdot 0.11^{n} \cdot 2^{n}$ on otherwise, subnormal

► If
$$e' = 2^{w_e} - 1$$
,
► $(-1)^s \cdot \infty$ if $m' - 0$

► $(-1)^{s} \cdot \infty$ if m' = 0, infinity ► Not-a-Number otherwise.



 $-2^{54} \times 206727 \approx -3.7 \times 10^{21}$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Semantics for the Finite Case

A floating-point operator shall behave as if it was first computing the infinitely-precise value and then rounding it so that it fits in the destination floating-point format.

(ロ) (同) (三) (三) (三) (○) (○)

Semantics for the Finite Case

A floating-point operator shall behave as if it was first computing the infinitely-precise value and then rounding it so that it fits in the destination floating-point format.

Rounding of a real number *x*:



Overflows are not considered when defining rounding: exponents are supposed to have no upper bound!

Partial Specification

Same as with integers, we specify FP operations so that no overflow occurs.

```
type bin32
function of_bin32(x: bin32): real
axiom finite_bin32: forall x: bin32. ???
function rnd...(x: real): real
axiom about_rnd...: ???
function bin32_sub(x: bin32, y: bin32): bin32
requires abs(rnd...(of_bin32(x) - of_bin32(y))) ≤ ...
ensures of_bin32(result) =
rnd(of_bin32(x) - of_bin32(y))
```

Floating-point numbers as a subset \mathbb{F} of real numbers:

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

neither infinities nor NaNs,

Floating-point numbers as a subset \mathbb{F} of real numbers:

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

- neither infinities nor NaNs,
- no signed zeros,

Floating-point numbers as a subset $\mathbb F$ of real numbers:

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

- neither infinities nor NaNs,
- no signed zeros,
- no upper bound.

Floating-point numbers as a subset \mathbb{F} of real numbers:

- neither infinities nor NaNs,
- no signed zeros,
- no upper bound.

 $\mathbb{F} = \{ m \cdot 2^e \in \mathbb{R}; \ |m| < 2^p \land e \ge e_{\min} \}$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Floating-point numbers as a subset \mathbb{F} of real numbers:

- neither infinities nor NaNs,
- no signed zeros,
- no upper bound.

 $\mathbb{F} = \{ m \cdot 2^{e} \in \mathbb{R}; \ |m| < 2^{p} \land e \ge e_{\min} \}$

Canonical representation:

▶ either $2^{p-1} \le |m| < 2^p$ and $e \ge e_{\min}$, normal ▶ or $|m| < 2^{p-1}$ and $e = e_{\min}$. subnormal

Usual Properties: Representation and Successors

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Given a representable number $x = m_x \cdot 2^{e_x} \ge 0$,

1. $y = (m_x + 1) \cdot 2^{e_x} \in \mathbb{F}$,

2. $m_x \cdot 2^{e_x}$ canonic $\Rightarrow \exists z \in \mathbb{F}, x < z < y$.

Usual Properties: Representation and Successors

Given a representable number $x = m_x \cdot 2^{e_x} \ge 0$, 1. $y = (m_x + 1) \cdot 2^{e_x} \in \mathbb{F}$, 2. $m_x \cdot 2^{e_x}$ canonic $\Rightarrow \exists z \in \mathbb{F}, x < z < y$.

Proof:

1. Hyp:
$$0 \le m_x < 2^p$$
 et $e_x \ge e_{\min}$.
If $|m_x + 1| < 2^p$, then $y = (m_x + 1) \cdot 2^{e_x} \in \mathbb{F}$.
Otherwise $m_x + 1 = 2^p$, so $y = 1 \cdot 2^{e_x + p} \in \mathbb{F}$.

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Usual Properties: Representation and Successors

Given a representable number $x = m_x \cdot 2^{e_x} \ge 0$,

1. $y = (m_x + 1) \cdot 2^{e_x} \in \mathbb{F}$,

2. $m_x \cdot 2^{e_x}$ canonic $\Rightarrow \exists z \in \mathbb{F}, x < z < y$.

Proof:

1. Hyp:
$$0 \le m_x < 2^p$$
 et $e_x \ge e_{\min}$.
If $|m_x + 1| < 2^p$, then $y = (m_x + 1) \cdot 2^{e_x} \in \mathbb{F}$.
Otherwise $m_x + 1 = 2^p$, so $y = 1 \cdot 2^{e_x + p} \in \mathbb{F}$.

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

2. Hyp:
$$2^{p-1} \le m_x < 2^p$$
 or $e_x = e_{\min}$.
If $m_x \cdot 2^{e_x} < m_z \cdot 2^{e_z} < (m_x + 1) \cdot 2^{e_x}$,
then $e_z > e_x$ and $m_z > 2^{e_x - e_z} m_x \ge 2m_x$.

Usual Properties: Rounding Modes

Faithful rounding:

- $\nabla(x) = \max\{y \in \mathbb{F} \mid y \leq x\},\$
- $\Delta(x) = \min\{y \in \mathbb{F} \mid y \ge x\},\$
- either $\operatorname{rnd}(x) = \nabla(x)$ or $\operatorname{rnd}(x) = \Delta(x)$.

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Usual Properties: Rounding Modes

Faithful rounding:

- $\nabla(x) = \max\{y \in \mathbb{F} \mid y \leq x\},\$
- $\Delta(x) = \min\{y \in \mathbb{F} \mid y \ge x\},\$
- either $\operatorname{rnd}(x) = \nabla(x)$ or $\operatorname{rnd}(x) = \Delta(x)$.

Idempotency:

 $\forall x \in \mathbb{F}, \ \operatorname{rnd}(x) = x$



Usual Properties: Rounding Modes

Faithful rounding:

- $\nabla(x) = \max\{y \in \mathbb{F} \mid y \leq x\},\$
- $\Delta(x) = \min\{y \in \mathbb{F} \mid y \ge x\},\$
- either $\operatorname{rnd}(x) = \nabla(x)$ or $\operatorname{rnd}(x) = \Delta(x)$.

Idempotency:

 $\forall x \in \mathbb{F}, \ \operatorname{rnd}(x) = x$

Local monotonicity:

 $\forall x, y \in \mathbb{R}, y \in [\operatorname{rnd}(x), x] \Rightarrow \operatorname{rnd}(y) = \operatorname{rnd}(x)$

A D F A 同 F A E F A E F A Q A

$\forall x, y \in \mathbb{R}, x \leq y \Rightarrow \operatorname{rnd}(x) \leq \operatorname{rnd}(y)$

▲□▶▲圖▶▲≣▶▲≣▶ ≣ の�?

$$\forall x, y \in \mathbb{R}, x \leq y \Rightarrow \operatorname{rnd}(x) \leq \operatorname{rnd}(y)$$

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Proof:

▶ If $\nabla(x) < \nabla(y)$, 1. $x < \nabla(y)$ by definition of $\nabla(x)$, 2. $\operatorname{rnd}(x) \leq \Delta(x) \leq \nabla(y) \leq \Box(y)$.

$$\forall x, y \in \mathbb{R}, x \leq y \Rightarrow \operatorname{rnd}(x) \leq \operatorname{rnd}(y)$$

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Proof:

▶ If $\nabla(x) < \nabla(y)$, 1. $x < \nabla(y)$ by definition of $\nabla(x)$, 2. $\operatorname{rnd}(x) \leq \Delta(x) \leq \nabla(y) \leq \Box(y)$.

• If $\nabla(x) \geq \nabla(y)$,

$$\forall x, y \in \mathbb{R}, x \leq y \Rightarrow \operatorname{rnd}(x) \leq \operatorname{rnd}(y)$$

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Proof:

▶ If $\nabla(x) < \nabla(y)$, 1. $x < \nabla(y)$ by definition of $\nabla(x)$, 2. $\operatorname{rnd}(x) \leq \Delta(x) \leq \nabla(y) \leq \Box(y)$.

▶ If
$$\nabla(x) \ge \nabla(y)$$
,
1. $\nabla(x) = \nabla(y)$ by definition of $\nabla(y)$,

$$\forall x, y \in \mathbb{R}, x \leq y \Rightarrow \operatorname{rnd}(x) \leq \operatorname{rnd}(y)$$

Proof:

▶ If $\nabla(x) < \nabla(y)$, 1. $x < \nabla(y)$ by definition of $\nabla(x)$, 2. $\operatorname{rnd}(x) \leq \Delta(x) \leq \nabla(y) \leq \Box(y)$.

• If
$$\nabla(x) \geq \nabla(y)$$
,

- 1. $\nabla(x) = \nabla(y)$ by definition of $\nabla(y)$,
- 2. $\Delta(x) = \Delta(y)$ by idempotency or successor,

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

$$\forall x, y \in \mathbb{R}, x \leq y \Rightarrow \operatorname{rnd}(x) \leq \operatorname{rnd}(y)$$

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Proof:

If $\nabla(x) < \nabla(y)$,

x < $\nabla(y)$ by definition of $\nabla(x)$,
rnd(x) ≤ ∆(x) ≤ $\nabla(y)$ ≤ □(y).

• If
$$\nabla(x) \geq \nabla(y)$$
,

- 1. $\nabla(x) = \nabla(y)$ by definition of $\nabla(y)$,
- 2. $\Delta(x) = \Delta(y)$ by idempotency or successor,
- 3. if $rnd(y) = \Delta(y)$, then $rnd(x) \le rnd(y)$,

 $\forall x, y \in \mathbb{R}, x \leq y \Rightarrow \operatorname{rnd}(x) \leq \operatorname{rnd}(y)$

Proof:

• If $\nabla(x) < \nabla(y)$,

1. $x < \nabla(y)$ by definition of $\nabla(x)$,

2. $\operatorname{rnd}(x) \leq \Delta(x) \leq \nabla(y) \leq \Box(y)$.

• If $\nabla(x) \geq \nabla(y)$,

- 1. $\nabla(x) = \nabla(y)$ by definition of $\nabla(y)$,
- 2. $\Delta(x) = \Delta(y)$ by idempotency or successor,
- 3. if $rnd(y) = \Delta(y)$, then $rnd(x) \le rnd(y)$,
- 4. otherwise rnd(x) = rnd(y) by local monotonicity.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Monotonicity:

 $\forall x, y \in \mathbb{R}, x \leq y \Rightarrow \operatorname{rnd}(x) \leq \operatorname{rnd}(y)$

Ordering with respect to representable numbers:

 $\forall x \in \mathbb{F}, \ \forall y \in \mathbb{R}, \ x \leq y \Rightarrow x \leq \operatorname{rnd}(y)$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Usual Properties: Round-Off Errors

Rounding to nearest:

For all $x \in$, there are ε and δ such that

 $\operatorname{rnd}(x) = x \cdot (1 + \varepsilon) + \delta$ and $|\varepsilon| \le 2^{-p}$ and $|\delta| \le 2^{e_{\min}-1}$

Moreover, $\delta = 0$ or $\varepsilon = 0$.


Rounding to nearest: For all $x \in$, there are ε and δ such that

 $\operatorname{rnd}(x) = x \cdot (1 + \varepsilon) + \delta$ and $|\varepsilon| \le 2^{-\rho}$ and $|\delta| \le 2^{e_{\min}-1}$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Moreover, $\delta = 0$ or $\varepsilon = 0$.

1. Hyp:
$$0 < x \notin \mathbb{F}$$
.
 $\nabla(x) = m \cdot 2^e$ and $\Delta(x) = (m+1) \cdot 2^{e+1}$.

Rounding to nearest: For all $x \in$, there are ε and δ such that

 $\operatorname{rnd}(x) = x \cdot (1 + \varepsilon) + \delta$ and $|\varepsilon| \le 2^{-\rho}$ and $|\delta| \le 2^{e_{\min}-1}$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Moreover, $\delta = 0$ or $\varepsilon = 0$.

1. Hyp:
$$0 < x \notin \mathbb{F}$$
.
 $\nabla(x) = m \cdot 2^e$ and $\Delta(x) = (m+1) \cdot 2^{e+1}$.
2. $|\operatorname{rnd}(x) - x| \leq (\Delta(x) - \nabla(x))/2 = 2^{e-1}$.

Rounding to nearest: For all $x \in$, there are ε and δ such that

 $\operatorname{rnd}(x) = x \cdot (1 + \varepsilon) + \delta$ and $|\varepsilon| \le 2^{-p}$ and $|\delta| \le 2^{e_{\min}-1}$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Moreover, $\delta = 0$ or $\varepsilon = 0$.

1. Hyp:
$$0 < x \notin \mathbb{F}$$
.
 $\nabla(x) = m \cdot 2^e$ and $\Delta(x) = (m+1) \cdot 2^{e+1}$.
2. $|\operatorname{rnd}(x) - x| \leq (\Delta(x) - \nabla(x))/2 = 2^{e-1}$.
 \blacktriangleright If $\nabla(x)$ is subnormal, $e = e_{\min}$.
 $\varepsilon = 0$ and $\delta = \operatorname{rnd}(x) - x$ so $|\delta| \leq 2^{e_{\min}-1}$

Rounding to nearest: For all $x \in$, there are ε and δ such that

 $\operatorname{rnd}(x) = x \cdot (1 + \varepsilon) + \delta$ and $|\varepsilon| \le 2^{-p}$ and $|\delta| \le 2^{e_{\min}-1}$

Moreover, $\delta = 0$ or $\varepsilon = 0$.

1. Hyp:
$$0 < x \notin \mathbb{F}$$
.
 $\nabla(x) = m \cdot 2^e$ and $\Delta(x) = (m+1) \cdot 2^{e+1}$.
2. $|\operatorname{rnd}(x) - x| \leq (\Delta(x) - \nabla(x))/2 = 2^{e-1}$.
 \bullet If $\nabla(x)$ is subnormal, $e = e_{\min}$.
 $\varepsilon = 0$ and $\delta = \operatorname{rnd}(x) - x$ so $|\delta| \leq 2^{e_{\min}-1}$.
 \bullet If $\nabla(x)$ is normal, $2^{p-1} \leq m$.
 $\delta = 0$ and $\varepsilon = (\operatorname{rnd}(x) - x)/x$ so
 $|\varepsilon| \leq 2^{e-1}/(2^{p-1} \cdot 2^e) = 2^{-p}$.

Rounding to nearest: For all $x \in$, there are ε and δ such that

 $\operatorname{rnd}(x) = x \cdot (1 + \varepsilon) + \delta$ and $|\varepsilon| \le 2^{-p}$ and $|\delta| \le 2^{e_{\min}-1}$

Moreover, $\delta = 0$ or $\varepsilon = 0$.



Rounding to nearest: For all $x \in$, there are ε and δ such that

 $\operatorname{rnd}(x) = x \cdot (1 + \varepsilon) + \delta$ and $|\varepsilon| \le 2^{-p}$ and $|\delta| \le 2^{e_{\min}-1}$

Moreover, $\delta = 0$ or $\varepsilon = 0$.

Directed rounding: For all $x \in$, there are ε and δ such that

 $\operatorname{rnd}(x) = x \cdot (1 + \varepsilon) + \delta$ and $|\varepsilon| \le 2^{-p+1}$ and $|\delta| \le 2^{e_{\min}}$

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

Moreover, $\delta = 0$ or $\varepsilon = 0$.

Usual Properties: Subnormal Addition

Sums in the subnormal range are representable:

 $\forall x, y \in \mathbb{F}, |x + y| \le 2^{e_{\min} + \rho} \Rightarrow x + y \in \mathbb{F}$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

Usual Properties: Subnormal Addition

Sums in the subnormal range are representable:

$$orall x,y\in \mathbb{F},\;|x+y|\leq 2^{e_{\mathsf{min}}+
ho} \Rightarrow x+y\in \mathbb{F}$$

Proof:

1.
$$x = m_x \cdot 2^{e_x}$$
 and $y = m_y \cdot 2^{e_y}$.
2. $m = m_x \cdot 2^{e_x - e_{\min}} + m_y \cdot 2^{e_y - e_{\min}}$ and $x + y = m \cdot 2^{e_{\min}}$.
3. $|m| \le 2^p$ so $x + y \in \mathbb{F}$.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

Usual Properties: Subnormal Addition

Sums in the subnormal range are representable:

 $\forall x, y \in \mathbb{F}, |x + y| \le 2^{e_{\min} + \rho} \Rightarrow x + y \in \mathbb{F}$

Proof:

1.
$$x = m_x \cdot 2^{e_x}$$
 and $y = m_y \cdot 2^{e_y}$.
2. $m = m_x \cdot 2^{e_x - e_{\min}} + m_y \cdot 2^{e_y - e_{\min}}$ and $x + y = m \cdot 2^{e_{\min}}$.
3. $|m| \le 2^p$ so $x + y \in \mathbb{F}$.

Round-off error for addition:

 $\forall x, y \in \mathbb{F}, \exists \varepsilon, \circ (x + y) = (x + y) \cdot (1 + \varepsilon) \text{ and } |\varepsilon| \le 2^{-p}$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Outline

Handling Machine Integers

Floating-Point Computations

Numerical Analysis

Automation

Numerical Algorithms

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ● ● ●

Numerical Errors

Given two real numbers u and v,

• absolute error: u - v, |u - v|

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

• relative error: $u/v - 1, \ldots$

Numerical Errors

Given two real numbers u and v,

- absolute error: u v, |u v|
- relative error: $u/v 1, \ldots$

Combining errors:

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Numerical Errors

Given two real numbers u and v,

- absolute error: u v, |u v|
- relative error: $u/v 1, \ldots$

Combining errors:

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Remark: rnd(u) - v = (rnd(u) - u) + (u - v)

Notations:

- a mathematical function f(x),
- a floating-point program $\tilde{f}(x)$,
- the infinitely-precise evaluation $\hat{f}(x)$ of $\tilde{f}(x)$.

▲□▶ ▲□▶ ▲三▶ ▲三▶ - 三 - のへで

Notations:

- a mathematical function f(x),
- a floating-point program $\tilde{f}(x)$,
- the infinitely-precise evaluation $\hat{f}(x)$ of $\tilde{f}(x)$.

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

Definitions:

- forward error: $\tilde{f}(x) f(x)$,
 - round-off error: $\tilde{f}(x) \hat{f}(x)$
 - method error: $\hat{f}(x) f(x)$

Notations:

- a mathematical function f(x),
- a floating-point program $\tilde{f}(x)$,
- the infinitely-precise evaluation $\hat{f}(x)$ of $\tilde{f}(x)$.

Definitions:

- forward error: $\tilde{f}(x) f(x)$,
 - round-off error: $\tilde{f}(x) \hat{f}(x)$
 - method error: $\hat{f}(x) f(x)$
- backward error: $\tilde{x} x$

with \tilde{x} closest from x such that $f(\tilde{x}) = \tilde{f}(x)$

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

Notations:

- a mathematical function f(x),
- a floating-point program $\tilde{f}(x)$,
- the infinitely-precise evaluation $\hat{f}(x)$ of $\tilde{f}(x)$.

Definitions:

- forward error: $\tilde{f}(x) f(x)$,
 - round-off error: $\tilde{f}(x) \hat{f}(x)$
 - method error: $\hat{f}(x) f(x)$
- backward error: $\tilde{x} x$

with \tilde{x} closest from x such that $f(\tilde{x}) = \tilde{f}(x)$

Remark: $\tilde{f}(x) - f(x) \simeq (\tilde{x} - x) \times \frac{\partial f}{\partial x}$. In other words: forward err \simeq backward err \times condition num.

Evaluating $\sum_i a_i \cdot x^i$:

```
function Horner
  (a:map binary32, n:int, x:binary32)
body
  let ref y := binary32 cst(0.) in
```

```
let ref i := n in
for i = 0 to n - 1 do
    y := binary32_add(binary32_mul(y, x), a[i]);
done;
```

У

Outline

Handling Machine Integers

Floating-Point Computations

Numerical Analysis

Automation

Numerical Algorithms



Using Ghost Variables for Model Values

```
function det (a b c d: binary32, aM bM cM dM: real):
  (binary32, real)
body
  let t1 := binary32_mul(a, d) in
  let \pm 1M := aM * dM in
  let t2 := binary32 mul(b, c) in
  let t2M := bM * cM in
  let t3 := binary32 sub(t1, t2) in
  let \pm 3M := \pm 1M - \pm 2M in
  (t3, t3M)
```

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Forward error: property about t3 - t3M or t3/t3M - 1.

Implicit Model Values

```
function of_bin32(x: binary32): real
function model_of(x: binary32): real
```

```
function binary32_add(x y: binary32): binary32
requires
   abs(rnd...(of_bin32(x) + of_bin32(y))) ≤
    max_binary32
ensures
   of_bin32(result) =
    rnd(of_bin32(x) + of_bin32(y)) ∧
   model of(result) = model of(x) + model of(y)
```

Abstract Interpretation

Domains for floating-point variables:

- for the computed value x,
- for the infinitely-precise value \hat{x} ,

(ロ) (同) (三) (三) (三) (○) (○)

• for the absolute error $x - \hat{x}$,

▶

Abstract Interpretation

Domains for floating-point variables:

- for the computed value x,
- for the infinitely-precise value \hat{x} ,
- for the absolute error $x \hat{x}$,

▶ ...

Naive domains:

- ▶ $[\underline{x}, \overline{x}]$ such that $x \in [\underline{x}, \overline{x}]$, ex: $rnd(x + y) \in [rnd(\underline{x} + \underline{y}), rnd(\overline{x} + \overline{y})]$,
- no domain for \hat{x} ,

•
$$\delta_x$$
 such that $|x - \hat{x}| \le \delta_x$,
ex: $\delta_{x+y} = \delta_x + \delta_y + 2^{-p} \max(\overline{x} + \overline{y}, -(\underline{x} + \underline{y}))$

Outline

Handling Machine Integers

Floating-Point Computations

Numerical Analysis

Automation

Numerical Algorithms

◆□ ▶ ◆□ ▶ ◆三 ▶ ◆□ ▶ ◆□ ●

Newton's Iterated Square Root

```
function fp_sqrt_init(x:binary64) : binary64
  requires 0.5 < x < 2;
  ensures abs(result - 1/sqrt(x)) \leq 2^{-6} + 1/sqrt(x);
function fp sqrt(x:binary64) : binary64
  requires 0.5 < x < 2;
  ensures abs (result - sqrt(x)) \leq 2^{-43} * sqrt(x);
body
  let ref t := fp sqrt init(x) in
  for i = 1 \pm 0.3 do
   t := 0.5 * t * (3 - t * t * x);
  done;
  t. * x
```

Quadratic Convergence

For all *u* and *x*:

$$0.5u(3-u^2x)\sqrt{x}-1=-(1.5+0.5(u\sqrt{x}-1))\times(u\sqrt{x}-1)^2$$

Loop iterations:

$$t_{n+1}\sqrt{x} - 1 \simeq 0.5t_n(3 - t_n^2 x)\sqrt{x} - 1 \simeq -1.5(t_n\sqrt{x} - 1)^2$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Round-off error at step *n* vanishes at step n + 1.

Accurate Summation

```
Computing \sum_i x_i:
s := x[0];
e := 0.;
for i = 1 to n - 1 do
  y := x[i];
  t := s + y;
  u := t - y;
  r := (s - u) + (y - (t - u));
  s := t;
  e := e + r;
done;
s' := s + e;
```

◆□> ◆□> ◆豆> ◆豆> ・豆・ のへぐ

Accurate Summation

```
Computing \sum_i x_i:
s := x[0];
e := 0.;
for i = 1 to n - 1 do
  y := x[i];
  t := s + y;
  u := t - y;
  r := (s - u) + (y - (t - u));
  s := t;
  e := e + r;
done;
s' := s + e;
```

Naive sum

Accurate Summation

```
Computing \sum_i x_i:
s := x[0];
e := 0.;
for i = 1 to n - 1 do
  y := x[i];
  t := s + y;
  u := t - y;
  r := (s - u) + (y - (t - u));
  s := t;
  e := e + r;
done;
s' := s + e;
```

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Error-free addition: t + r = s + y

Sterbenz: $\forall x, y \in F, x/2 \le y \le 2x \Rightarrow \operatorname{rnd}(x - y) = x - y$

Sterbenz: $\forall x, y \in F, x/2 \le y \le 2x \Rightarrow \operatorname{rnd}(x - y) = x - y$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

• error of addition: $\forall x, y \in F$, $rnd(x + y) - (x + y) \in F$

Sterbenz: $\forall x, y \in F, x/2 \le y \le 2x \Rightarrow \operatorname{rnd}(x - y) = x - y$

A D F A 同 F A E F A E F A Q A

- error of addition: $\forall x, y \in F$, $rnd(x + y) (x + y) \in F$
- ▶ fast twosum: $\forall x, y \in F$, $|x| \ge |y| \Rightarrow s + e = x + y$ with $s = \operatorname{rnd}(x + y)$ and $e = \operatorname{rnd}(y - \operatorname{rnd}(s - x))$

Sterbenz: $\forall x, y \in F, x/2 \le y \le 2x \Rightarrow \operatorname{rnd}(x - y) = x - y$

- error of addition: $\forall x, y \in F$, $rnd(x + y) (x + y) \in F$
- ▶ fast twosum: $\forall x, y \in F$, $|x| \ge |y| \Rightarrow s + e = x + y$ with $s = \operatorname{rnd}(x + y)$ and $e = \operatorname{rnd}(y - \operatorname{rnd}(s - x))$
- ▶ twosum: $\forall x, y \in F$, s + e = x + ywith $s = \operatorname{rnd}(x + y)$ and $u = \operatorname{rnd}(s - y)$ and $e = \operatorname{rnd}(\operatorname{rnd}(x - u) + \operatorname{rnd}(y - \operatorname{rnd}(s - u)))$

Payne & Hanek's Argument Reduction

```
Reducing x \ge 2^{31} to 0 \le y \le \pi/4 for circular functions:
```

```
function reduce(x:binary32): (binary32, int)
requires 2^31 ≤ x
ensures exists 1:int.
   abs((result + k * pi/4) - (x + 1 * 2*pi)) ≤ 2^-25
body
let x' = binary64_of_binary32 x in
let t = x' * 1.273239545... in
let k = trunc(t) in
let y = (t - k) * 0.785398163... in
(binary32_of_binary64(y), k)
```

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Note: computations are performed with binary64.

Payne & Hanek's Argument Reduction

```
Reducing x \ge 2^{31} to 0 \le y \le \pi/4 for circular functions:
```

```
function reduce(x:binary32): (binary32, int)
requires 2^31 ≤ x
ensures exists 1:int.
    abs((result + k * pi/4) - (x + 1 * 2*pi)) ≤ 2^-25
body
let x' = binary64_of_binary32 x in
let t = x' * 0.02323954474... in
let k = trunc(t) in
let y = (t - k) * 0.785398163... in
(binary32_of_binary64(y), k)
```

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Note: computations are performed with binary64.