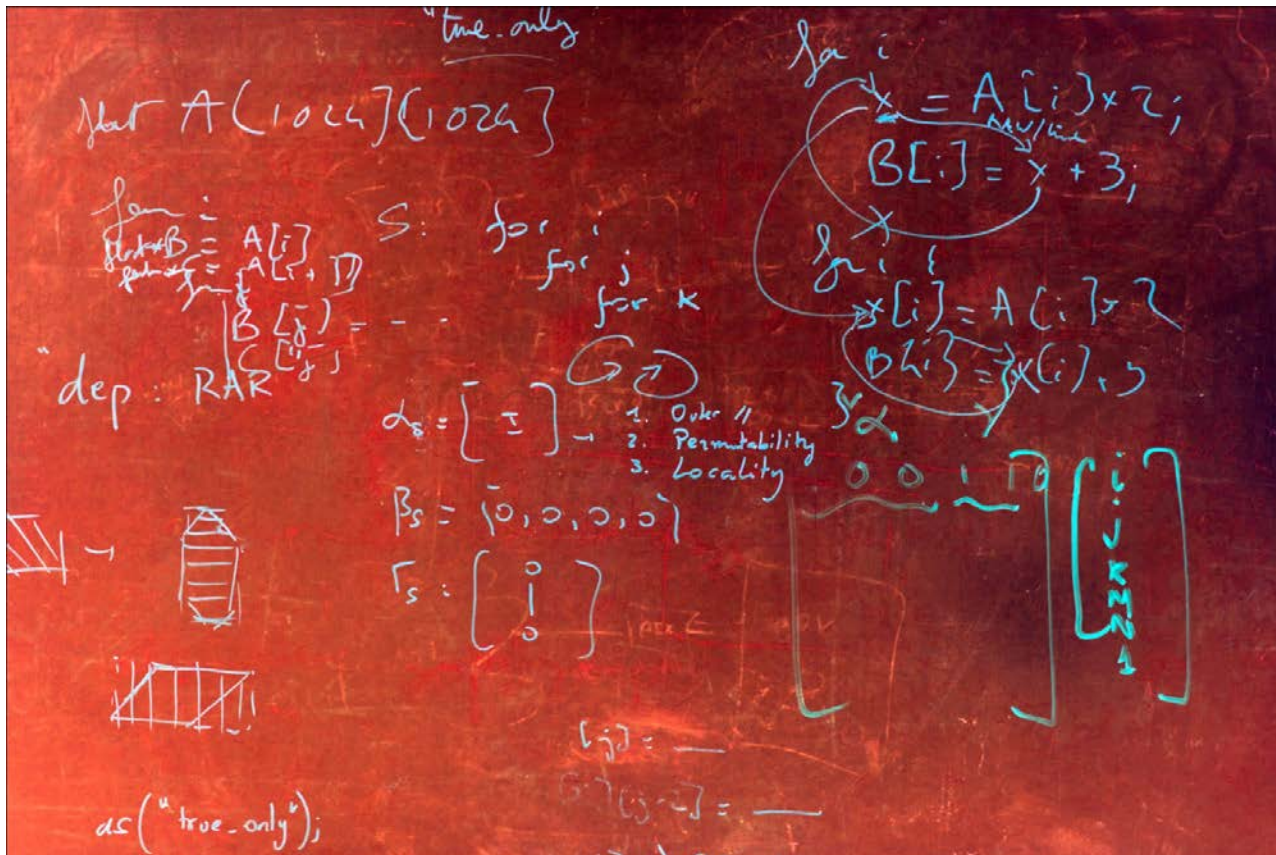


A Unified Coq Framework for Verifying C Programs with Floating-Point Computations

Tahina Ramananandro, Paul Mountcastle,
Benoît Meister, Richard Lethin



The High-Level Problem

Goal: energy-efficient implementations of radar algorithms

- Naïve implementations consume time and energy
- Ideas: compute in lower-precision floating-point and/or with further approximations
- Approximations and their floating-point implementations introduce some error in the result
- How to compute some implementation **error bound**?
- How can we **trust** this error bound?

Our Achievements

VCFloat: a Coq library for handling floating-point computations in the verification of C programs

- Automatically compute real-number expressions with rounding error terms and their correctness proofs

Use case: SAR backprojection with linear interpolation

- Introduce approximations for square root and sine
- Tune between single- and double-precision floating-points
- Compute error bounds wrt. "ideal" mathematical real-number algorithm
- **Formal proof** of correctness using the Coq proof assistant
- Energy measurements: ~10-20% saved on Intel Haswell

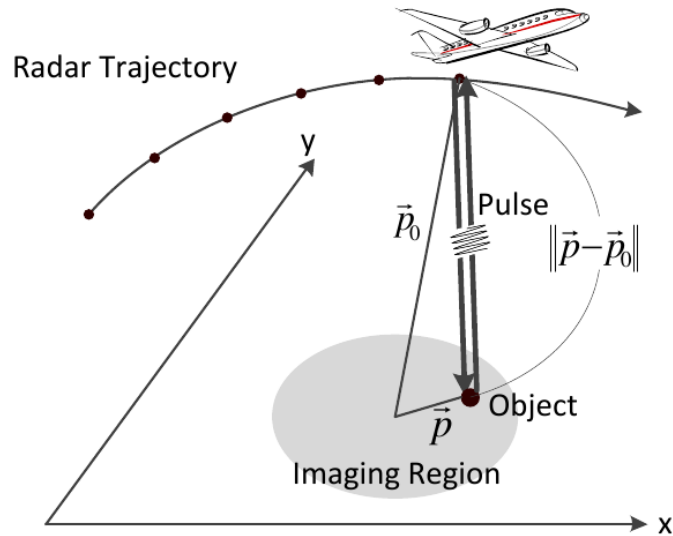
This Presentation

- Certified error bounds for energy-efficient radar image processing
- Our Coq framework: VCFloat
- Demo
- Conclusions

Certified Error Bounds for

RADAR IMAGE PROCESSING

Synthetic Aperture Radar (SAR) Backprojection

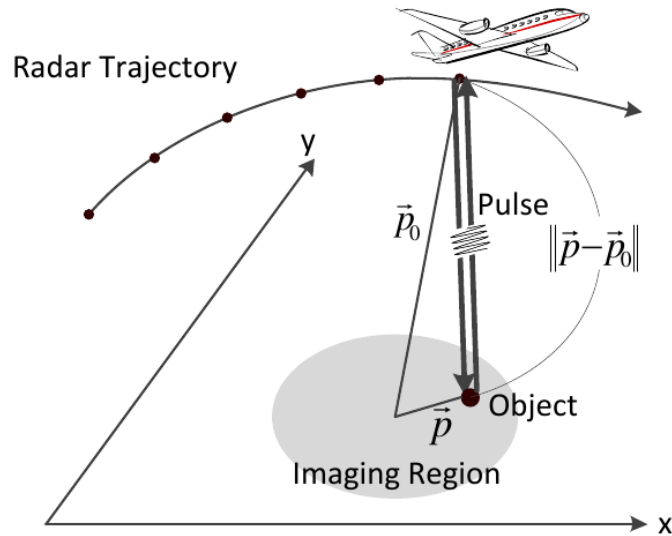


```
for all pixels  $x$  do
  for all pulses  $p$  do
     $R$  = distance between platform and pixel  $x$  for pulse  $p$ 
     $s$  = sample interpolated from pulse  $p$  in neighborhood of range  $R$ 
    Apply phase correction to sample  $s$  based on range  $R$ 
    Accumulate sample  $s$  into pixel  $x$ 
  end for
end for
```

Real-number algorithm

Figure from Park et al. *Efficient Backprojection-Based Synthetic Aperture Radar Computation with Many-Core Processors*, SC 2012

SAR Backprojection



```

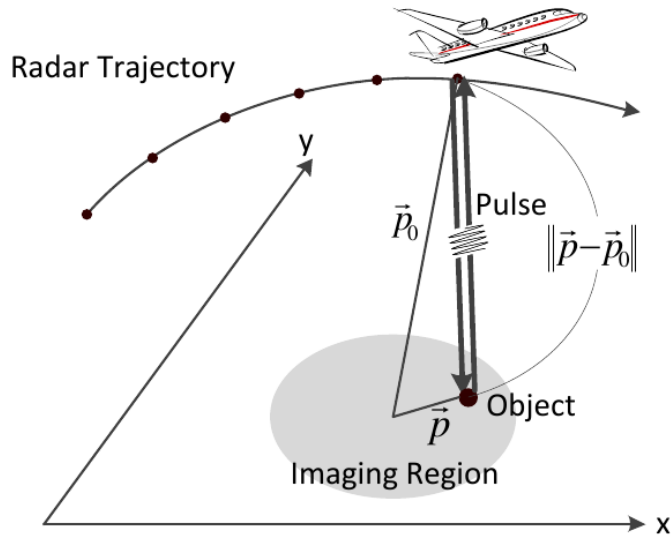
for  $y := 0$  to  $BP\_NPIX\_Y - 1$  do
   $py := (y + \frac{1-BP\_NPIX\_Y}{2}) \times dx dy$ 
  for  $x := 0$  to  $BP\_NPIX\_X - 1$  do
     $px := (x + \frac{1-BP\_NPIX\_X}{2}) \times dx dy$ 
     $\underline{image}[y][x] := 0 \in \mathbb{C}$ 
    for  $p := 0$  to  $N\_PULSES - 1$  do
       $r := \|\underline{platpos}[p] - (px, py, z[p][y][x])\|$ 
       $\underline{bin} := (r - r_0) / dr$ 
       $\underline{sample} := \underline{binSample}(N\_RANGE\_UPSAMPLED, \underline{data}[p], \underline{bin})$ 
       $\underline{matchedFilter} := \exp(2i \times ku \times r)$ 
       $\underline{image}[y][x] := \underline{image}[y][x] + \underline{sample} \times \underline{matchedFilter}$ 
    end for
  end for
end for
return  $\underline{image}$ 

```

Real-number algorithm

Figure from Park et al. *Efficient Backprojection-Based Synthetic Aperture Radar Computation with Many-Core Processors*, SC 2012

SAR Backprojection



```

for y := 0 to BP_NPIX_Y - 1 do
  py := (y +  $\frac{1 - BP\_NPIX\_Y}{2}$ ) × dx dy
  for x := 0 to BP_NPIX_X - 1 do
    px := (x +  $\frac{1 - BP\_NPIX\_X}{2}$ ) × dx dy
    image[y][x] := 0 ∈ ℂ
    for p := 0 to N_PULSES - 1 do
      r :=  $\|\text{platpos}[p] - (px, py, z[p][y][x])\|$ 
      bin := (r - r0) / dr
      sample := binSample(N_RANGE_UPSAMPLED, data[p], bin)
      matchedFilter :=  $\exp(2i \times ku \times r)$ 
      image[y][x] := image[y][x] + sample × matchedFilter
    end for
  end for
end for
return image

```

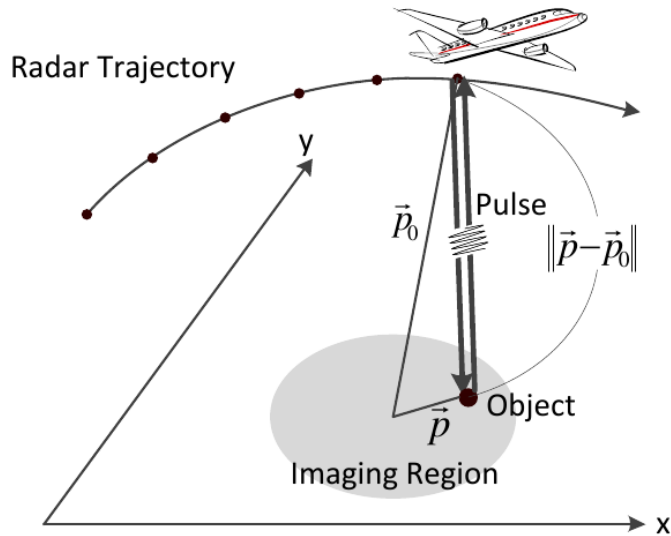
square root

sine

Real-number algorithm

Figure from Park et al. *Efficient Backprojection-Based Synthetic Aperture Radar Computation with Many-Core Processors*, SC 2012

SAR Backprojection



```

for y := 0 to BP_NPIX_Y - 1 do
  py := (y +  $\frac{1-BP\_NPIX\_Y}{2}$ ) × dx dy
  for x := 0 to BP_NPIX_X - 1 do
    px := (x +  $\frac{1-BP\_NPIX\_X}{2}$ ) × dx dy
    image[y][x] := 0 ∈ ℂ
    for p := 0 to N_PULSES - 1 do
      r := ||platpos[p] - (px, py, z[p][y][x])||
      bin := (r - r0)/dr
      sample := binSample(N_RANGE_UPSAMPLED, data[p], bin)
      matchedFilter := exp(2i × ku × r)
      image[y][x] := image[y][x] + sample × matchedFilter
    end for
  end for
end for
return image

```

square root

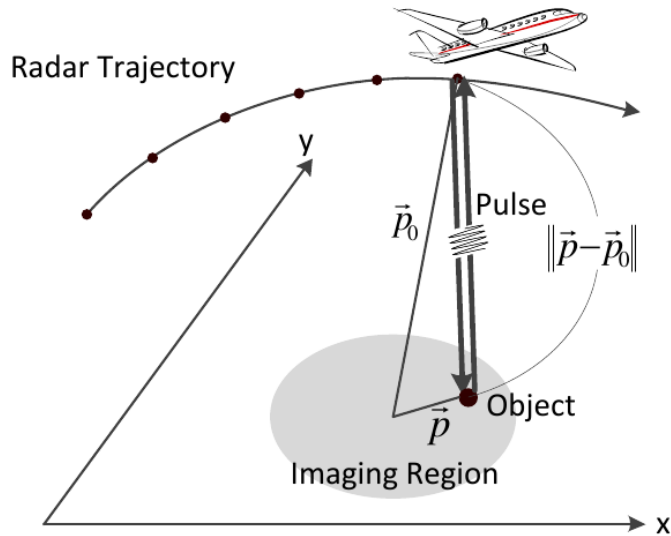
sine

floating-point computations

Implementation

Figure from Park et al. *Efficient Backprojection-Based Synthetic Aperture Radar Computation with Many-Core Processors*, SC 2012

SAR Backprojection



```
void backprojection
(int const BP_NPIX_X, int const BP_NPIX_Y,
 int const N_PULSES, ...,
 float const **data_r, float const **data_i,
 double** image_r, double** image_i, ...) {
  for (int y = 0; y < BP_NPIX_Y, ++y) {
    ...
    for (int x = 0; x < BP_NPIX_X, ++x) {
      ...
      for (int p = 0; p < N_PULSES, ++p) {
        ...
        double ... = ... sqrt(...) ... ;
        ...
        double ... = ... sin(...) ... ;
        ...
      }
    }
  }
}
```

square root

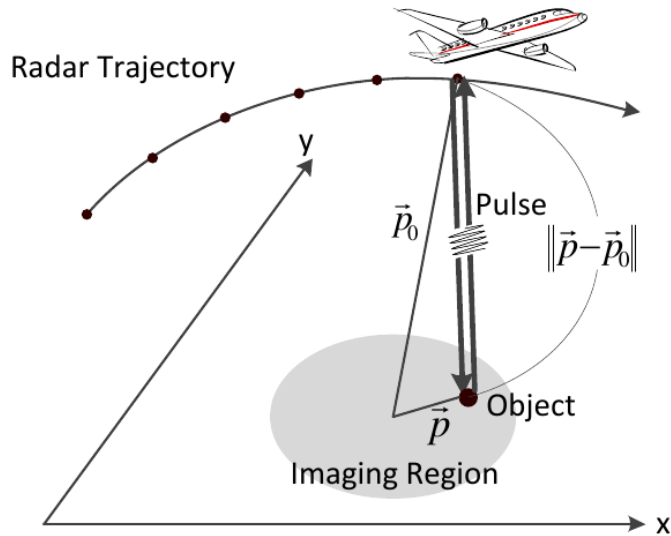
sine

floating-point computations

C Implementation

Figure from Park et al. *Efficient Backprojection-Based Synthetic Aperture Radar Computation with Many-Core Processors*, SC 2012

SAR Backprojection



```
void backprojection
(int const BP_NPIX_X, int const BP_NPIX_Y,
 int const N_PULSES, ...,
 float const **data_r, float const **data_i,
 double** image_r, double** image_i, ...) {
  for (int y = 0; y < BP_NPIX_Y, ++y) {
```

```
    ...
    for (int x = 0; x < BP_NPIX_X, ++x) {
      ...
      for (int p = 0; p < N_PULSES, ++p) {
        ...
        double ... = ... approx_sqrt(...) ... ;
        ...
        double ... = ... approx_sin(...) ... ;
        ...
      }
    }
  }
```

approximate
square root

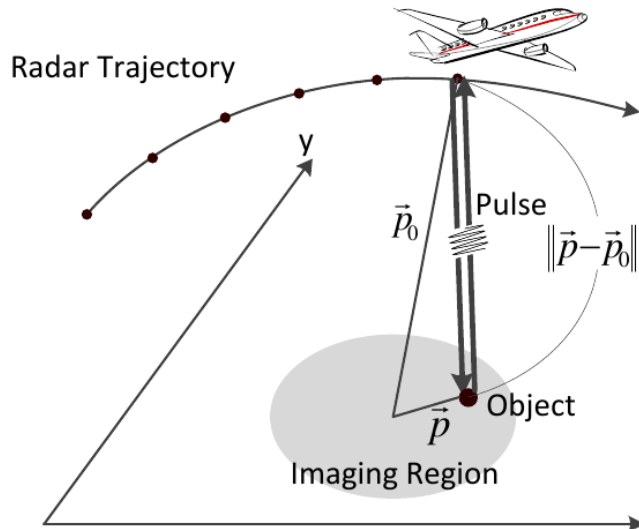
approximate
sine

floating-point computations

C Implementation

Figure from Park et al. *Efficient Backprojection-Based Synthetic Aperture Radar Computation with Many-Core Processors*, SC 2012

SAR Backprojection



```
void backprojection
(int const BP_NPIX_X, int const BP_NPIX_Y,
 int const N_PULSES, ...,
 float const **data_r, float const **data_i,
 float** image_r, float** image_i, ...) {
  for (int y = 0; y < BP_NPIX_Y, ++y) {
    ...
    for (int x = 0; x < BP_NPIX_X, ++x) {
      ...
      for (int p = 0; p < N_PULSES, ++p) {
        ...
        double ... = ... approx_sqrt(...) ... ;
        float ... = ... approx_sin(...) ... ;
        ...
      }
    }
  }
}
```

Precision tuning

floating-point computations

approximate square root

approximate sine

C Implementation
(~ 150 lines)

Figure from Park et al. *Efficient Backprojection-Based Synthetic Aperture Radar Computation with Many-Core Processors*, SC 2012

Image Error Analysis

Maximize Signal-Noise Ratio:

$$SNR := \frac{\|image_0\|^2}{\|image - image_0\|^2}$$

Find an upper bound on the denominator

- Absolute error bound is enough

Error sources:

- Method errors introduced by approximation
- Rounding errors introduced by floating-point computations

Final Correctness Statement (slightly simplified)

$\forall P \text{ ` (HYPS: SARHypotheses P) } m$

$(Hm: \text{holds } m (P \ ++$

$\text{Pperm_int bir oir (BP_NPIX_X } \times \text{ BP_NPIX_Y) } \ ++$

$\text{Pperm_int bii oii (BP_NPIX_X } \times \text{ BP_NPIX_Y)),$

Hypotheses

$\exists m' , \text{star Clight.step2}$

$(\text{Callstate fn_sar_backprojection ...}) (\text{Returnstate Vundef Kstop } m') \wedge$

$\exists \text{image_r image_i,}$

$\text{holds } m (P \ ++$

$\text{Parray_int image_r bir oir (BP_NPIX_X } \times \text{ BP_NPIX_Y) } \ ++$

$\text{Parray_int image_i bii oii (BP_NPIX_X } \times \text{ BP_NPIX_Y)) } \wedge$

$\forall y, (y < \text{BP_NPIX_Y}) \% \text{nat} \rightarrow \forall x, (x < \text{BP_NPIX_X}) \% \text{nat} \rightarrow$

$\text{let } ir := \text{image_r } (y \times \text{BP_NPIX_X} + x) \text{ in}$

$\text{let } ii := \text{image_i } (y \times \text{BP_NPIX_X} + x) \text{ in}$

$\text{is_finite } _ _ ir = \text{true} \wedge \text{is_finite } _ _ ii = \text{true} \wedge$

$\text{let } (tr, ti) := \text{SARBackProj.sar_backprojection } y \ x \text{ in}$

$\text{Rabs (B2R } _ _ ir - tr) \leq \text{pixel_bound} \wedge$

$\text{Rabs (B2R } _ _ ii - ti) \leq \text{pixel_bound.}$

Conclusions

Final Correctness Statement (slightly simplified)

$\forall P$ (HYPS: SARHypotheses P) m | Hypotheses on input data

(Hm: holds m (P ++

Pperm_int bir oir (BP_NPIX_X \times BP_NPIX_Y) ++

Pperm_int bii oii (BP_NPIX_X \times BP_NPIX_Y)),

Memory contents
and permissions

$\exists m'$, star Clight.step2

(Callstate fn_sar_backprojection ...) (Returnstate Vundef Kstop m') \wedge

\exists image_r image_i,

holds m (P ++

Parray_int image_r bir oir (BP_NPIX_X \times BP_NPIX_Y) ++

Parray_int image_i bii oii (BP_NPIX_X \times BP_NPIX_Y)) \wedge

$\forall y, (y < BP_NPIX_Y)\%nat \rightarrow \forall x, (x < BP_NPIX_X)\%nat \rightarrow$

let ir := image_r (y \times BP_NPIX_X + x) in

let ii := image_i (y \times BP_NPIX_X + x) in

is_finite __ ir = true \wedge is_finite __ ii = true \wedge

let (tr, ti) := SARBackProj.sar_backprojection y x in

Rabs (B2R __ ir - tr) \leq pixel_bound \wedge

Rabs (B2R __ ii - ti) \leq pixel_bound.

Conclusions

Final Correctness Statement (slightly simplified)

$\forall P$ (HYPS: SARHypotheses P) m

(Hm: holds m (P ++

Pperm_int bir oir (BP_NPIX_X × BP_NPIX_Y) ++

Pperm_int bii oii (BP_NPIX_X × BP_NPIX_Y)),

Hypotheses

$\exists m'$, star Clight.step2

(Callstate fn_sar_backprojection ...) (Returnstate Vundef Kstop m') \wedge

C code runs

\exists image_r image_i,

holds m (P ++

Parray_int image_r bir oir (BP_NPIX_X × BP_NPIX_Y) ++

Parray_int image_i bii oii (BP_NPIX_X × BP_NPIX_Y) \wedge

$\forall y, (y < BP_NPIX_Y)\%nat \rightarrow \forall x, (x < BP_NPIX_X)\%nat \rightarrow$

let ir := image_r (y × BP_NPIX_X + x) in

let ii := image_i (y × BP_NPIX_X + x) in

is_finite __ ir = true \wedge is_finite __ ii = true \wedge FP does not overflow

let (tr, ti) := SARBackProj.sar_backprojection y x in

Rabs (B2R __ ir - tr) \leq pixel_bound \wedge

Rabs (B2R __ ii - ti) \leq pixel_bound.

Total implementation error bound
(approximation + rounding)

computed at proof-building time

Memory
contents

Polynomial Approximations of Sine

Built-in hardware sine is costly in energy and time

- Replace core sine with a polynomial approximation
 - Use convex optimization
 - Compute coefficients with **unverified** numerical tools
 - Do not trust the results, use Coq to prove an error bound
- Naïve argument reduction is enough for SAR
 - Errors due to approximation of π and roundings
 - Lower than implementation error for core computation

Adaptive Approximate Square Root

- Replace square root with 2-degree Taylor polynomial
 - Taylor-Lagrange inequality bounds method error
- Valid only in a convergence disc
 - Outside, use accurate hardware square root
 - **Adaptive algorithm:** Re-center the disc as needed

Precision Results

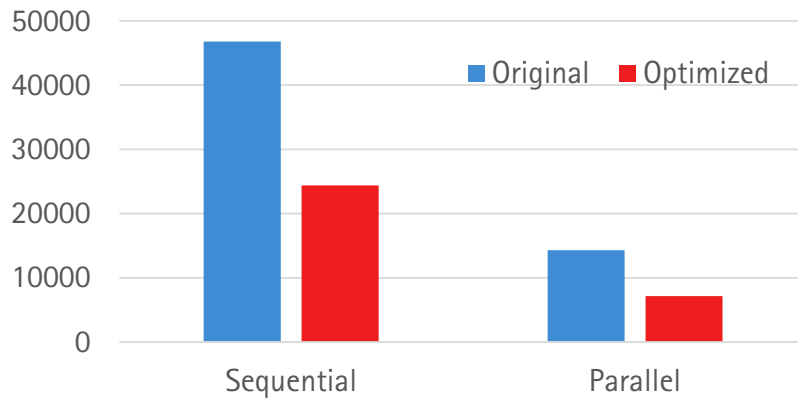
- Input data bounds from DARPA PERFECT suite
- Error grows with image size
- No statistical reasoning about data

Norm	Interpol.	Sine	Final sum	Small dB	Med. dB	Large dB
Double	Double	Double	Doub/Sing	-2	4	10
Double	Single	Doub/Sing	Doub/Sing	3	15	27
Adaptive	Double	Double	Doub/Sing	12	24	37
Adaptive	Single	Doub/Sing	Doub/Sing	14	26	39
Double	Single	Approx.	Doub/Sing	17	29	40
Adaptive	Single	Approx.	Doub/Sing	20	32	44
Double	Single	Doub/Sing	Single	34	52	70
Adaptive	Single	Approx.	Single	35	53	70

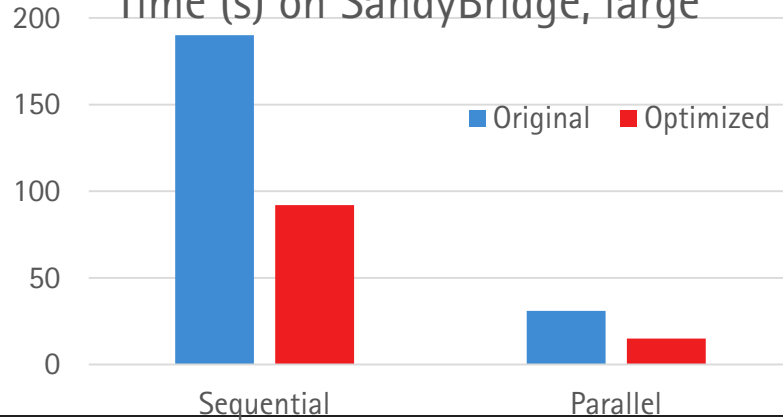
Performance Measurements for Optimized C Code

- Intel SandyBridge: direct energy measurements
- Intel Haswell: energy model unknown, time instead

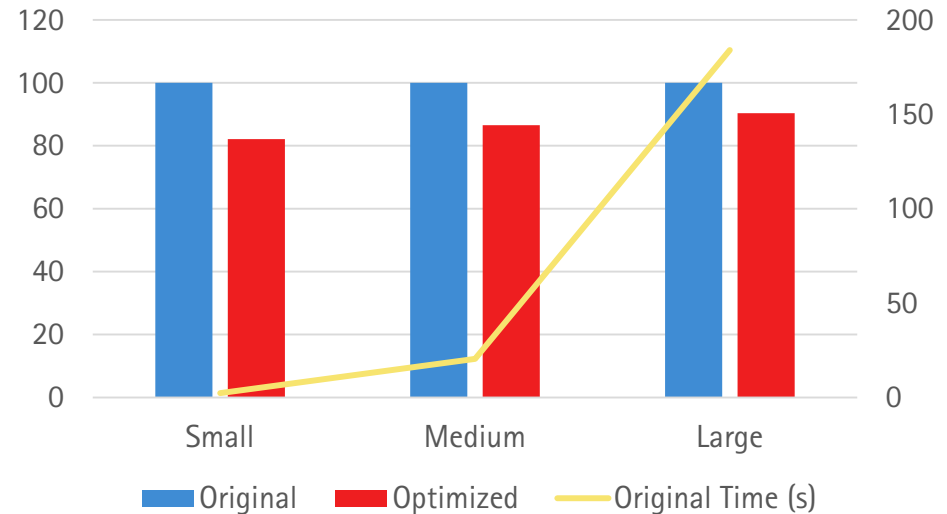
Energy (J) on SandyBridge, large



Time (s) on SandyBridge, large



Time % on Haswell, parallel



SAR proof: facts and figures

C code size: 150 lines

Proof size:

- Previous all-manual proof: 26k lines, no connection with C
- Thanks to VCFloat: reduced to 12k lines
 - 5k lines of spec (loop invariants), 7k lines of proof
 - ~2k lines of proof for real-number reasoning
 - Remaining part due to C language constructs, could be further reduced when integrating with Verifiable C (Appel et al. 2014)

Proof building/checking time:

- 1 hour (4-core Intel Core i7, 2.10 GHz, 4 Gb RAM)
- Mostly due to interval computations

Formal Verification of Floating-Point Computations in C Programs

OUR COQ FRAMEWORK: VCFLOAT

Final Correctness Statement (slightly simplified)

$\forall P$ (HYPS: SARHypotheses P) m

(Hm: holds m (P ++

Pperm_int bir oir (BP_NPIX_X × BP_NPIX_Y) ++

Pperm_int bii oii (BP_NPIX_X × BP_NPIX_Y)),

Hypotheses

$\exists m'$, star Clight.step2

(Callstate fn_sar_backprojection ...) (Returnstate Vundef Kstop m') \wedge

C code runs

\exists image_r image_i,

holds m (P ++

Parray_int image_r bir oir (BP_NPIX_X × BP_NPIX_Y) ++

Parray_int image_i bii oii (BP_NPIX_X × BP_NPIX_Y) \wedge

$\forall y, (y < BP_NPIX_Y)\%nat \rightarrow \forall x, (x < BP_NPIX_X)\%nat \rightarrow$

let ir := image_r (y × BP_NPIX_X + x) in

let ii := image_i (y × BP_NPIX_X + x) in

is_finite __ ir = true \wedge is_finite __ ii = true \wedge FP does not overflow

let (tr, ti) := SARBackProj.sar_backprojection y x in

Rabs (B2R __ ir - tr) \leq pixel_bound \wedge

Rabs (B2R __ ii - ti) \leq pixel_bound.

Total implementation error bound
(approximation + rounding)

computed at proof-building time



Final Correctness Statement (slightly simplified)

$\forall P \text{ ` (HYPS: SARHypotheses P) } m$

$(Hm: \text{holds } m (P \text{ ++}$

$\text{Pperm_int bir oir (BP_NPIX_X } \times \text{ BP_NPIX_Y) ++}$

$\text{Pperm_int bii oii (BP_NPIX_X } \times \text{ BP_NPIX_Y)),$

Hypotheses

$\exists m', \text{star Clight.step2}$

$(\text{Callstate fn_sar_backprojection ...}) (\text{Returnstate Vundef Kstop } m')$

C code runs

$\exists \text{image_r image_i,}$

$\text{holds } m (P \text{ ++}$

CompCert Clight

Memory contents

$\text{Parray_int image_r bir oir (BP_NPIX_X } \times \text{ BP_NPIX_Y) ++}$

$\text{Parray_int image_i bii oii (BP_NPIX_X } \times \text{ BP_NPIX_Y)) } \wedge$

$\forall y, (y < \text{BP_NPIX_Y}) \% \text{nat} \rightarrow \forall x, (x < \text{BP_NPIX_X}) \% \text{nat} \rightarrow$

$\text{let ir := image_r (y } \times \text{ BP_NPIX_X + x)) in}$

$\text{let ii := image_i (y } \times \text{ BP_NPIX_X + x)) in}$

$\text{is_finite __ ir = true } \wedge \text{ is_finite __ ii = true} \wedge$ FP does not overflow

$\text{let (tr, ti) := SARBackProj.sar_backprojection y x in}$

$\text{Rabs (B2R __ ir - tr) } \leq \text{pixel_bound } \wedge$

$\text{Rabs (B2R __ ii - ti) } \leq \text{pixel_bound.}$

Total implementation error bound (approximation + rounding) computed at proof-building time



Flocq

Our Design Choices: Which Formal Methods?

Verification using the Coq proof assistant

Correctness fully embedded in Coq using existing libraries:

- CompCert Clight (Blazy & Leroy, J. Autom. Reason. 2009)
 - Formal semantics of a deterministic sequential subset of C
- Flocq (Boldo & Melquiond, ARITH 2009)
 - Formalization of floating-point numbers
- Coq standard library
 - Formalization of real numbers

Our Design Choices: Which Formal Methods?

We use Coq + CompCert Clight + Flocq.

Advantages for trust:

- Unified verification framework
 - OK to combine proof libraries
- Formalization in the Gallina mathematical language of Coq
 - Can be trusted more easily than practical implementations (e.g. Fluctuat, Frama-C/Why3, etc.)
- Coq is the only setting where C, floating-point and real numbers are trustworthily mixed together

Our Approach and our Trusted Computing Base

We use Coq + CompCert Clight + Flocq.

- What do we need to trust?
 - Coq's underlying logic is sound
 - Implementation of Coq is sound wrt. Coq's logic
 - Coq standard library real numbers are consistent and faithful
 - Clight is faithful wrt. the corresponding subset of ISO C99
 - Flocq is faithful wrt. IEEE 754-2008 floating-point numbers
- Formalizations in the Gallina mathematical language
- Can be assessed more easily than practical implementations of verification tools

Verification of C Floating-Point Expressions

```
2.0f * (float) x - 3.0;
```

C floating-point
expression



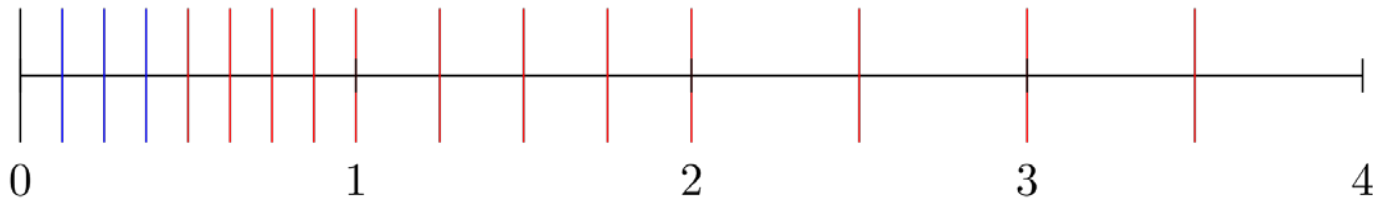
?

Real-number
semantics

Floating-Point Numbers

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

- A binary operation $a \mp b$ is not computed exactly
- Rounded from its ideal value
 - Example rounding mode: rounding to nearest
- What is the shape of the **rounding error**?



Floating-Point Numbers

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

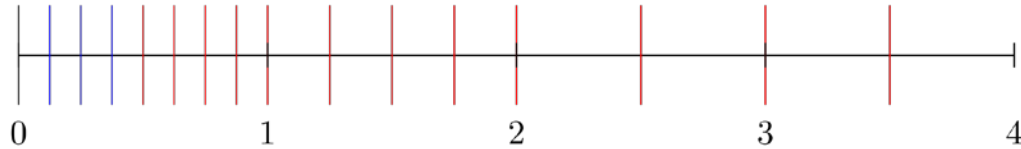
$$\pm m \cdot 2^e \quad 0 \leq m < 2^{\text{prec}}, \quad e_{\min} \leq e \leq e_{\max}$$

Floating-Point Numbers

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

$$\pm m \cdot 2^e \mid 0 \leq m < 2^3 = 8; -3 \leq e \leq -1$$

Example: prec = 3, emin = -3, emax = -1

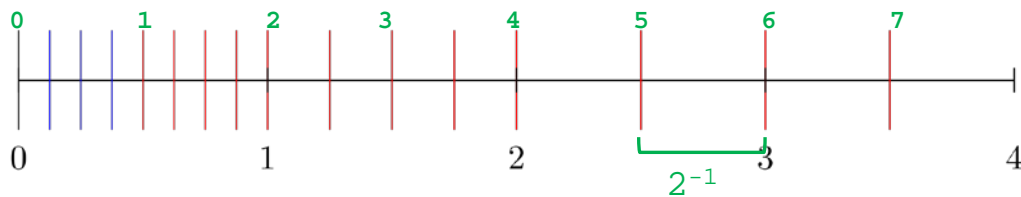
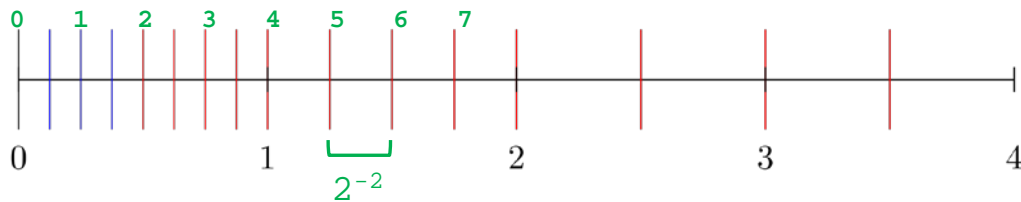
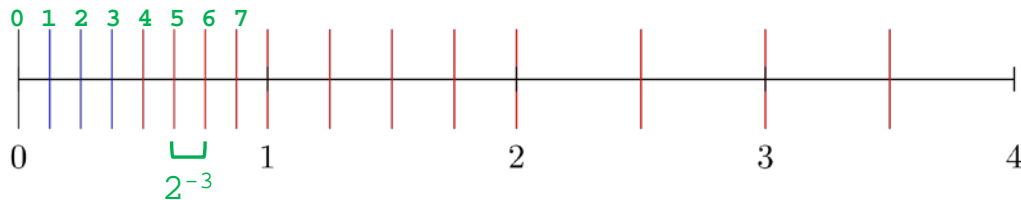


Floating-Point Numbers

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

$$\pm m \cdot 2^e \mid 0 \leq m < 2^3 = 8; -3 \leq e \leq -1$$

Example: prec = 3, emin = -3, emax = -1

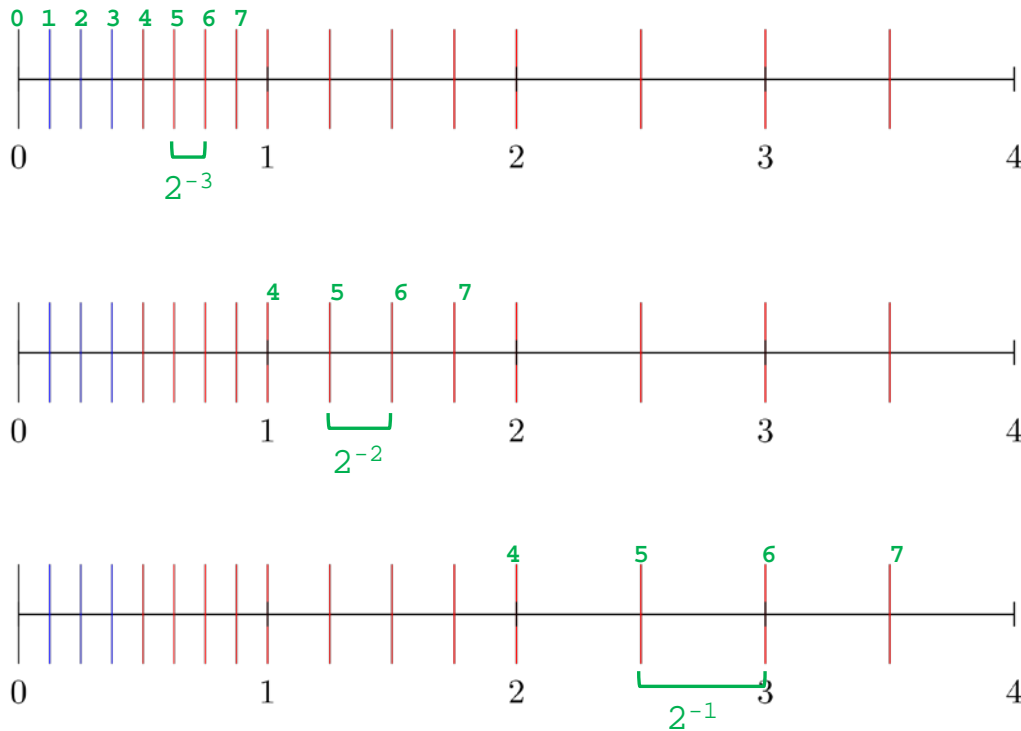


Floating-Point Numbers

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

$$\pm m \cdot 2^e \mid 0 \leq m < 2^3 = 8; -3 \leq e \leq -1$$

Example: prec = 3, emin = -3, emax = -1

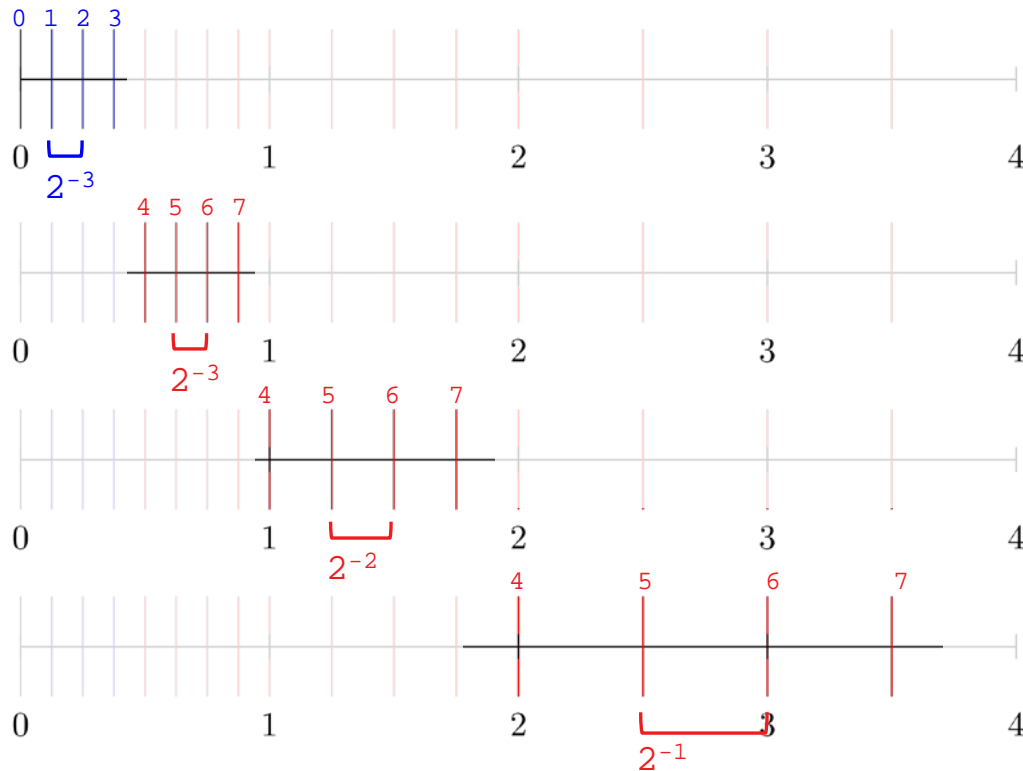


Floating-Point Numbers

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

$$\pm m \cdot 2^e \mid 0 \leq m < 2^3 = 8; -3 \leq e \leq -1$$

Example: prec = 3, emin = -3, emax = -1

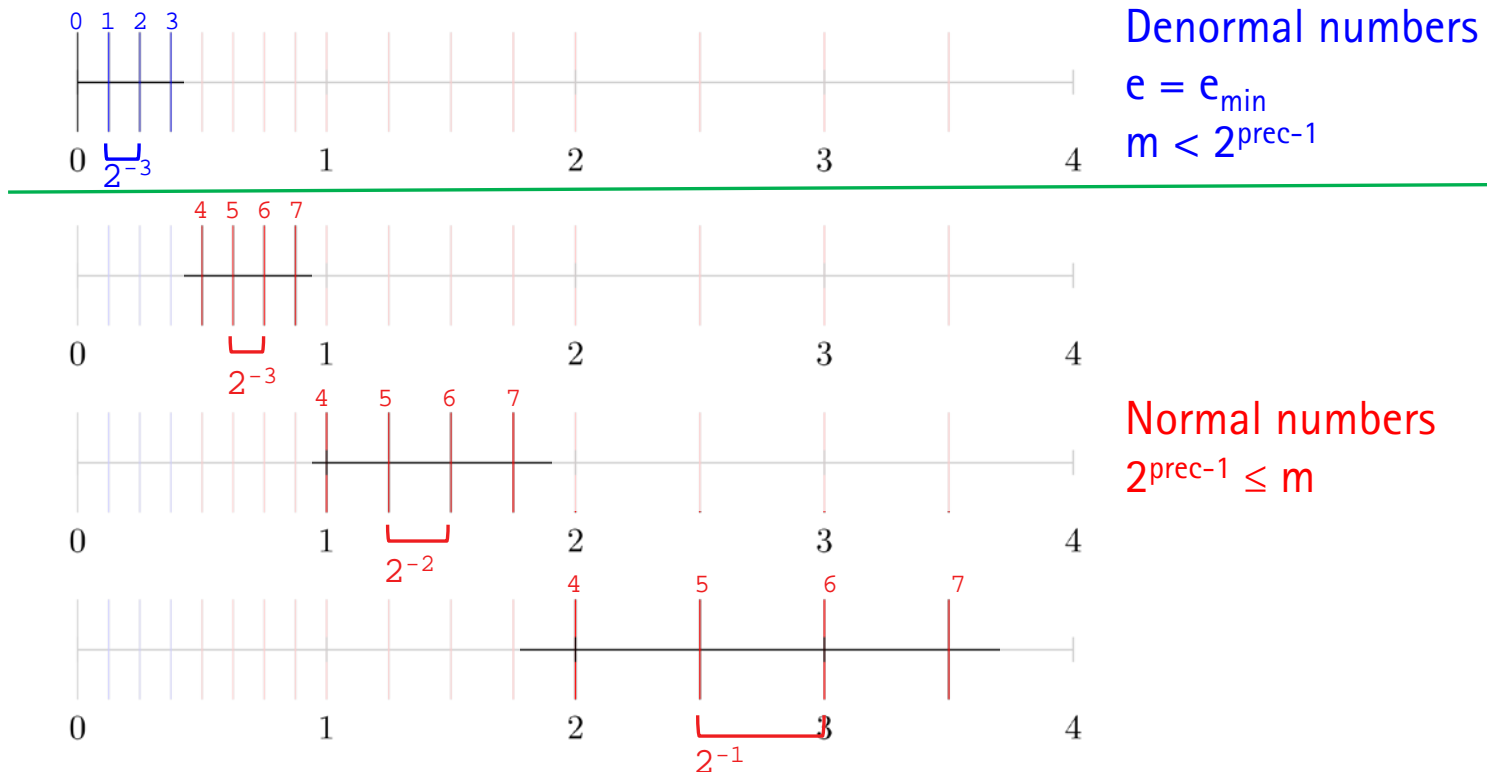


Floating-Point Numbers

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

$$\pm m \cdot 2^e \mid 0 \leq m < 2^3 = 8; -3 \leq e \leq -1$$

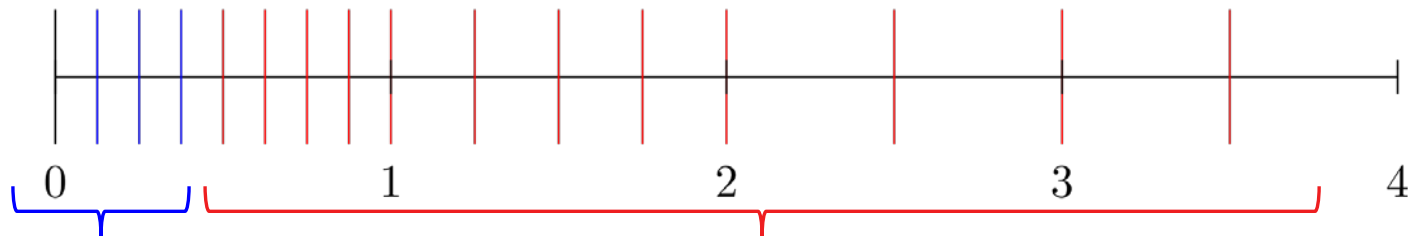
Example: prec = 3, emin = -3, emax = -1



Floating-Point Numbers

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

$$\pm m \cdot 2^e \quad 0 \leq m < 2^{\text{prec}}, \quad e_{\min} \leq e \leq e_{\max}$$



Denormal numbers

$$e = e_{\min}$$
$$m < 2^{\text{prec}-1}$$

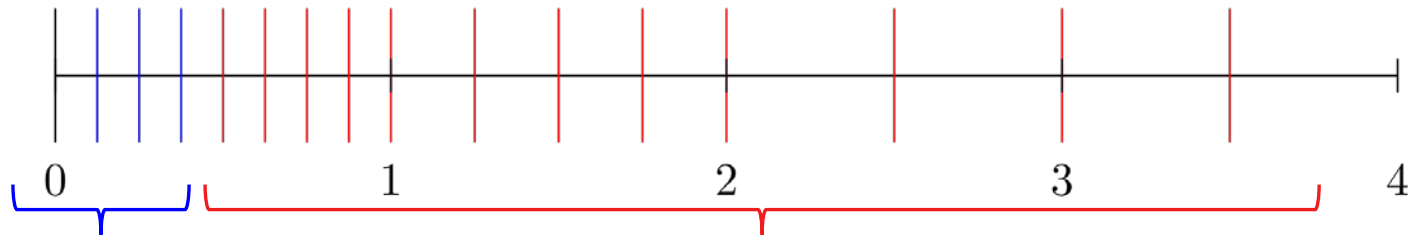
Normal numbers

$$2^{\text{prec}-1} \leq m$$

Floating-Point Numbers and Rounding Errors

IEEE 754-2008 modelled by Flocq (Boldo et al. 2009)

- A binary operation $a \text{ T } b$ is not computed exactly
- Rounded from its ideal value
 - Rounding mode: rounding to nearest, ties to even mantissa



Denormal:
 $(a \text{ T } b) + c$
 $|c| \leq 2^{\text{emin}-1}$

Normal:
 $(a \text{ T } b) (1 + d)$
 $|d| \leq 2^{-\text{prec}}$

General case: $(a \text{ T } b) (1 + d) + c$
with $c*d = 0$

Optimized Rounding Errors

- Normal numbers: $(a \mp b) (1 + d)$, if $|a \mp b|$ large enough and no overflow
- Denormal numbers: $(a \mp b) + e$, if $|a \mp b|$ small enough
- Sterbenz subtraction: $(a - b)$ if $a/2 \leq b \leq 2a$
- Multiply by power of 2 is always exact (unless overflow)
- Divide by power of 2 is exact if no gradual underflow

Flocq: correctness of floating-point arithmetic

Theorem Bplus_correct :

forall plus_nan m x y,

is_finite x = true ->

is_finite y = true ->

if $\text{Rlt_bool } (\text{Rabs } (\text{round radix2 fexp } (\text{round_mode } m) (\text{B2R } x + \text{B2R } y))) (\text{bpow radix2 emax})$ then

$\text{B2R } (\text{Bplus plus_nan } m \ x \ y) = \text{round radix2 fexp } (\text{round_mode } m) (\text{B2R } x + \text{B2R } y) \wedge$

$\text{is_finite } (\text{Bplus plus_nan } m \ x \ y) = \text{true} \wedge$

$\text{Bsign } (\text{Bplus plus_nan } m \ x \ y) =$

match Rcompare (B2R x + B2R y) 0 with

| Eq => match m with mode_DN => orb (Bsign x) (Bsign y)

| _ => andb (Bsign x) (Bsign y) end

| Lt => true

| Gt => false

end

else

$(\text{B2FF } (\text{Bplus plus_nan } m \ x \ y) = \text{binary_overflow } m (\text{Bsign } x) \wedge \text{Bsign } x = \text{Bsign } y).$

No overflow

Theorem relative_error_ex :

forall x,

$(\text{bpow } \text{emin} \leq \text{Rabs } x) \% R \leftrightarrow$

exists eps,

$(\text{Rabs } \text{eps} < \text{bpow } (-p + 1)) \% R \wedge \text{round beta fexp rnd } x = (x * (1 + \text{eps})) \% R.$

Normal numbers

Flocq: correctness of floating-point arithmetic

Theorem Bplus_correct :

forall plus_nan m x y,

is_finite x = true ->

is_finite y = true ->

if $\text{Rlt_bool } (\text{Rabs } (\text{round_radix2_fexp } (\text{round_mode } m) (\text{B2R } x + \text{B2R } y))) (\text{bpow_radix2_emax})$ then

$\text{B2R } (\text{Bplus_nan } m \ x \ y) = \text{round_radix2_fexp } (\text{round_mode } m) (\text{B2R } x + \text{B2R } y) \wedge$

is_fini

Bsign

matc

| Eo

| Lt

| Lt

| Gt

end

else

(B2FF

No overflow

Better tackle them
automatically

Theorem relative_error_ex :

forall x,

$(\text{bpow_emin} \leq \text{Rabs } x) \% \text{R} \leftrightarrow$

exists eps,

$(\text{Rabs } \text{eps} < \text{bpow } (-p + 1)) \% \text{R} \wedge \text{round_beta_fexp_rnd } x = (x * (1 + \text{eps})) \% \text{R}.$

Normal numbers

Rounding Error Terms

Optimized cases

- Normal numbers: $(a \text{ T } b) (1 + d)$, if $|a \text{ T } b|$ large enough and no overflow
- Denormal numbers: $(a \text{ T } b) + e$, if $|a \text{ T } b|$ small enough
- Sterbenz subtraction: $(a - b)$ if $a/2 \leq b \leq 2a$
- Multiply by power of 2 is always exact (unless overflow)
- Divide by power of 2 is exact if no gradual underflow

Our VCFloat approach:

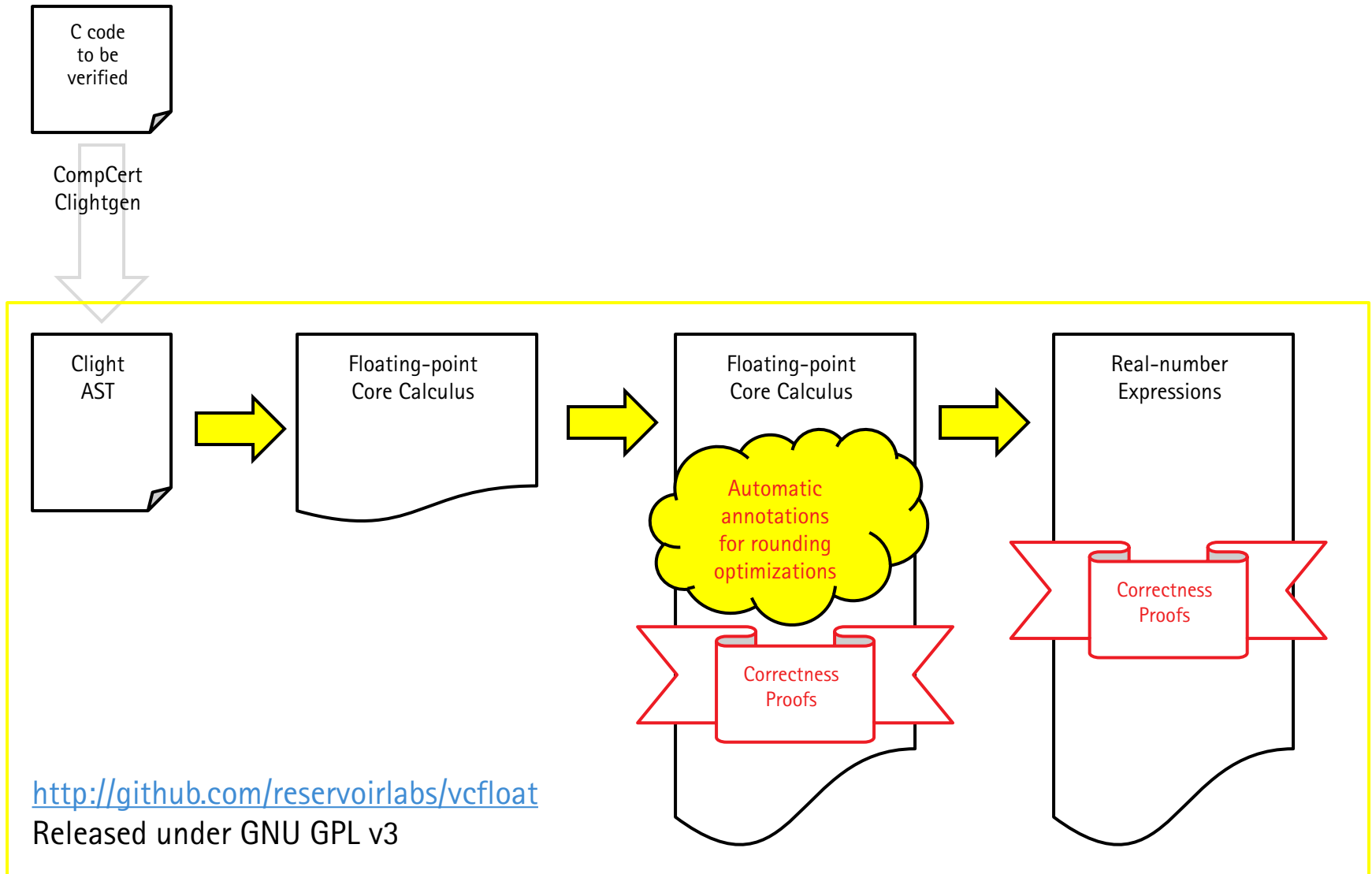
- Automatically generate validity conditions
- Automatically check them on the fly
- Add annotations for optimized rounding

Verification of Rounding Error Terms

Use Coq-Interval (Melquiond 2015) to automatically check validity conditions

- Automatic certified interval arithmetic
- Reduce correlation issues:
 - Bisection (branch-and-bound)
 - Automatic differentiation
 - Taylor models
- Used for all rounding errors
- All computations within Coq: consumes most proof checking time and memory in overall proof
- Stress test

Our Verification Framework: VCFloat



Verification of Rounding Error Terms

```
2.0f * (float) x - 3.0;
```

C floating-point
expression

Verification of Rounding Error Terms

`2.0f * (float) x - 3.0;`

C floating-point
expression

$(2_{(24,128)} \otimes [x]_{(24,128)}) \ominus 3_{(53,1024)}$

Core floating-point
expression

Verification of Rounding Error Terms

`2.0f * (float) x - 3.0;`

C floating-point
expression

$(2_{(24,128)} \otimes [x]_{(24,128)}) \ominus 3_{(53,1024)}$

Core floating-point
expression

Assume x in $[1, 2]$

Verification of Rounding Error Terms

`2.0f * (float) x - 3.0;`

C floating-point
expression

$(2_{(24,128)} \otimes [x]_{(24,128)}) \ominus 3_{(53,1024)}$

Core floating-point
expression

Assume x in $[1, 2]$

$(2_{(24,128)} \otimes [x]_{(24,128)}^{\text{Norm}}) \ominus 3_{(53,1024)}$

Annotated floating-point
expression

Because $2^{-125} \leq |x| < 2^{128}$

Verification of Rounding Error Terms

`2.0f * (float) x - 3.0;`

C floating-point expression

$(2_{(24,128)} \otimes [x]_{(24,128)}) \ominus 3_{(53,1024)}$

Core floating-point expression

Assume x in $[1, 2]$

$(2^1 \otimes \cancel{\otimes} [x]_{(24,128)}^{\text{Norm}}) \ominus 3_{(53,1024)}$

Annotated floating-point expression

Because $2^{-125} \leq |x| < 2^{128}$

And for all d in $[-2^{-24}, 2^{-24}]$, $|2 * (x * (1 + d))| < 2^{128}$

Verification of Rounding Error Terms

`2.0f * (float) x - 3.0;`

C floating-point expression

$(2_{(24,128)} \otimes [x]_{(24,128)}) \ominus 3_{(53,1024)}$

Core floating-point expression

Assume x in $[1, 2]$

$(2^1 \otimes [x]_{(24,128)}^{\text{Norm}}) \ominus 3_{(53,1024)}^{\text{Sterbenz}}$

Annotated floating-point expression

Because $2^{-125} \leq |x| < 2^{128}$

And for all d in $[-2^{-24}, 2^{-24}]$, $|2 * (x * (1 + d))| < 2^{128}$

And for all d in $[-2^{-24}, 2^{-24}]$, $3/2 \leq 2 * (x * (1 + d)) \leq 3*2$

Verification of Rounding Error Terms

`2.0f * (float) x - 3.0;`

C floating-point expression

$(2_{(24,128)} \otimes [x]_{(24,128)}) \ominus 3_{(53,1024)}$

Core floating-point expression

Assume x in $[1, 2]$

$(2^1 \otimes [x]_{(24,128)}^{\text{Norm}}) \ominus 3_{(53,1024)}^{\text{Sterbenz}}$

Annotated floating-point expression

Because $2^{-125} \leq |x| < 2^{128}$

And for all d in $[-2^{-24}, 2^{-24}]$, $|2 * (x * (1 + d))| < 2^{128}$

And for all d in $[-2^{-24}, 2^{-24}]$, $3/2 \leq 2 * (x * (1 + d)) \leq 3 * 2$

$2 \times (x \times (1 + \delta)) - 3$
 For some δ in $[-2^{-24}, 2^{-24}]$

Real-number expression with error terms

Formal Verification of Error Bounds

DEMO

Formal Verification of Error Bounds

CONCLUSIONS

Conclusion and Future Work

C programs with floating-point computations can now be fully verified within Coq with a TCB smaller than ever:

- Implementation of Coq, faithfulness of Clight and Flocq

Energy-efficient approximate implementations can be allowed in critical applications, once their error bounds are certified correct

Floating-point steps are now automated, but overall more practical improvements are still desirable:

- Integer handling
- C control flow: Integrate into Verifiable C

Thank you!

Coq library and proofs available online:

- <http://github.com/reservoirlabs/vcfloat>
- Mostly released under GNU GPL v3
- Currently based on Coq 8.5beta2, will be adapted to 8.5rc1
- Stress test for Flocq, Coq-Interval, and computations within Coq

For more information:

- ramananandro@reservoir.com