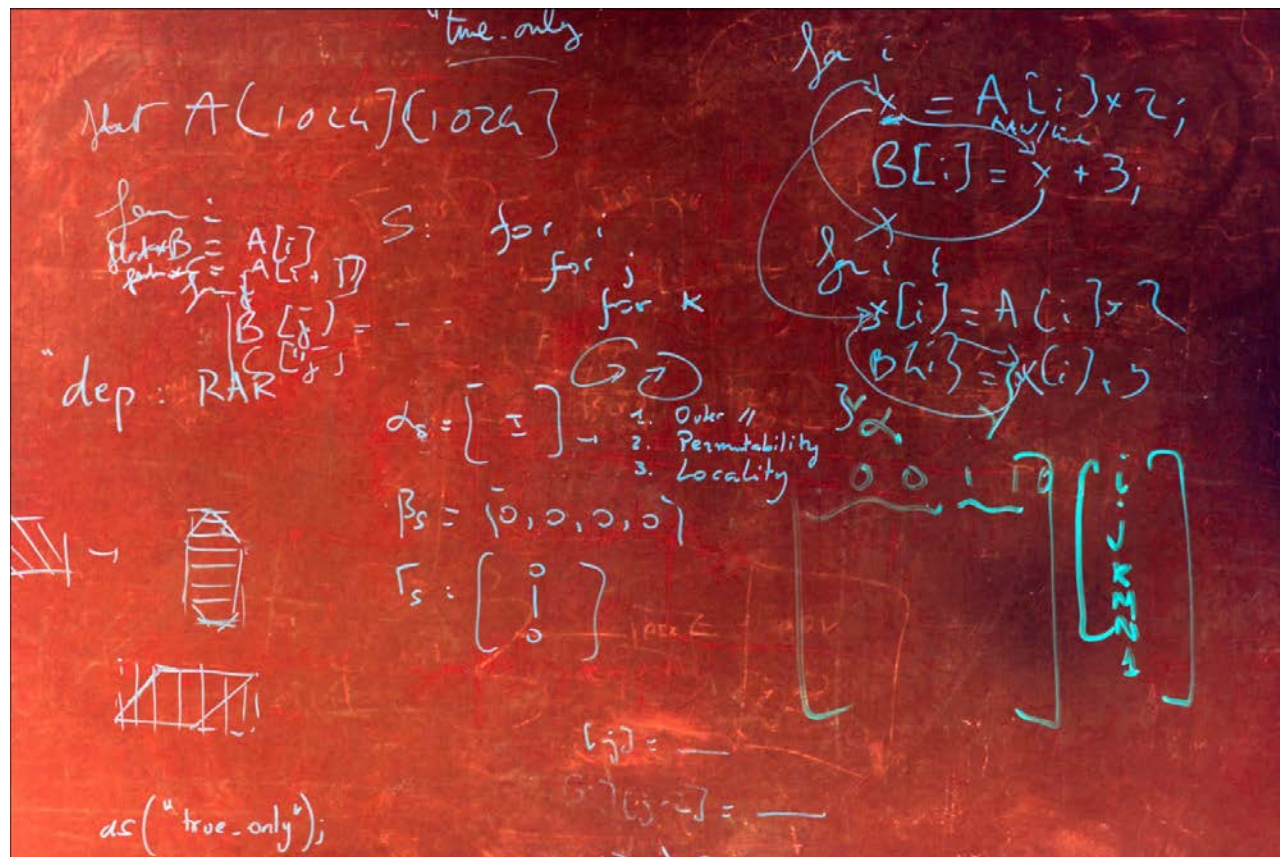# Formal Verification of C Programs with Floating-Point Computations
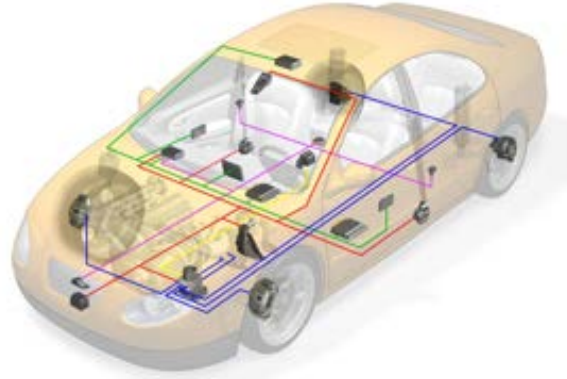## Certified Error Bounds for Signal Processing

**Tahina Ramananandro**

Paul Mountcastle, Benoît Meister, Richard Lethin

# Overview



- Sensor systems data processing
  - Positioning, obstacle avoidance, radar imaging, etc.
- Problems
  - Numerical optimizations for energy/time efficiency
  - Correctness and accuracy
  - Which guarantees on the actual software code?

# Overview

Can we reason about sensor systems:

- With new performance optimizations for data processing
- In a very **deep** way, with strong correctness guarantees down to the actual code?

Our contributions

- VCFloat: proof library and tactics to verify C programs with floating-point computations
- Example use case: a radar algorithm and its C implementation

Use cases

- Cyber-physical systems
- Lightweight UAVs (copters), cars
- Military, transportation, medicine, etc.

# The High-Level Problem

Goal: energy-efficient implementations of numerical algorithms

- Naïve implementations consume time and energy
- Main ideas for performance improvement:
  - compute in lower-precision floating-point
  - introduce approximations
- Problem: **uncertainty** introduced by errors 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
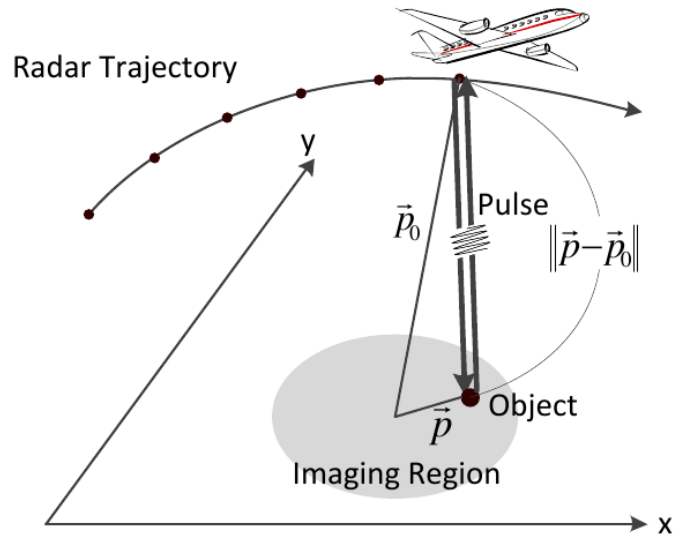- Energy measurements: ∼10-20% saved on Intel Haswell

## 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
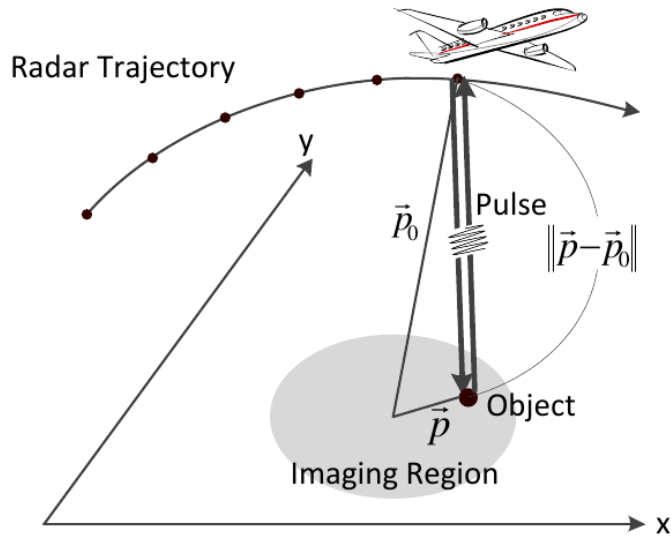


for all pixels $x$ do
    for all pulses $p$ do
        $R$ = distance between platform and pixel $x$ for pulse $p$
        $s$ = sample interpolated from pulse $p$ in neighborhood of range $R$
        Apply phase correction to sample $s$ based on range $R$
        Accumulate sample $s$ into pixel $x$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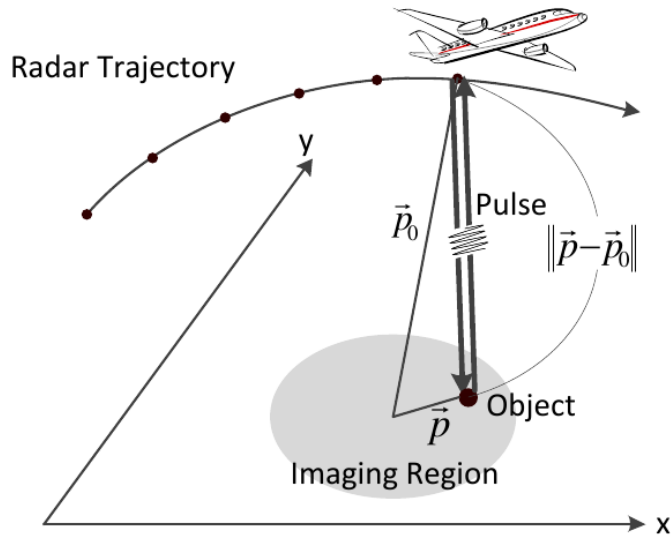
# SAR Backprojection



```
for y := 0 to BP_NPIX_Y − 1 do
    py := (y + (1−BP_NPIX_Y)/2) × dxdy
    for x := 0 to BP_NPIX_X − 1 do
        px := (x + (1−BP_NPIX_X)/2) × dxdy
        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
```

Real−number algorithm

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

# SAR Backprojection



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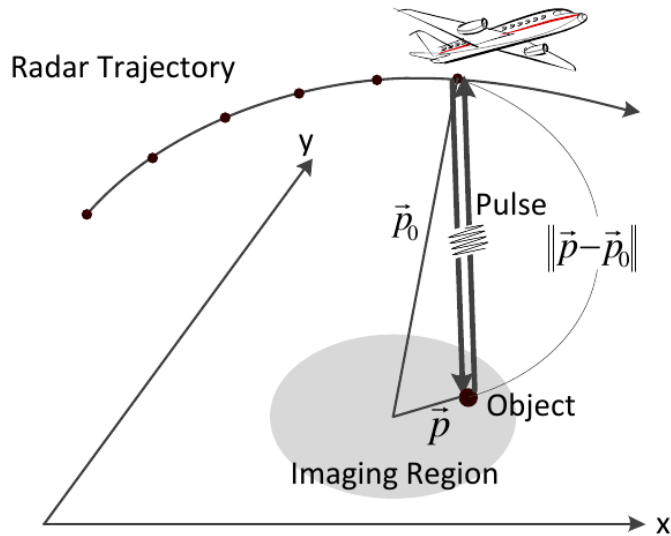$$\textbf{return } \underline{image}$$

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



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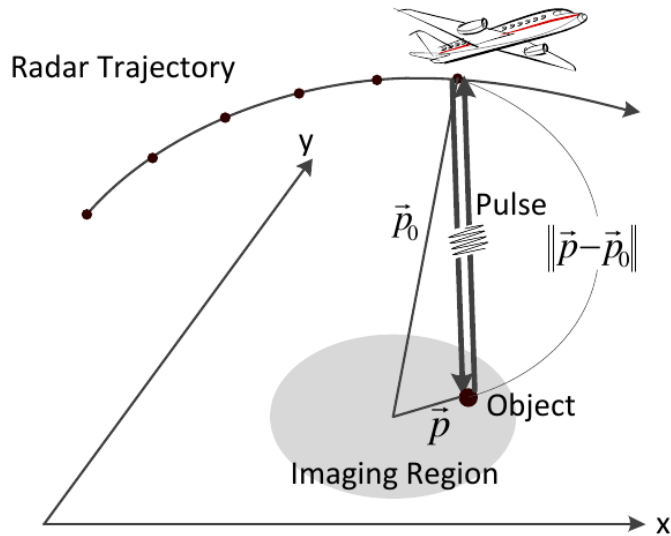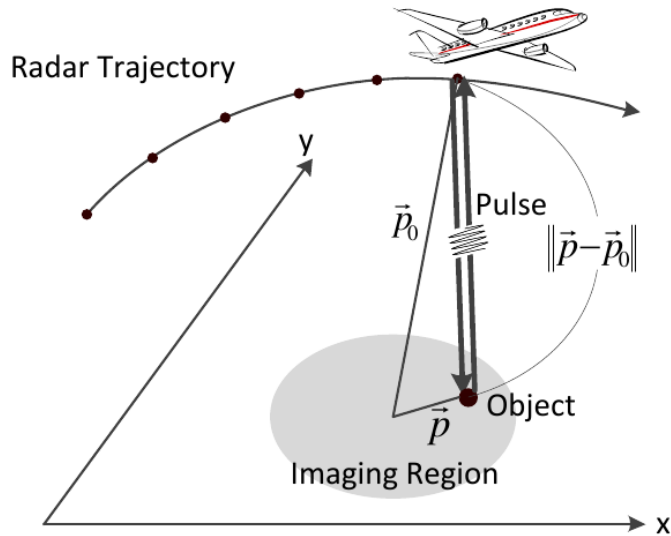$$\textbf{end for}$$
$$\textbf{return } image$$

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



Radar Trajectory

y

$\vec{P}_0$  Pulse

$\|\vec{p}-\vec{p}_0\|$

Object

$\vec{p}$

Imaging Region

x

square root

sine

floating-point computations

```
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(…) … ;
    …
   }
  }
 }
}
```

C Implementation

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

# SAR Backprojection



approximate
square root

```
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(…) … ;
       …
     }
   }
 }
}
```
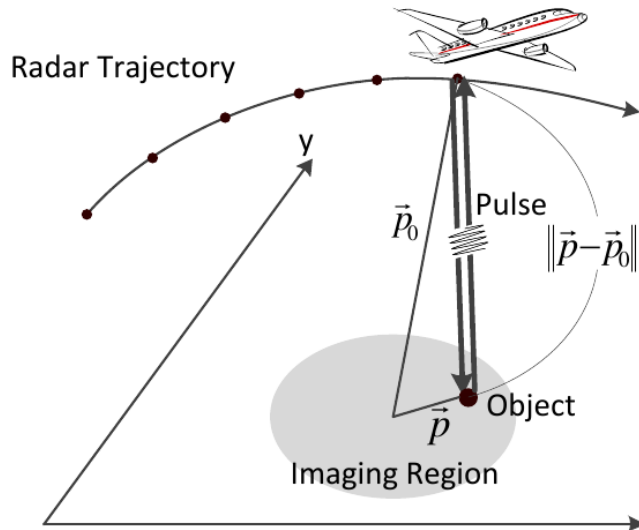
approximate
sine

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



approximate square root

approximate sine

```c
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

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:

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

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)

∀ P ˋ(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**

∃ m' , star Clight.step2
        (Callstate fn_sar_backprojection ...) (Returnstate Vundef Kstop m') ∧
∃ 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)) ∧
     ∀ y, (y < BP_NPIX_Y)%nat → ∀ x, (x < BP_NPIX_X)%nat →
       let ir := image_r (y × BP_NPIX_X + x))   in
       let ii := image_i (y × BP_NPIX_X + x))   in
       is_finite _ _ ir = true ∧ is_finite _ _ ii = true ∧
       let (tr, ti) := SARBackProj.sar_backprojection y x     in
       Rabs (B2R _ _ ir – tr) ≤ pixel_bound ∧
       Rabs (B2R _ _ ii – ti) ≤ pixel_bound.

Coq

**Conclusions**

# Final Correctness Statement (slightly simplified)

∀ P `(HYPS: SARHypotheses P)  m  | Hypotheses on input data
    (Hm: holds m (P ++
            Pperm_int bir oir (BP_NPIX_X × BP_NPIX_Y) ++  | Memory contents
            Pperm_int bii oii (BP_NPIX_X × BP_NPIX_Y)),    | and permissions

∃ m' , star Clight.step2
            (Callstate fn_sar_backprojection ...) (Returnstate Vundef Kstop m') ∧
∃ 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)) ∧
        ∀ y, (y < BP_NPIX_Y)%nat → ∀ x, (x < BP_NPIX_X)%nat →
            let ir := image_r (y × BP_NPIX_X + x))   in
            let ii := image_i (y × BP_NPIX_X + x))   in
            is_finite _ _ ir = true ∧ is_finite _ _ ii = true ∧
            let (tr, ti) := SARBackProj.sar_backprojection y x     in
            Rabs (B2R _ _ ir – tr) ≤ pixel_bound ∧
            Rabs (B2R _ _ ii – ti) ≤ pixel_bound.

Coq

Conclusions

# Final Correctness Statement (slightly simplified)

∀ 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**

∃ m', star Clight.step2
        (Callstate fn_sar_backprojection ...) (Returnstate Vundef Kstop m') ∧

**C code runs**

∃ 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)) ∧

**Memory contents**

∀ y, (y < BP_NPIX_Y)%nat → ∀ x, (x < BP_NPIX_X)%nat →
      let ir := image_r (y × BP_NPIX_X + x))   in
      let ii := image_i (y × BP_NPIX_X + x))   in
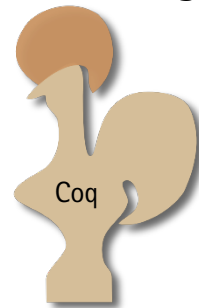      is_finite _ _ ir = true ∧ is_finite _ _ ii = true ∧

**FP does not overflow**

      let (tr, ti) := SARBackProj.sar_backprojection y x     in
      Rabs (B2R _ _ ir – tr) ≤ pixel_bound ∧
      Rabs (B2R _ _ ii – ti) ≤ pixel_bound.

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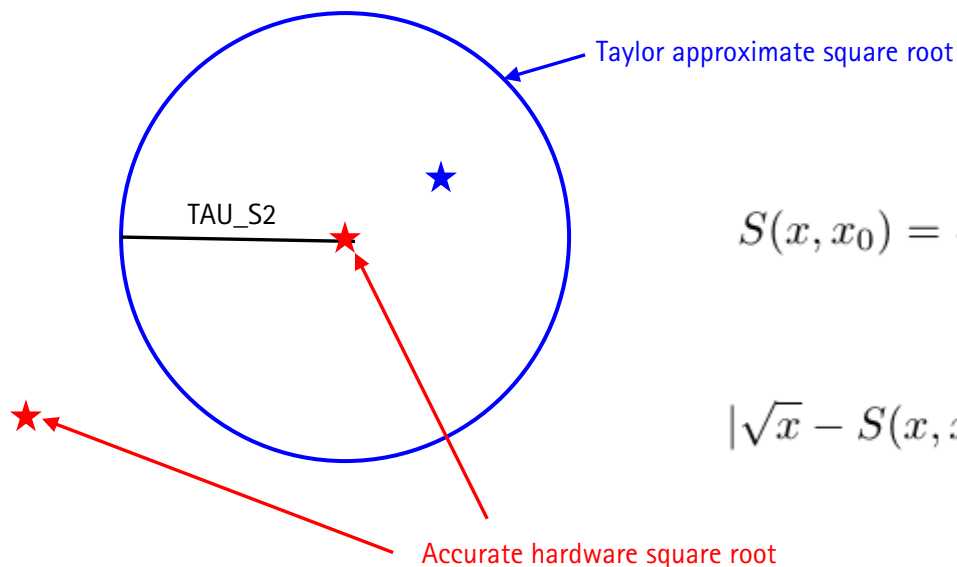
Coq

# 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 $\pi$ 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

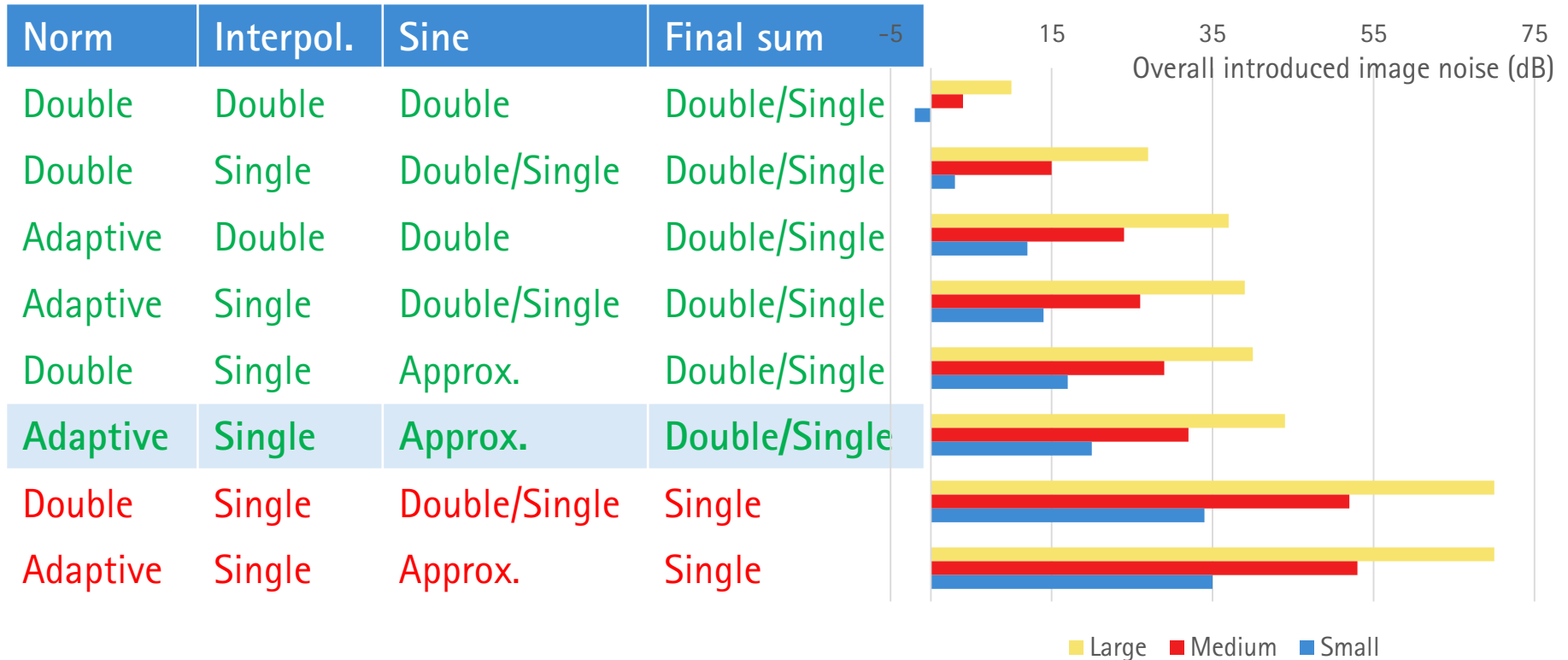Taylor approximate square root

TAU_S2

$$S(x, x_0) = \sqrt{x_0} + \frac{x - x_0}{2\sqrt{x_0}} - \frac{(x - x_0)^2}{8(\sqrt{x_0})^3}$$

$$|\sqrt{x} - S(x, x_0)| \leq \frac{(\text{TAU\_S2})^3}{16 \times (x_0 - \text{TAU\_S2})^{5/2}}$$

Accurate hardware square root
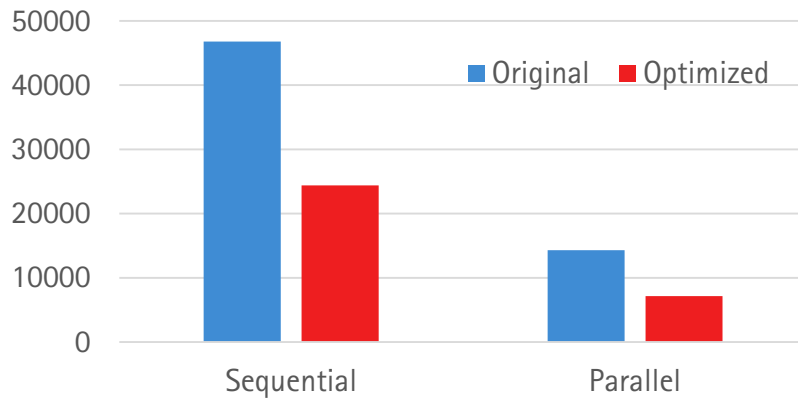
# Precision Results

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

| Norm | Interpol. | Sine | Final sum |
|------|-----------|------|-----------|
| Double | Double | Double | Double/Single |
| Double | Single | Double/Single | Double/Single |
| Adaptive | Double | Double | Double/Single |
| Adaptive | Single | Double/Single | Double/Single |
| Double | Single | Approx. | Double/Single |
| **Adaptive** | **Single** | **Approx.** | **Double/Single** |
| Double | Single | Double/Single | Single |
| Adaptive | Single | Approx. | Single |

Overall introduced image noise (dB)

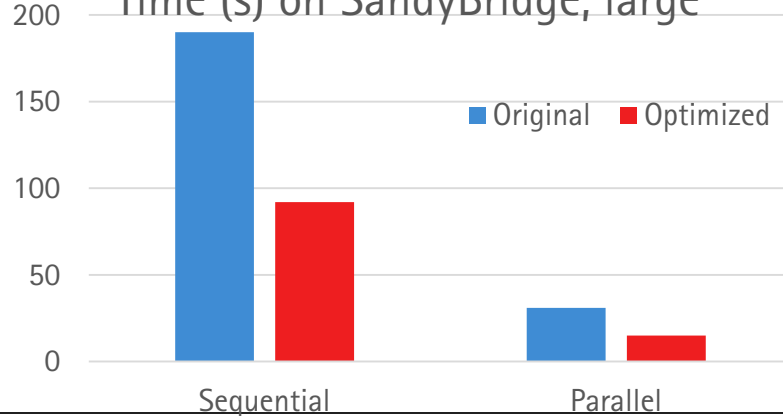-5    15    35    55    75

Large  Medium  Small

# Performance Measurements for Optimized C Code

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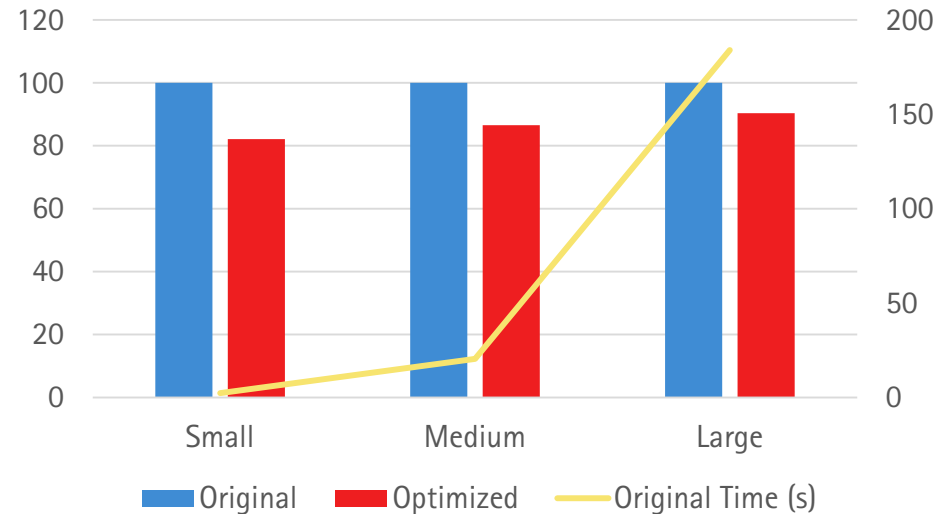- Intel Haswell: energy model unknown, time instead



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)

∀ 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**

∃ m' , star Clight.step2
      (Callstate fn_sar_backprojection ...) (Returnstate Vundef Kstop m') ∧

**C code runs**

∃ 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)) ∧

**Memory contents**

∀ y, (y < BP_NPIX_Y)%nat → ∀ x, (x < BP_NPIX_X)%nat →
    let ir := image_r (y × BP_NPIX_X + x))   in
    let ii := image_i (y × BP_NPIX_X + x))   in
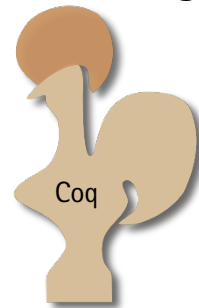    is_finite _ _ ir = true ∧ is_finite _ _ ii = true ∧

**FP does not overflow**

    let (tr, ti) := SARBackProj.sar_backprojection y x     in
    Rabs (B2R _ _ ir – tr) ≤ pixel_bound ∧
    Rabs (B2R _ _ ii – ti) ≤ pixel_bound.

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

Coq

# Final Correctness Statement (slightly simplified)

∀ 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**

∃ m', star Clight.step2
   (Callstate fn_sar_backprojection ...) (Returnstate Vundef Kstop m') ∧

**C code runs**

## CompCert Clight

∃ 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)) ∧
∀ y, (y < BP_NPIX_Y)%nat → ∀ x, (x < BP_NPIX_X)%nat →
   let ir := image_r (y × BP_NPIX_X + x))   in
   let ii := image_i (y × BP_NPIX_X + x))   in

**Memory contents**

is_finite _ _ ir = true ∧ is_finite _ _ ii = true ∧

**FP does not overflow**

   let (tr, ti) := SARBackProj.sar_backprojection y x     in
   Rabs (B2R _ _ ir – tr) ≤ pixel_bound ∧
   Rabs (B2R _ _ ii – ti) ≤ pixel_bound.

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

## Flocq

Coq

# 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 T 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)

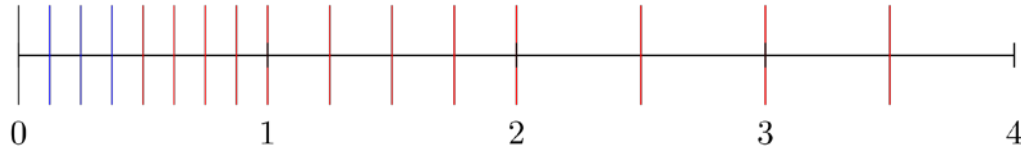$$\pm m \cdot 2^e \qquad\qquad 0 \leq m < 2^{prec}, \ e_{min} \leq e \leq e_{max}$$

# Floating-Point Numbers

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

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

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

# Floating-Point Numbers

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

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

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

# Floating-Point Numbers

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

$$\pm m \cdot 2^e \mid 0 \leq m < 2^3 = 8; \; -3 \leq e \leq -1$$
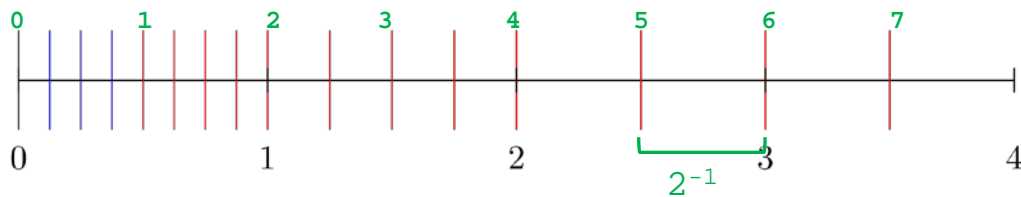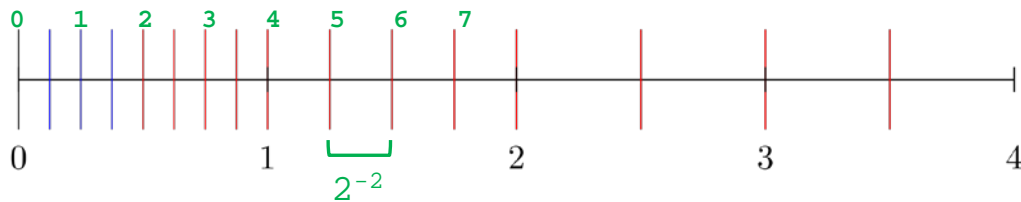
Example: prec = 3, emin = –3, emax = –1

# Floating-Point Numbers

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

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

Example: prec = 3, emin = –3, emax = –1
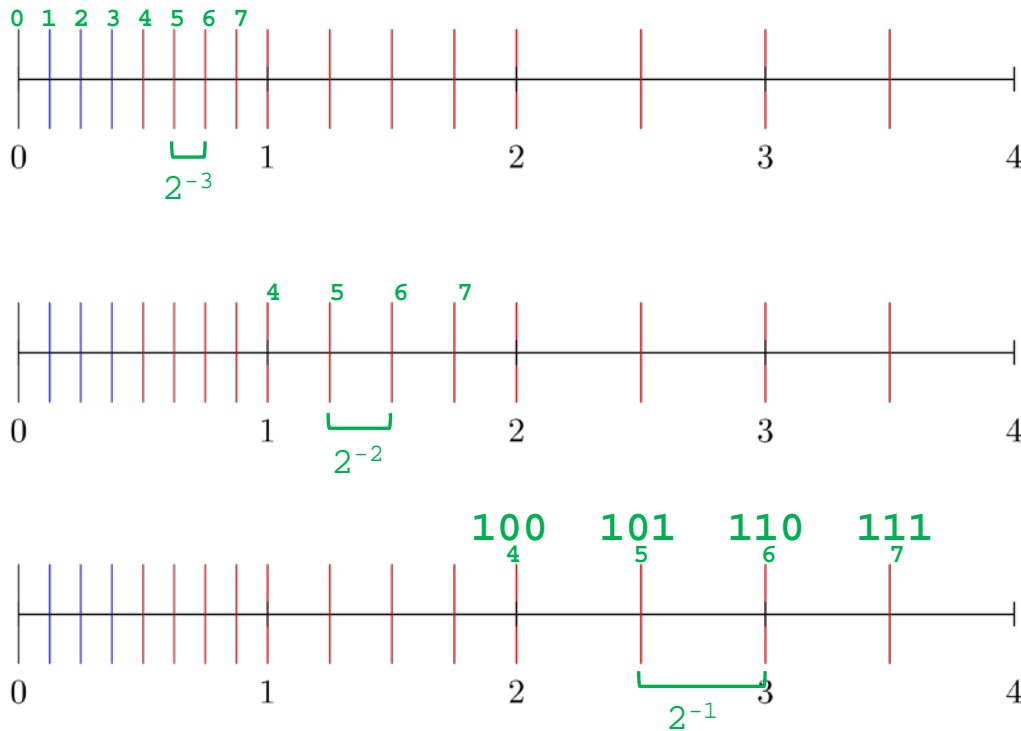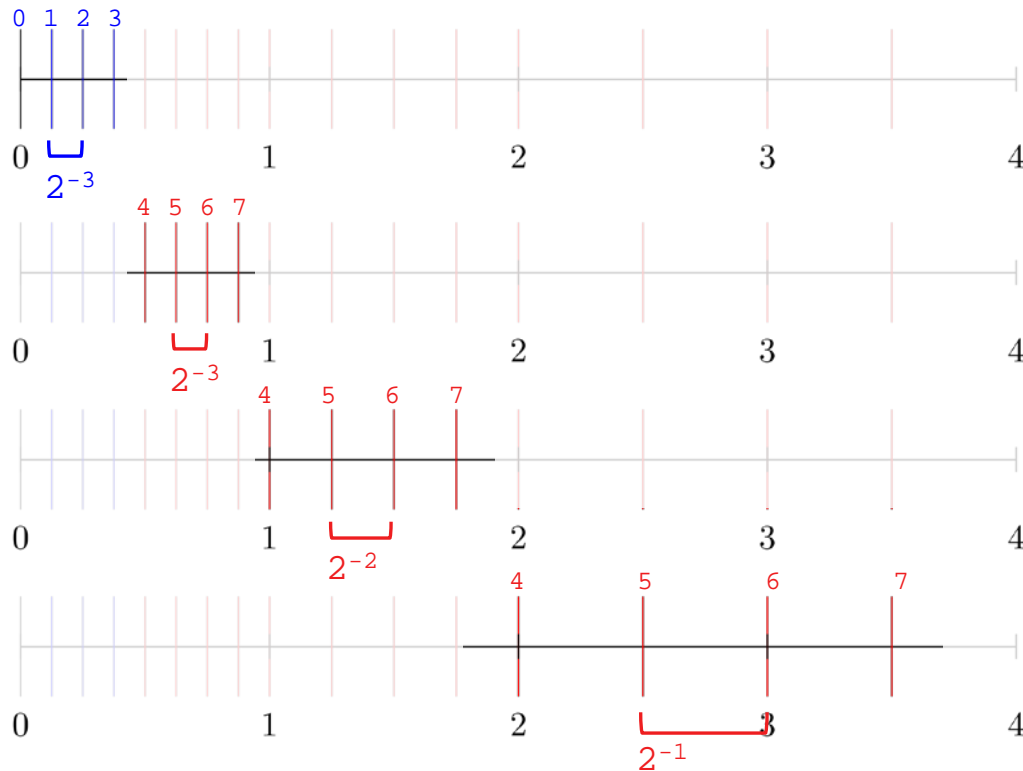
# Floating-Point Numbers

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



Denormal numbers
$e = e_{min}$
$m < 2^{prec-1}$

Normal numbers
$2^{prec-1} \leq m$

# Floating-Point Numbers

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

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



Denormal numbers
$e = e_{min}$
$m < 2^{prec-1}$

Normal numbers
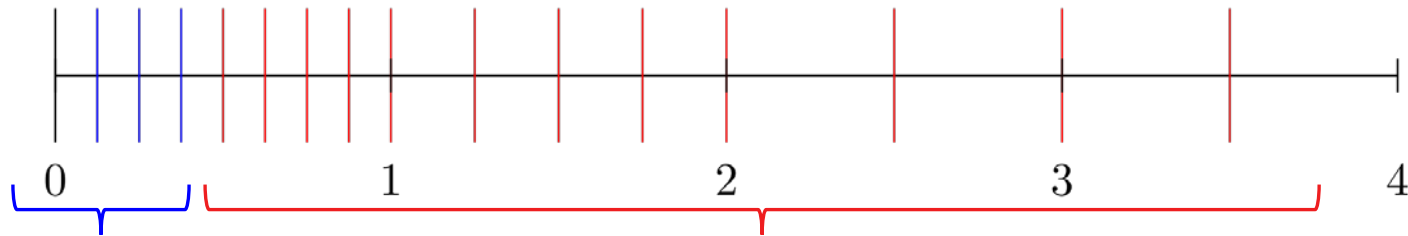$2^{prec-1} \leq m$

# Floating-Point Numbers and Rounding Errors

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

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



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

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

General case: (a T b) (1 + d) + c
with c*d = 0

# Optimized Rounding Errors

- Normal numbers: (a T b) (1 + d), if |a T b| large enough and no overflow
- Denormal numbers: (a T b) + e,  if |a T b| small enough
- Sterbenz subtraction: (a − b)     if a/2 <= b <= 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 Rlt_bool (Rabs (round radix2 fexp (round_mode m) (B2R x + B2R y))) (bpow radix2 emax) then
  B2R (Bplus plus_nan m x y) = round radix2 fexp (round_mode m) (B2R x + B2R y) /\
  is_finite (Bplus plus_nan m x y) = true /\
  Bsign (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
  (B2FF (Bplus plus_nan m x y) = binary_overflow m (Bsign x) /\ Bsign x = Bsign y).

**No overflow**

Theorem relative_error_ex :
 forall x,
 (bpow emin <= Rabs x)%R ->
 exists eps,
 (Rabs eps < bpow (-p + 1))%R /\ round beta fexp rnd x = (x * (1 + 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 Rlt_bool (Rabs (round radix2 fexp (round_mode m) (B2R x + B2R y))) (bpow radix2 emax) then
    B2R (Bplus plus_nan m x y) = round radix2 fexp (round_mode m) (B2R x + B2R y) /\
    is_fini
    Bsign
      matc
        | Eq
        | Lt
        | Gt
      end
  else
    (B2FF

*No overflow*

## Better automate their use

Theorem relative_error_ex :
  forall x,
  (bpow emin <= Rabs x)%R ->
  exists eps,
  (Rabs eps < bpow (-p + 1))%R /\ round beta fexp rnd x = (x * (1 + eps))%R.

*Normal numbers*

# Rounding Error Terms

Optimized cases

- Normal numbers: (a T b) (1 + d), if |a T b| large enough and no overflow
- Denormal numbers: (a T b) + e, if |a T b| small enough
- Sterbenz subtraction: (a − b)  if a/2 <= b <= 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



http://github.com/reservoirlabs/vcfloat
Released under GNU GPL v3

# 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

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

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

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

Core floating-point expression

Assume x in [1, 2]

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

Annotated floating-point expression

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

# 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]

$$( \; 2^1 \quad \textbf{✗} \; [x]_{(24,128)}^{\text{Norm}} ) \ominus 3_{(53,1024)}$$

Annotated floating-point expression

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

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

# Verification of Rounding Error Terms

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

C floating-point expression

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

Core floating-point expression

Assume x in [1, 2]

$\left( \; 2^1 \quad \textcolor{red}{\times}^{\text{Norm}} [x]_{(24,128)} \right) \textcolor{red}{\ominus}^{\text{Sterbenz}} 3_{(53,1024)}$

Annotated floating-point expression

Because 2ˆ–125 ≤ |x| < 2ˆ128

And for all d in [–2ˆ–24, 2ˆ–24], |2 * (x * (1 + d))| < 2ˆ128

And for all d in [–2ˆ–24, 2ˆ–24], 3/2 ≤ 2 * (x * (1 + d))  ≤ 3*2

# Verification of Rounding Error Terms

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

C floating-point expression

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

Core floating-point expression

Assume x in [1, 2]

$$\left( \quad 2^1 \quad \times \quad [x]^{Norm}_{(24,128)}\right) \ominus^{Sterbenz} 3_{(53,1024)}$$

Annotated floating-point expression

Because 2^-125 ≤ |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 ≤ 2 * (x * (1 + d)) ≤ 3*2

$$2 \times (x \times (1 + \delta)) - 3$$

For some $\delta$ in $[-2^{-24}, 2^{-24}]$

Real-number expression with error terms

Formal Verification of Error Bounds

# DEMO

Formal Verification of Error Bounds

# CONCLUSIONS

# Ongoing Work

- Floating-point steps are now automated
- Overall more practical improvements ongoing:
  - Integer handling
  - C control flow: Integrate into Verifiable C (Princeton)
  - More statistical support for hypotheses on input data to tighten error bounds

# Conclusion and Opportunities

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

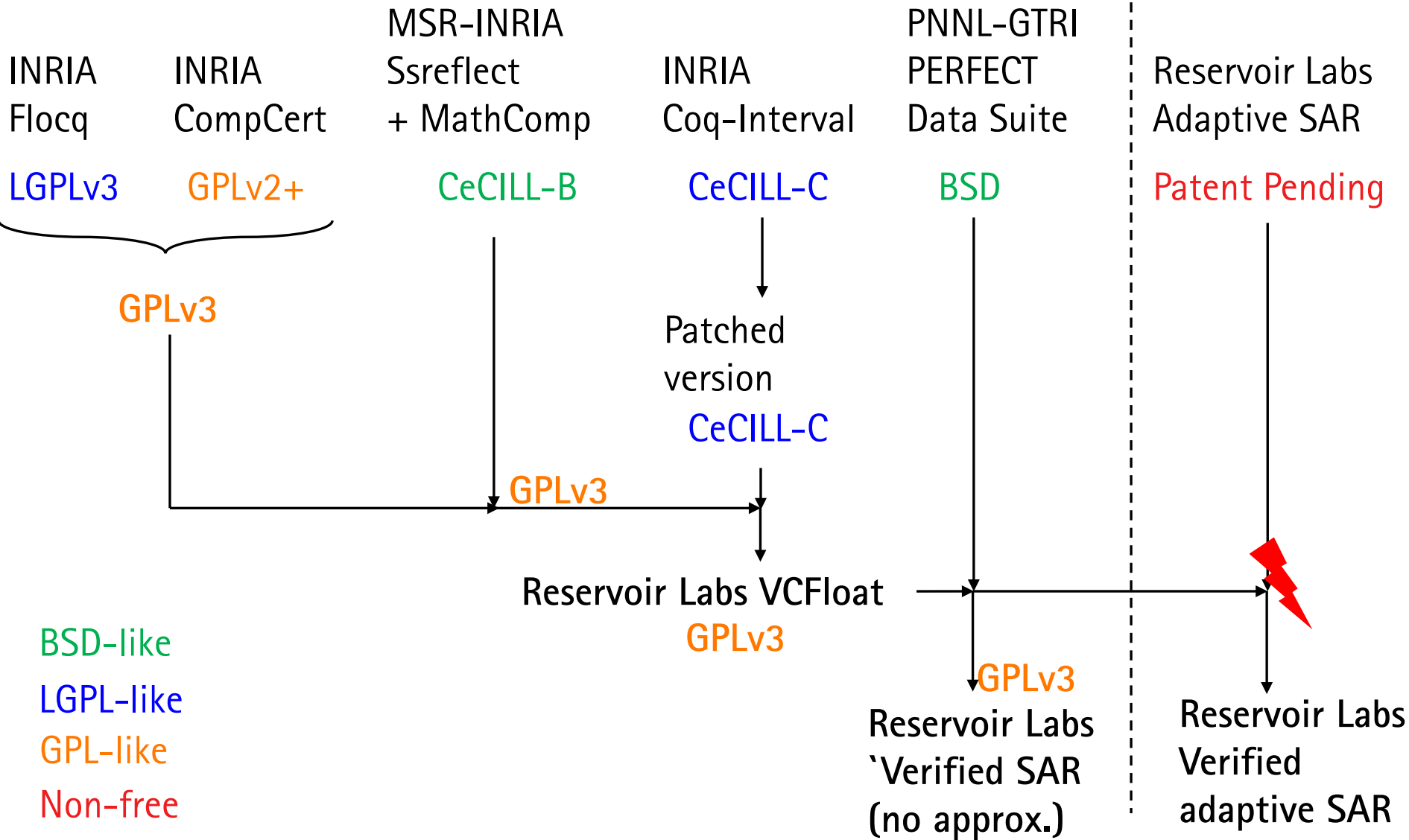- Implementation of Coq, faithfulness of Clight and Flocq

Error bounds certified with VCFloat
allow the use of energy-efficient approximate
implementations in critical applications

Ready for use in research and industry

- NSF DeepSpec
- End-to-end verification of cyber-physical systems

# Proofs, Code and Data Rights



INRIA Flocq — LGPLv3

INRIA CompCert — GPLv2+

MSR-INRIA Ssreflect + MathComp — CeCILL-B

INRIA Coq-Interval — CeCILL-C

PNNL-GTRI PERFECT Data Suite — BSD

Reservoir Labs Adaptive SAR — Patent Pending

GPLv3

Patched version CeCILL-C

GPLv3

Reservoir Labs VCFloat GPLv3

GPLv3
Reservoir Labs `Verified SAR (no approx.)

Reservoir Labs Verified adaptive SAR

BSD-like
LGPL-like
GPL-like
Non-free

# Thank you!

Coq library and proofs available online:

- http://github.com/reservoirlabs/vcfloat
- Stress test for Flocq, Coq-Interval, and computations within Coq
- Conference paper published at ACM/SIGPLAN CPP 2016
  - Ramananandro, Mountcastle, Meister, Lethin.
    *A Unified Coq Framework for Verifying C Programs with Floating-Point Computations*

For more information:

- ramananandro@reservoir.com
- http://www.reservoir.com/wp-publications/ramananandro2016