# On Automatically Proving the Correctness of `math.h` Implementations

**Wonyeol Lee**[1,2]        Rahul Sharma[3]        Alex Aiken[1]

[1] Stanford University        [2] KAIST        [3] Microsoft Research

# Our Goal

$$log\ x$$

mathematical
specification

# Our Goal

$$log\ x$$

mathematical
specification

```
<log>
...
mulpd   %xmm2,  %xmm6
addsd   %xmm0,  %xmm5
addsd   %xmm7,  %xmm3
addsd   %xmm5,  %xmm1
...
```

`math.h` implementation

# Our Goal

$$log\ x$$

mathematical
specification

```
<log>
...
mulpd   %xmm2,  %xmm6
addsd   %xmm0,  %xmm5
addsd   %xmm7,  %xmm3
addsd   %xmm5,  %xmm1
...
```

`math.h` implementation

# Our Goal

infinite bits

fixed bits

$$log\ x$$

mathematical specification

```
<log>
...
mulpd  %xmm2, %xmm6
addsd  %xmm0, %xmm5
addsd  %xmm7, %xmm3
addsd  %xmm5, %xmm1
...
```

`math.h` implementation

# Our Goal

infinite bits

$$log\ x$$

mathematical
specification

$\neq$

```
<log>
...
mulpd   %xmm2, %xmm6
addsd   %xmm0, %xmm5
addsd   %xmm7, %xmm3
addsd   %xmm5, %xmm1
...
```
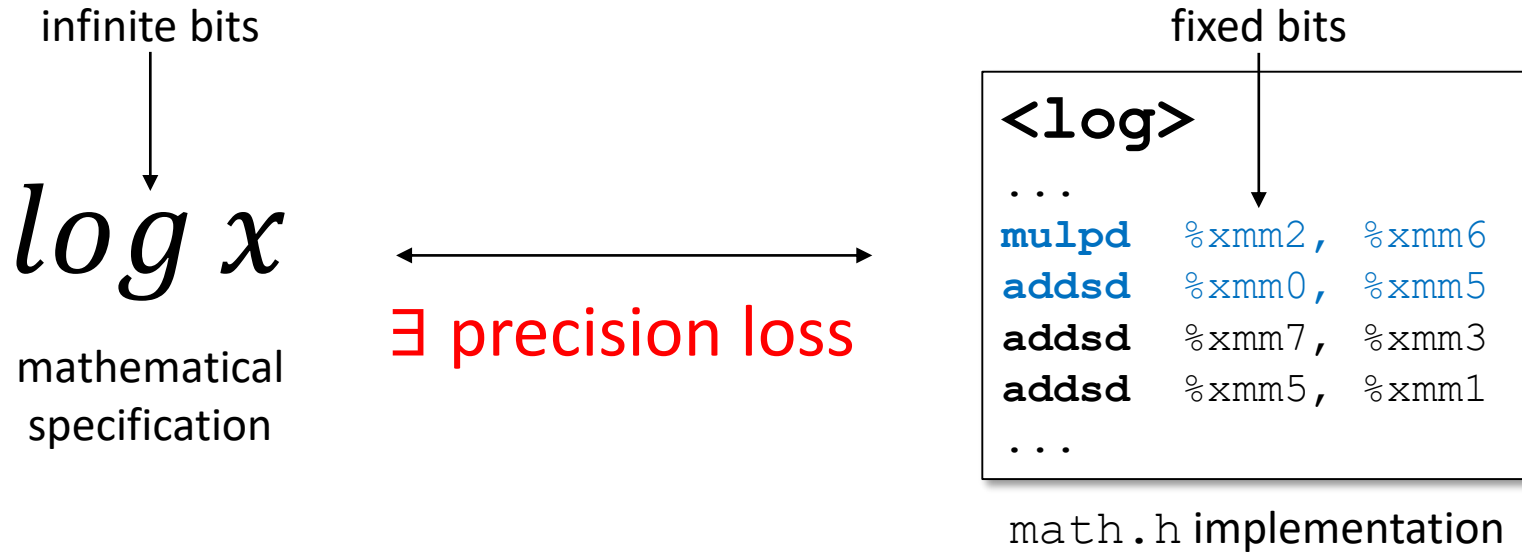
fixed bits

`math.h` implementation

# Our Goal

infinite bits

$$log\ x$$

mathematical
specification

∃ precision loss

fixed bits

```
<log>
...
mulpd   %xmm2, %xmm6
addsd   %xmm0, %xmm5
addsd   %xmm7, %xmm3
addsd   %xmm5, %xmm1
...
```

`math.h` implementation

# Our Goal

infinite bits

fixed bits

$$log\ x$$

mathematical
specification

∃ precision loss

```
<log>
...
mulpd   %xmm2,  %xmm6
addsd   %xmm0,  %xmm5
addsd   %xmm7,  %xmm3
addsd   %xmm5,  %xmm1
...
```
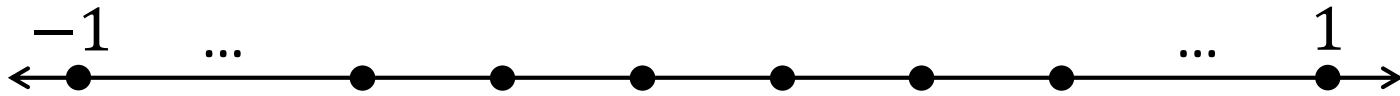
`math.h` implementation
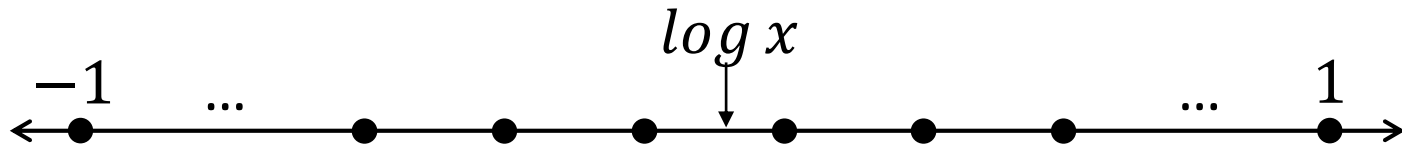
Industry standard implementations of `math.h` claim:

"precision loss is less than 1 ulp"

# Our Goal



Industry standard implementations of `math.h` claim:

"precision loss is less than 1 ulp"

# Our Goal
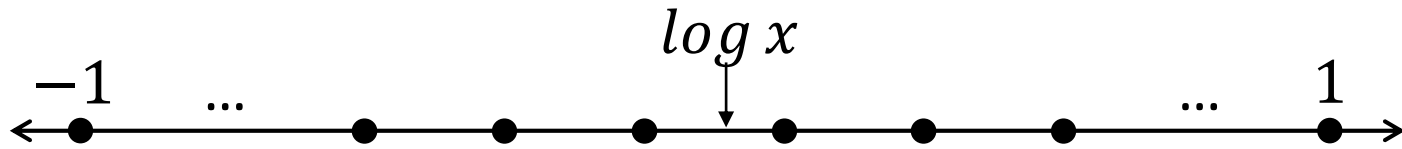
$$log \; x$$

$-1$    ...    •    •    •    ↓    •    •    •    ...    $1$

Industry standard implementations of `math.h` claim:

"precision loss is less than 1 ulp"

# Our Goal

**log** has precision loss of <span style="color:red">< 1 ulp</span>:
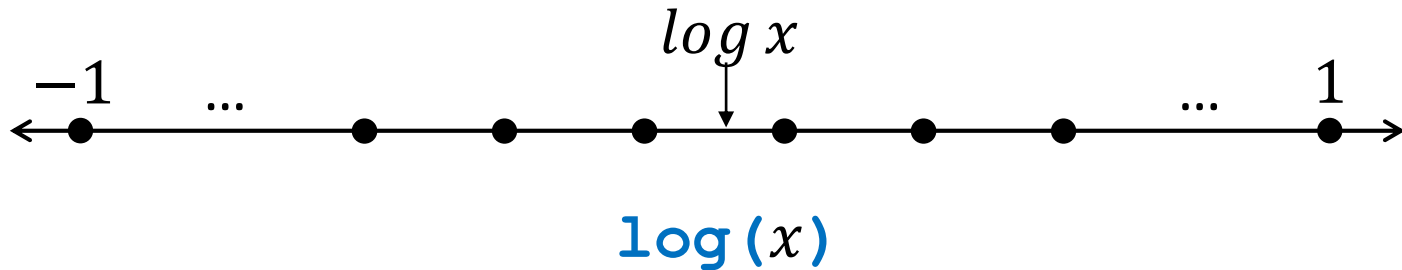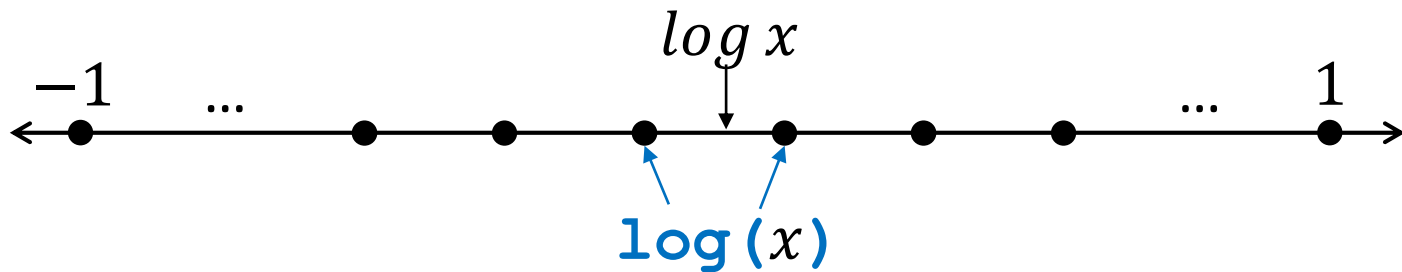


Industry standard implementations of `math.h` claim:

"precision loss is less than <span style="color:red">1 ulp</span>"

# Our Goal

**log** has precision loss of <span style="color:red">< 1 ulp</span>:



Industry standard implementations of `math.h` claim:

"precision loss is less than <span style="color:red">1 ulp</span>"

# Our Goal

**log** has precision loss of < 1 ulp:



Industry standard implementations of `math.h` claim:

"precision loss is less than 1 ulp"

# Our Goal

**log** has precision loss of < 1 ulp:



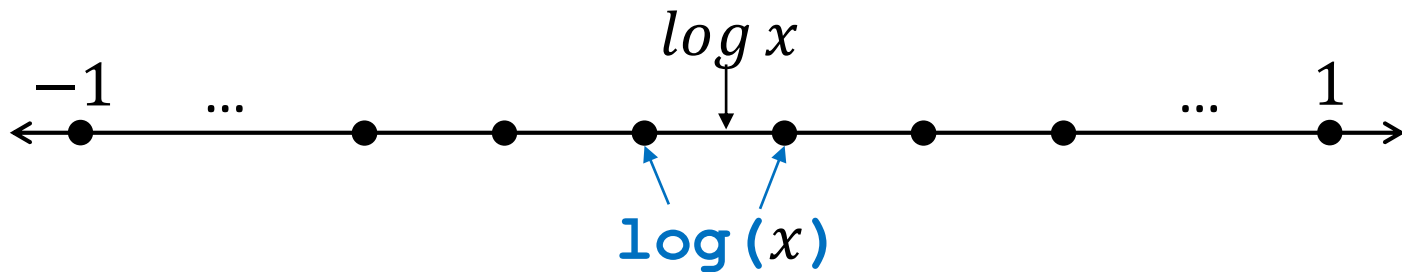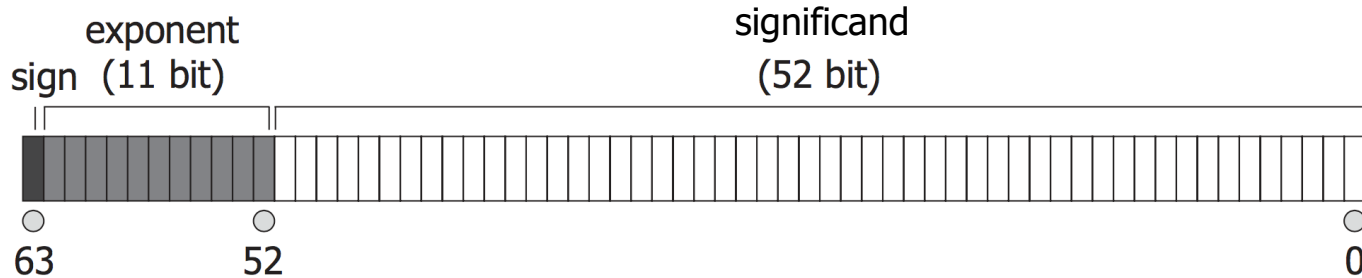Goal: Prove this claim automatically!

Industry standard implementations of `math.h` claim:

"precision loss is less than 1 ulp"

# Floating-Point Numbers/Operations



sign | exponent (11 bit) | significand (52 bit)

63    52    0

- Example:

$$1 \quad 01111111111 \quad 1100\cdots00_{(2)}$$

$$= (-1)^{1} \cdot 2^{1023 - 1023} \cdot 1.110\cdots00_{(2)}$$

# Floating-Point Numbers/Operations



- Example:

$$1 \quad 01111111111 \quad 1100\cdots00 \,_{(2)}$$

$$= (-1)^{1} \cdot 2^{1023-1023} \cdot 1.110\cdots00_{(2)}$$

- Rounding errors

# Floating-Point Numbers/Operations



- Example:

$$= (-1)^1 \cdot 2^{1023 - 1023} \cdot 1.110\cdots00_{(2)}$$

- Rounding errors
  - Floating-point arithmetic $\neq$ Real arithmetic

floating-point operation

$$1 \oplus (2^{100} \ominus 2^{100}) = 1 \neq 0 = (1 \oplus 2^{100}) \ominus 2^{100}$$

# Floating-Point Numbers/Operations

sign | exponent (11 bit) | significand (52 bit)

63      52      0

- Example: 

$$\underbrace{1}_{} \; \underbrace{01111111111}_{} \; \underbrace{1100\cdots00}_{} {}_{(2)}$$

$$= (-1)^{1} \; \cdot \; 2^{1023-1023} \; \cdot \; 1.110\cdots00_{(2)}$$

- Rounding errors

  - Floating-point arithmetic ≠ Real arithmetic

    floating-point operation

$$1 \oplus (2^{100} \ominus 2^{100}) = 1 \neq 0 = (1 \oplus 2^{100}) \ominus 2^{100}$$

  - Floating-point implementations often have precision loss

# Ulp Error

- Typically used to measure accuracy of numeric libraries

# Ulp Error

- Typically used to measure accuracy of numeric libraries

- Definition:



- $a =$ exact value, $b =$ approximate value

# Ulp Error

- Typically used to measure accuracy of numeric libraries

- Definition:



- $a =$ exact value, $b =$ approximate value
- $\text{ulp}(a) =$ gap between two doubles that surround $a$

# Ulp Error

- Typically used to measure accuracy of numeric libraries

- Definition:



- $a$ = exact value, $b$ = approximate value
- $\mathrm{ulp}(a)$ = gap between two doubles that surround $a$

$$\mathrm{ErrUlp}(a, b) = \frac{|a - b|}{\mathrm{ulp}(a)}$$

# Ulp Error

- Typically used to measure accuracy of numeric libraries

- Definition:



- $a =$ exact value, $b =$ approximate value
- $\mathrm{ulp}(a) =$ gap between two doubles that surround $a$

$$\mathrm{ErrUlp}(a, b) = \frac{|a - b|}{\mathrm{ulp}(a)}$$

# Problem Statement

$$log\ x$$

mathematical
specification $f$

```
<log>
...
mulpd   %xmm2, %xmm6
addsd   %xmm0, %xmm5
addsd   %xmm7, %xmm3
addsd   %xmm5, %xmm1
...
```

math.h implementation $P$

# Problem Statement

$$log\ x$$

mathematical
specification $f$

$(0, 2^{1024})$

input interval $X$

```
<log>
...
mulpd   %xmm2,  %xmm6
addsd   %xmm0,  %xmm5
addsd   %xmm7,  %xmm3
addsd   %xmm5,  %xmm1
...
```

`math.h` implementation $P$

# Problem Statement

$$log\ x$$

mathematical
specification $f$

$\xleftrightarrow{\quad (0, 2^{1024}) \quad}$

input interval $X$

```
<log>
...
mulpd   %xmm2,  %xmm6
addsd   %xmm0,  %xmm5
addsd   %xmm7,  %xmm3
addsd   %xmm5,  %xmm1
...
```

`math.h` implementation $P$

- Problem: Find a small $\Theta > 0$ automatically such that

$$\text{ErrUlp}(f(x), P(x)) \leq \Theta \quad \text{for all } x \in X$$

# Problem Statement

$$log\ x$$

mathematical
specification $f$

$(0, 2^{1024})$

input interval $X$

```
<log>
...
mulpd   %xmm2,  %xmm6
addsd   %xmm0,  %xmm5
addsd   %xmm7,  %xmm3
addsd   %xmm5,  %xmm1
...
```

`math.h` implementation $P$

- Problem: Find a small $\Theta > 0$ automatically such that

$$\mathrm{ErrUlp}(f(x), P(x)) \leq \Theta \quad \text{for all } x \in X$$

i.e., prove a bound on the maximum precision loss

# Problem Statement

$$log\ x$$

mathematical
specification $f$

$(0, 2^{1024})$
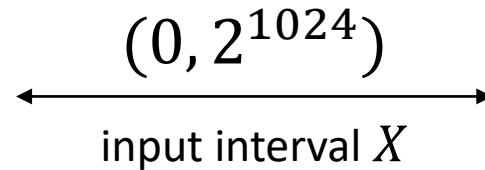
input interval $X$

```
<log>
...
mulpd   %xmm2,  %xmm6
addsd   %xmm0,  %xmm5
addsd   %xmm7,  %xmm3
addsd   %xmm5,  %xmm1
...
```

`math.h` implementation $P$

- Problem: Find a small $\Theta > 0$ automatically such that

$$\mathrm{ErrUlp}(f(x), P(x)) \leq \Theta \quad \text{for all } x \in X$$

i.e., prove a bound on the maximum precision loss

Goal: want to find $\Theta < 1$

# Previous Work

- Machine-checkable proofs
  - Harrison used HOL Light to prove [FMCAD'00]:

    "Intel's `sin` has precision loss of $< 0.574$ ulps"

# Previous Work

- Machine-checkable proofs
  - Harrison used HOL Light to prove [FMCAD'00]:

    "Intel's `sin` has precision loss of $< 0.574$ ulps"

  - "Construction of these proofs often requires considerable persistence." [FMSD'00]

# Previous Work

- Machine-checkable proofs
  - Harrison used HOL Light to prove [FMCAD'00]:

    "Intel's `sin` has precision loss of $< 0.574$ ulps"

  - "Construction of these proofs often requires considerable persistence." [FMSD'00]

- Automatic techniques
  - Astree [PLDI'03], Fluctuat [FMICS'09], Gappa [TOMS'10], MathSAT [FMCAD'12], Rosa [POPL'14], FPTaylor [FM'15], Lee/Sharma/Aiken [PLDI'16], $\cdots$

# Previous Work

- Machine-checkable proofs
  - Harrison used HOL Light to prove [FMCAD'00]:

    "Intel's `sin` has precision loss of $< 0.574$ ulps"

  - "Construction of these proofs often requires considerable persistence." [FMSD'00]

- Automatic techniques
  - Astree [PLDI'03], Fluctuat [FMICS'09], Gappa [TOMS'10], MathSAT [FMCAD'12], Rosa [POPL'14], FPTaylor [FM'15], Lee/Sharma/Aiken [PLDI'16], ⋯
  - None of them can prove $< 1$ ulp error bound automatically

# $(1 + \epsilon)$ Property

- A standard way to model rounding errors

# $(1 + \epsilon)$ Property

- A standard way to model rounding errors

- <u>Theorem</u>  For any double $a$ and $b$, and $* \in \{+, -, \times, /\}$,

$$a \circledast b = (a * b)(1 + \delta) \quad \text{for some } |\delta| < \epsilon$$

where $\epsilon = 2^{-53}$

# $(1 + \epsilon)$ Property

- A standard way to model rounding errors

- <u>Theorem</u>  For any double $a$ and $b$, and $* \in \{+, -, \times, /\}$,

$$a \circledast b = (a * b)(1 + \delta) \quad \text{for some } |\delta| < \epsilon$$

where $\epsilon = 2^{-53}$

That is,



$a \circledast b$

# $(1 + \epsilon)$ Property

- A standard way to model rounding errors

- <u>Theorem</u>  For any double $a$ and $b$, and $* \in \{+, -, \times, /\}$,

$$a \circledast b = (a * b)(1 + \delta) \quad \text{for some } |\delta| < \epsilon$$

  where $\epsilon = 2^{-53}$

  That is,

# $(1 + \epsilon)$ Property

- A standard way to model rounding errors

- <u>Theorem</u>  For any double $a$ and $b$, and $* \in \{+, -, \times, /\}$,

$$a \circledast b = (a * b)(1 + \delta) \quad \text{for some } |\delta| < \epsilon$$

  where $\epsilon = 2^{-53}$

  That is,

  $$\{(a * b)(1 + \delta) : |\delta| < \epsilon\}$$



$a * b$

$a \circledast b$

0                                                                    1

# $(1 + \epsilon)$ Property

- A standard way to model rounding errors

- <u>Theorem</u>  For any double $a$ and $b$, and $* \in \{+, -, \times, /\}$,

$$a \circledast b = (a * b)(1 + \delta) \quad \text{for some } |\delta| < \epsilon$$
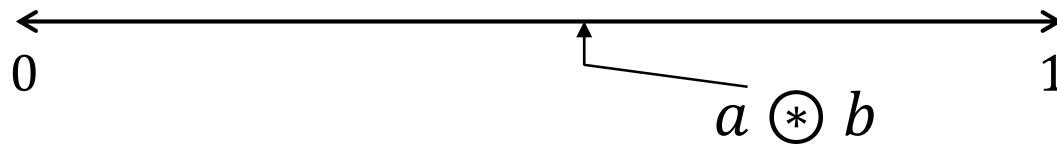
  where $\epsilon = 2^{-53}$

  That is,

  $$\{(a * b)(1 + \delta) : |\delta| < \epsilon\}$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1+x} \quad \text{over} \quad X = [-0.1, 0.1]$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1+x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \longleftarrow \quad \text{implementation of } f(x)$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1 + x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \longleftarrow \quad \text{implementation of } f(x)$$

1. Construct abstraction $A_{\vec{\delta}}$ of $e$:

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1+x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \longleftarrow \quad \text{implementation of } f(x)$$

1. Construct abstraction $A_{\vec{\delta}}$ of $e$:

$$A_{\vec{\delta}}(x) = 1 \oplus (0.5 \otimes x)$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1 + x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \longleftarrow \quad \text{implementation of } f(x)$$

1. Construct abstraction $A_{\vec{\delta}}$ of $e$:

$$A_{\vec{\delta}}(x) = 1 + (0.5 \times x)$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1+x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \longleftarrow \quad \text{implementation of } f(x)$$

1. Construct abstraction $A_{\vec{\delta}}$ of $e$:

$$A_{\vec{\delta}}(x) = \big(1 + (0.5 \times x)(1 + \delta_1)\big)(1 + \delta_2)$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1 + x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \longleftarrow \quad \text{implementation of } f(x)$$

1. Construct abstraction $A_{\vec{\delta}}$ of $e$:

$$A_{\vec{\delta}}(x) = \big(1 + (0.5 \times x)(1 + \delta_1)\big)(1 + \delta_2)$$

   - $A_{\vec{\delta}}$ is a sound abstraction of $e$ (or $e \sqsubseteq A_{\vec{\delta}}$):

$$e(x) \in \{A_{\vec{\delta}}(x) : |\delta_1|, |\delta_2| < \epsilon\} \quad \text{for all } x \in X$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1+x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \longleftarrow \quad \text{implementation of } f(x)$$

2. Compute a bound on relative error of $e$:

$$\Theta_{rel} = \max_{x \in [-0.1, 0.1], \, \delta_1, \delta_2 \in [-\epsilon, \epsilon]} \left| \frac{f(x) - A_{\vec{\delta}}(x)}{f(x)} \right|$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1+x} \quad \text{over} \quad X = [-0.1, 0.1]$$
$$e = 1 \oplus (0.5 \otimes x) \quad \leftarrow \quad \text{implementation of } f(x)$$

2. Compute a bound on relative error of $e$:

$$\Theta_{rel} = \max_{x \in [-0.1, 0.1], \, \delta_1, \delta_2 \in [-\epsilon, \epsilon]} \left| \frac{f(x) - A_{\vec{\delta}}(x)}{f(x)} \right|$$

$$= \max_{x \in [-0.1, 0.1], \, \delta_1, \delta_2 \in [-\epsilon, \epsilon]} \left| \frac{\sqrt{1+x} - \left(1 + (0.5 \times x)(1+\delta_1)\right)(1+\delta_2)}{\sqrt{1+x}} \right|$$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1 + x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \leftarrow \quad \text{implementation of } f(x)$$

2. Compute a bound on relative error of $e$:

$$\Theta_{rel} = \max_{x \in [-0.1, 0.1], \, \delta_1, \delta_2 \in [-\epsilon, \epsilon]} \left| \frac{f(x) - A_{\vec{\delta}}(x)}{f(x)} \right|$$

$$= \max_{x \in [-0.1, 0.1], \, \delta_1, \delta_2 \in [-\epsilon, \epsilon]} \left| \frac{\sqrt{1+x} - \left(1 + (0.5 \times x)(1 + \delta_1)\right)(1 + \delta_2)}{\sqrt{1+x}} \right|$$

---

Theorem  $\text{ErrUlp}(a, b) \leq \text{ErrRel}(a, b) / \epsilon$

# Standard Error Analysis

- Example:

$$f(x) = \sqrt{1 + x} \quad \text{over} \quad X = [-0.1, 0.1]$$

$$e = 1 \oplus (0.5 \otimes x) \quad \longleftarrow \quad \text{implementation of } f(x)$$

2. Compute a bound on relative error of $e$:

$$\Theta_{rel} = \max_{x \in [-0.1, 0.1],\ \delta_1, \delta_2 \in [-\epsilon, \epsilon]} \left| \frac{f(x) - A_{\vec{\delta}}(x)}{f(x)} \right|$$

$$= \max_{x \in [-0.1, 0.1],\ \delta_1, \delta_2 \in [-\epsilon, \epsilon]} \left| \frac{\sqrt{1+x} - \left(1 + (0.5 \times x)(1 + \delta_1)\right)(1 + \delta_2)}{\sqrt{1+x}} \right|$$

3. Compute a bound on ulp error of $e$:

$$\Theta_{ulp} = \Theta_{rel} / \epsilon$$

$$\boxed{\underline{\text{Theorem}} \quad \text{ErrUlp}(a, b) \leq \text{ErrRel}(a, b) / \epsilon}$$

# Standard Error Analysis: Limitation

# Standard Error Analysis: Limitation

Intel's `log`

$10^{14}$



ulp error

input value

based on standard
error analysis

error bounds from [PLDI'16]

# Standard Error Analysis: Limitation

Intel's `log`

$10^{14}$

ulp error

3 ulps

input value

based on standard
error analysis

—— error bounds from [PLDI'16]
- - - 1 ulp

# Standard Error Analysis: Limitation

Intel's `log`

$10^{14}$



ulp error

$\left[ \dfrac{4095}{4096}, 1 \right)$

3 ulps

input value

based on standard
error analysis

— error bounds from [PLDI'16]

--- 1 ulp

# Why the Loose Error Bound?

# Why the Loose Error Bound?

- Floating-point operations are sometimes **exact**

# Why the Loose Error Bound?

- Floating-point operations are sometimes **exact**

$$0.125 \otimes x \quad = \quad 0.125 \times x \qquad \text{if } |0.125x| \geq 2^{-1022}$$

$$x \ominus 1 \quad = \quad x - 1 \qquad \text{if } 0.5 \leq x \leq 2$$

# Why the Loose Error Bound?

- Floating-point operations are sometimes **exact**

$$0.125 \otimes x \quad = \quad 0.125 \times x \qquad \text{if } |0.125x| \geq 2^{-1022}$$

$$x \ominus 1 \quad = \quad x - 1 \qquad \text{if } 0.5 \leq x \leq 2$$

- Standard error analysis constructs imprecise abstractions

# Why the Loose Error Bound?

- Floating-point operations are sometimes **exact**

$$0.125 \otimes x \quad = \quad 0.125 \times x \qquad \text{if } |0.125x| \geq 2^{-1022}$$
$$x \ominus 1 \quad = \quad x - 1 \qquad \text{if } 0.5 \leq x \leq 2$$

- Standard error analysis constructs <span style="color:red">imprecise</span> abstractions

$$0.125 \otimes x \qquad\qquad\qquad\qquad \text{if } |0.125x| \geq 2^{-1022}$$

$$x \ominus 1 \qquad\qquad\qquad\qquad \text{if } 0.5 \leq x \leq 2$$

# Why the Loose Error Bound?

- Floating-point operations are sometimes **exact**

$$0.125 \otimes x \quad = \quad 0.125 \times x \qquad \text{if } |0.125x| \geq 2^{-1022}$$

$$x \ominus 1 \quad = \quad x - 1 \qquad \text{if } 0.5 \leq x \leq 2$$

- Standard error analysis constructs <span style="color:red">imprecise</span> abstractions

$$0.125 \otimes x \quad \sqsubseteq \quad 0.125 \times x \qquad \text{if } |0.125x| \geq 2^{-1022}$$

$$x \ominus 1 \quad \sqsubseteq \quad x - 1 \qquad \text{if } 0.5 \leq x \leq 2$$

# Why the Loose Error Bound?

- Floating-point operations are sometimes **exact**

$$0.125 \otimes x \quad = \quad 0.125 \times x \qquad \text{if } |0.125x| \geq 2^{-1022}$$
$$x \ominus 1 \quad = \quad x - 1 \qquad \text{if } 0.5 \leq x \leq 2$$

- Standard error analysis constructs <span style="color:red">imprecise</span> abstractions

$$0.125 \otimes x \quad \sqsubseteq \quad 0.125 \times x \qquad \text{if } |0.125x| \geq 2^{-1022}$$
$$0.125 \otimes x \quad \sqsubseteq \quad (0.125 \times x)(1 + \delta)$$

$$x \ominus 1 \quad \sqsubseteq \quad x - 1 \qquad \text{if } 0.5 \leq x \leq 2$$
$$x \ominus 1 \quad \sqsubseteq \quad (x - 1)(1 + \delta')$$

# Analysis of `log`

- For $x \in \left[\frac{4095}{4096}, 1\right)$, `log` computes

# Analysis of `log`

- For $x \in \left[\frac{4095}{4096}, 1\right),$ `log` computes

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad (\approx x - 1)$$

and returns

$$r(x) - \frac{1}{2}r(x)^2 + \cdots + \frac{1}{7}r(x)^7$$

# Analysis of `log`

- For $x \in \left[\frac{4095}{4096}, 1\right)$, `log` computes

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad (\approx x - 1)$$

and returns

$$\boxed{r(x)} - \frac{1}{2}r(x)^2 + \cdots + \frac{1}{7}r(x)^7$$

# Analysis of `log`

- For $x \in \left[\frac{4095}{4096}, 1\right)$, `log` computes

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad (\approx x - 1)$$

and returns

$$\boxed{r(x)} - \frac{1}{2}r(x)^2 + \cdots + \frac{1}{7}r(x)^7$$

- Roughly speaking,

$$(\text{precision loss of } \texttt{log}) \geq (\text{precision loss of } r(x))$$

# Analysis of `log`

- For $x \in \left[ \frac{4095}{4096}, 1 \right)$, `log` computes

$$r(x) = \left[ \left( (2 \otimes x) \ominus \frac{255}{128} \right) \otimes \frac{1}{2} \right] \oplus \left[ \left( \frac{255}{128} \otimes \frac{1}{2} \right) \ominus 1 \right] \quad (\approx x - 1)$$

and returns

$$\boxed{r(x)} - \frac{1}{2} r(x)^2 + \cdots + \frac{1}{7} r(x)^7$$

- Roughly speaking,

$$(\text{precision loss of } \texttt{log}) \geq \boxed{(\text{precision loss of } r(x))}$$

# Standard Analysis of `log`

$$r(x) = \left[ \left( (2 \otimes x) \ominus \frac{255}{128} \right) \otimes \frac{1}{2} \right] \oplus \left[ \left( \frac{255}{128} \otimes \frac{1}{2} \right) \ominus 1 \right] \quad \left( \frac{4095}{4096} \leq x < 1 \right)$$

- Abstraction of $r(x)$:

$$A_{\vec{\delta}}(x) = \left[ \left( (2 \times x)(1 + \delta_0) - \frac{255}{128} \right)(1 + \delta_1) \times \frac{1}{2} \right](1 + \delta_2) + \cdots$$

# Standard Analysis of `log`

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

- Abstraction of $r(x)$:

$$A_{\vec{\delta}}(x) = \left[\left((2 \times x)(1 + \delta_0) - \frac{255}{128}\right)(1 + \delta_1) \times \frac{1}{2}\right](1 + \delta_2) + \cdots$$

# Standard Analysis of `log`

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

- Abstraction of $r(x)$:

$$A_{\vec{\delta}}(x) = \left[\left((2 \times x)(1 + \delta_0) - \frac{255}{128}\right)(1 + \delta_1) \times \frac{1}{2}\right](1 + \delta_2) + \cdots$$

$$= (x - 1) + \left(x - \frac{255}{256}\right)\delta_1 + \cdots$$

# Standard Analysis of `log`

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \le x < 1\right)$$

- Abstraction of $r(x)$:

$$A_{\vec{\delta}}(x) = \left[\left((2 \times x)(1 + \delta_0) - \frac{255}{128}\right)(1 + \delta_1) \times \frac{1}{2}\right](1 + \delta_2) + \cdots$$

$$= (x - 1) + \left(x - \frac{255}{256}\right)\delta_1 + \cdots$$

- Bound on relative error of $r(x)$:

$$\max_{\frac{4095}{4096} \le x < 1, \, |\delta_i| < \epsilon} \left|\frac{A_{\vec{\delta}}(x) - (x - 1)}{x - 1}\right|$$

# Standard Analysis of `log`

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

- Abstraction of $r(x)$:

$$A_{\vec{\delta}}(x) = \left[\left((2 \times x)(1 + \delta_0) - \frac{255}{128}\right)(1 + \delta_1) \times \frac{1}{2}\right](1 + \delta_2) + \cdots$$

$$= (x - 1) + \boxed{\left(x - \frac{255}{256}\right)\delta_1} + \cdots$$

- Bound on relative error of $r(x)$:

$$\max_{\frac{4095}{4096} \leq x < 1,\ |\delta_i| < \epsilon} \left|\frac{A_{\vec{\delta}}(x) - (x - 1)}{x - 1}\right|$$

# Standard Analysis of `log`

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \le x < 1\right)$$

- Abstraction of $r(x)$:

$$A_{\vec{\delta}}(x) = \left[\left((2 \times x)(1 + \delta_0) - \frac{255}{128}\right)(1 + \delta_1) \times \frac{1}{2}\right](1 + \delta_2) + \cdots$$

$$= (x - 1) + \boxed{\left(x - \frac{255}{256}\right)\delta_1} + \cdots$$

- Bound on relative error of $r(x)$:

$$\max_{\frac{4095}{4096} \le x < 1, \, |\delta_i| < \epsilon} \left|\frac{A_{\vec{\delta}}(x) - (x - 1)}{x - 1}\right| \ge \max_{\frac{4095}{4096} \le x < 1} \left|\frac{x - \frac{255}{256}}{x - 1}\right| \epsilon$$

# Standard Analysis of `log`

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \le x < 1\right)$$

- Abstraction of $r(x)$:

$$A_{\vec{\delta}}(x) = \left[\left((2 \times x)(1 + \delta_0) - \frac{255}{128}\right)(1 + \delta_1) \times \frac{1}{2}\right](1 + \delta_2) + \cdots$$

$$= (x - 1) + \boxed{\left(x - \frac{255}{256}\right)\delta_1} + \cdots$$

- Bound on relative error of $r(x)$:

$$\max_{\frac{4095}{4096} \le x < 1,\ |\delta_i| < \epsilon} \left|\frac{A_{\vec{\delta}}(x) - (x - 1)}{x - 1}\right| \ge \max_{\frac{4095}{4096} \le x < 1} \boxed{\left|\frac{\left|x - \frac{255}{256}\right|}{x - 1}\right|\epsilon}$$

unbounded

# Standard Analysis of `log`

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

- Abstraction of $r(x)$:

$$A_{\vec{\delta}}(x) = \left[\left((2 \times x)(1 + \delta_0) - \frac{255}{128}\right)(1 + \delta_1) \times \frac{1}{2}\right](1 + \delta_2) + \cdots$$

$$= (x - 1) + \boxed{\left(x - \frac{255}{256}\right)\delta_1} + \cdots$$

$\boxed{10^{14} \text{ ulps for } \texttt{log} \text{ near } x = 1 \text{ [PLDI'16]}}$

- Bound on relative error of $r(x)$:

$$\max_{\frac{4095}{4096} \leq x < 1, |\delta_i| < \epsilon} \left|\frac{A_{\vec{\delta}}(x) - (x - 1)}{x - 1}\right| \geq \max_{\frac{4095}{4096} \leq x < 1} \boxed{\left|\frac{x - \frac{255}{256}}{x - 1}\right| \epsilon}$$

unbounded

# Precise Analysis of `log`

$$r(x) = \left[\left((2 \boxed{\otimes} x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \le x < 1\right)$$

# Precise Analysis of `log`

$$r(x) = \left[\left((2 \boxed{\otimes} x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

exact

# Precise Analysis of $\log$

exact

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

# Precise Analysis of `log`

exact

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

# Precise Analysis of `log`

exact

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \le x < 1\right)$$

# Precise Analysis of `log`

exact

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

- More precise abstraction of $r(x)$:

$$A'_{\vec{\delta}}(x) = (x - 1) + (x - 1)\delta'$$

# Precise Analysis of `log`

exact

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

- More precise abstraction of $r(x)$:

$$A'_{\vec{\delta}}(x) = (x - 1) + (x - 1)\delta'$$

- Tighter bound on relative error of $r(x)$:

$$\max_{\frac{4095}{4096} \leq x < 1, \, |\delta'| < \epsilon} \left|\frac{A'_{\vec{\delta}}(x) - (x - 1)}{x - 1}\right| = \max_{\frac{4095}{4096} \leq x < 1} \left|\frac{x - 1}{x - 1}\right| \epsilon$$

$$= \epsilon$$

# Precise Analysis of `log`

exact

$$r(x) = \left[ \left( (2 \otimes x) \ominus \frac{255}{128} \right) \otimes \frac{1}{2} \right] \oplus \left[ \left( \frac{255}{128} \otimes \frac{1}{2} \right) \ominus 1 \right] \quad \left( \frac{4095}{4096} \leq x < 1 \right)$$

- More precise abstraction of $r(x)$:

$$A'_{\vec{\delta}}(x) = (x - 1) + (x - 1)\delta'$$

- Tighter bound on relative error of $r(x)$:

$$\max_{\frac{4095}{4096} \leq x < 1, \, |\delta'| < \epsilon} \left| \frac{A'_{\vec{\delta}}(x) - (x - 1)}{x - 1} \right| = \max_{\frac{4095}{4096} \leq x < 1} \left| \frac{x - 1}{x - 1} \right| \epsilon = \epsilon$$

We prove a bound of
0.583 ulps for `log`

81/146

# Precise Analysis of `log`

exact

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \leq x < 1\right)$$

- To prove the 1 ulp error bound,

  - Need to construct more precise abstractions

-

$\frac{4095}{4096} \leq x < 1, |\delta\prime| < \epsilon$     $x - 1$     $\frac{4095}{4096} \leq x < 1$   $|x - 1|$

We prove a bound of
0.583 ulps for `log`     $= \quad \epsilon$

# Precise Analysis of `log`

exact

$$r(x) = \left[\left((2 \otimes x) \ominus \frac{255}{128}\right) \otimes \frac{1}{2}\right] \oplus \left[\left(\frac{255}{128} \otimes \frac{1}{2}\right) \ominus 1\right] \quad \left(\frac{4095}{4096} \le x < 1\right)$$

- 

To prove the 1 ulp error bound,

- Need to construct more precise abstractions

- Need to use floating-point exactness results

$\frac{4095}{4096} \le x < 1, |\delta\prime| < \epsilon$   $x - 1$   $\frac{4095}{4096} \le x < 1$   $|x - 1|$

$= \epsilon$

We prove a bound of
0.583 ulps for `log`

# Sterbenz's Theorem

- <u>Theorem</u> [Sterbenz, 1973]

$$\frac{1}{2}a \leq b \leq 2a \quad \Longrightarrow \quad a \ominus b = a - b$$

# Sterbenz's Theorem

- <u>Theorem</u> [Sterbenz, 1973]

$$\frac{1}{2}a \leq b \leq 2a \quad \Longrightarrow \quad a \ominus b = a - b$$

- Example: `log` for $x \in \left[\frac{4095}{4096}, 1\right)$

# Sterbenz's Theorem

- <u>Theorem</u> [Sterbenz, 1973]

$$\frac{1}{2}a \leq b \leq 2a \quad \Longrightarrow \quad a \ominus b = a - b$$

- Example: `log` for $x \in \left[\frac{4095}{4096}, 1\right)$

- Example: `sin` for $x \in \left[2\pi - \frac{\pi}{64}, 2\pi + \frac{\pi}{64}\right]$

  1. Compute $y = x \ominus 2\pi$
  2. Return $y - \frac{1}{3!}y^3 + \cdots + \frac{1}{9!}y^9 \ (\approx \sin y)$

# Sterbenz's Theorem

- <u>Theorem</u> [Sterbenz, 1973]

$$\frac{1}{2}a \le b \le 2a \quad \Longrightarrow \quad a \ominus b = a - b$$

- Example: `log` for $x \in \left[\frac{4095}{4096}, 1\right)$

- Example: `sin` for $x \in \left[2\pi - \frac{\pi}{64}, 2\pi + \frac{\pi}{64}\right]$

  1. Compute $y = x \ominus 2\pi$
  2. Return $y - \frac{1}{3!}y^3 + \cdots + \frac{1}{9!}y^9 \ (\approx \sin y)$

# Sterbenz's Theorem

- <u>Theorem</u> [Sterbenz, 1973]

$$\frac{1}{2}a \leq b \leq 2a \quad \implies \quad a \ominus b = a - b$$

- Example: `log` for $x \in \left[\frac{4095}{4096}, 1\right)$

- Example: `sin` for $x \in \left[2\pi - \frac{\pi}{64}, 2\pi + \frac{\pi}{64}\right]$

  1. Compute $y = x \ominus 2\pi$  exact
  2. Return $y - \frac{1}{3!}y^3 + \cdots + \frac{1}{9!}y^9 \ (\approx sin\,y)$

# Sterbenz's Theorem

- <u>Theorem</u> [Sterbenz, 1973]

$$\frac{1}{2}a \leq b \leq 2a \quad \Longrightarrow \quad a \ominus b = a - b$$

- Example: `log` for $x \in \left[\frac{4095}{4096}, 1\right)$

- Example: `sin` for $x \in \left[2\pi - \frac{\pi}{64}, 2\pi + \frac{\pi}{64}\right]$

  1. Compute $y = x \ominus 2\pi$   exact
  2. Return $y - \frac{1}{3!}y^3 + \cdots + \frac{1}{9!}y^9 \ (\approx sin\ y)$

precision loss of `sin` is small

# Sterbenz's Theorem

- How to apply the theorem automatically?

# Sterbenz's Theorem

- How to apply the theorem automatically?

$$e \ominus e' \quad \sqsubseteq \quad ?$$

# Sterbenz's Theorem

- How to apply the theorem automatically?

$$e \ominus e' \ \sqsubseteq \ \ ?$$

  - Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}$

# Sterbenz's Theorem

- How to apply the theorem automatically?

$$e \ominus e' \quad \sqsubseteq \quad ?$$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}$

- Precondition of theorem:

$$\frac{1}{2}e(x) \le e'(x) \le 2e(x) \qquad \text{for all } x \in X$$

# Sterbenz's Theorem

- How to apply the theorem automatically?

$$e \ominus e' \ \sqsubseteq \ ?$$

  - Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}$

$$\min_{x, \vec{\delta}} \left( A'_{\vec{\delta}}(x) - \frac{1}{2} A_{\vec{\delta}}(x) \right)$$

$$\max_{x, \vec{\delta}} \left( A'_{\vec{\delta}}(x) - 2 A_{\vec{\delta}}(x) \right)$$

  - Precondition of theorem:

$$\frac{1}{2} e(x) \leq e'(x) \leq 2e(x) \qquad \text{for all } x \in X$$

# Sterbenz's Theorem

- How to apply the theorem automatically?

$$e \ominus e' \sqsubseteq \; ?$$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \; e' \sqsubseteq A'_{\vec{\delta}}$

- Check: $\displaystyle\min_{x, \vec{\delta}} \left( A'_{\vec{\delta}}(x) - \frac{1}{2} A_{\vec{\delta}}(x) \right) \geq 0$ and

$$\max_{x, \vec{\delta}} \left( A'_{\vec{\delta}}(x) - 2 A_{\vec{\delta}}(x) \right) \leq 0$$

- Precondition of theorem:

$$\frac{1}{2} e(x) \leq e'(x) \leq 2 e(x) \qquad \text{for all } x \in X$$

# Sterbenz's Theorem

- How to apply the theorem automatically?

$$e \ominus e' \sqsubseteq \quad ?$$

  - Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}$

  - Check: $\displaystyle\min_{x,\vec{\delta}} \left( A'_{\vec{\delta}}(x) - \frac{1}{2} A_{\vec{\delta}}(x) \right) \geq 0$ and

    $\displaystyle\max_{x,\vec{\delta}} \left( A'_{\vec{\delta}}(x) - 2 A_{\vec{\delta}}(x) \right) \leq 0$ $\longleftarrow$ true

  - Precondition of theorem:

    $$\frac{1}{2} e(x) \leq e'(x) \leq 2 e(x) \qquad \text{for all } x \in X$$

# Sterbenz's Theorem

- How to apply the theorem automatically?

$$e \ominus e' \quad \sqsubseteq \quad ?$$

  - Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}$

  - Check: $\min\limits_{x,\vec{\delta}} \left( A'_{\vec{\delta}}(x) - \frac{1}{2} A_{\vec{\delta}}(x) \right) \geq 0 \quad$ and $\quad\longleftarrow$ true

    $\max\limits_{x,\vec{\delta}} \left( A'_{\vec{\delta}}(x) - 2 A_{\vec{\delta}}(x) \right) \leq 0$

    $\Longrightarrow \quad \frac{1}{2} A_{\vec{\delta}}(x) \leq A'_{\vec{\delta}}(x) \leq 2 A_{\vec{\delta}}(x) \quad$ for all $x \in X$, all $|\delta_i| < \epsilon$

  - Precondition of theorem:

    $$\frac{1}{2} e(x) \leq e'(x) \leq 2 e(x) \qquad \text{for all } x \in X$$

# Sterbenz's Theorem

- How to apply the theorem automatically?

$$e \ominus e' \ \sqsubseteq \ \ ?$$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}$

- Check: $\min\limits_{x,\vec{\delta}} \left( A'_{\vec{\delta}}(x) - \frac{1}{2} A_{\vec{\delta}}(x) \right) \geq 0$ and $\longleftarrow$ true

  $\max\limits_{x,\vec{\delta}} \left( A'_{\vec{\delta}}(x) - 2 A_{\vec{\delta}}(x) \right) \leq 0$

  $\Longrightarrow \quad \frac{1}{2} A_{\vec{\delta}}(x) \leq A'_{\vec{\delta}}(x) \leq 2 A_{\vec{\delta}}(x) \quad$ for all $x \in X$, all $|\delta_i| < \epsilon$

- Precondition of theorem:

  $\Longrightarrow \quad \frac{1}{2} e(x) \leq e'(x) \leq 2 e(x) \qquad$ for all $x \in X$

# Sterbenz's Theorem

- How to apply the theorem automatically?

<span style="color:red">can apply<br>the theorem</span> $\longrightarrow$ $e \ominus e' \ \sqsubseteq \ \ ?$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}$

- Check: $\min\limits_{x,\,\vec{\delta}} \left( A'_{\vec{\delta}}(x) - \frac{1}{2} A_{\vec{\delta}}(x) \right) \geq 0$ and

$\max\limits_{x,\,\vec{\delta}} \left( A'_{\vec{\delta}}(x) - 2 A_{\vec{\delta}}(x) \right) \leq 0$ $\longleftarrow$ <span style="color:red">true</span>

$\implies \quad \frac{1}{2} A_{\vec{\delta}}(x) \leq A'_{\vec{\delta}}(x) \leq 2 A_{\vec{\delta}}(x)$ for all $x \in X$, all $|\delta_i| < \epsilon$

- Precondition of theorem:

$\implies \quad \frac{1}{2} e(x) \leq e'(x) \leq 2 e(x)$ for all $x \in X$

# Sterbenz's Theorem

- How to apply the theorem automatically?

can apply
the theorem $\longrightarrow$ $e \ominus e' \;\sqsubseteq\; A_{\vec{\delta}}(x) - A'_{\vec{\delta}}(x)$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \; e' \sqsubseteq A'_{\vec{\delta}}$

- Check: $\displaystyle \min_{x, \vec{\delta}} \left( A'_{\vec{\delta}}(x) - \frac{1}{2} A_{\vec{\delta}}(x) \right) \geq 0$ and

  $\displaystyle \max_{x, \vec{\delta}} \left( A'_{\vec{\delta}}(x) - 2 A_{\vec{\delta}}(x) \right) \leq 0$ $\longleftarrow$ true

  $\implies \quad \frac{1}{2} A_{\vec{\delta}}(x) \leq A'_{\vec{\delta}}(x) \leq 2 A_{\vec{\delta}}(x)$ for all $x \in X$, all $|\delta_i| < \epsilon$

- Precondition of theorem:

  $\implies \quad \frac{1}{2} e(x) \leq e'(x) \leq 2 e(x)$ for all $x \in X$

# Sterbenz's Theorem

- How to apply the theorem automatically?

<span style="color:red">can apply the theorem</span> $\longrightarrow$ $e \ominus e' \ \sqsubseteq \ \textcolor{red}{A_{\vec{\delta}}(x) - A'_{\vec{\delta}}(x)}$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}$

- Check: $\min\limits_{x, \vec{\delta}} \left( A'_{\vec{\delta}}(x) - \frac{1}{2} A_{\vec{\delta}}(x) \right) \geq 0$ and

  $\max\limits_{x, \vec{\delta}} \left( A'_{\vec{\delta}}(x) - 2 A_{\vec{\delta}}(x) \right) \leq 0$ $\longleftarrow$ <span style="color:red">true</span>

  $\Longrightarrow \quad \frac{1}{2} A_{\vec{\delta}}(x) \leq A'_{\vec{\delta}}(x) \leq 2 A_{\vec{\delta}}(x)$ for all $x \in X$, all $|\delta_i| < \epsilon$

- Precondition of theorem:

  $\Longrightarrow \quad \frac{1}{2} e(x) \leq e'(x) \leq 2 e(x)$ for all $x \in X$

# Dekker's Theorem

- <u>Theorem</u> [Dekker, 1971]

$$|a| \geq |b| \quad \text{and} \quad r = a \oplus b \ominus a \ominus b$$

# Dekker's Theorem

- <u>Theorem</u> [Dekker, 1971]

$$|a| \geq |b| \ \text{ and } \ r = a \oplus b \ominus a \ominus b$$
$$\implies \ r = (a \oplus b) - (a + b)$$

# Dekker's Theorem

- <u>Theorem</u> [Dekker, 1971]

$$|a| \geq |b| \quad \text{and} \quad r = a \oplus b \ominus a \ominus b$$
$$\implies \quad r = (a \oplus b) - (a + b)$$

- Example: $a + b + c$, where $a \geq b \geq c > 0$

# Dekker's Theorem

- <u>Theorem</u> [Dekker, 1971]

$$|a| \geq |b| \quad \text{and} \quad r = a \oplus b \ominus a \ominus b$$
$$\implies \quad r = (a \oplus b) - (a + b)$$

- Example: $a + b + c$, where $a \geq b \geq c > 0$

  - Naïve ways: $(a \oplus b) \oplus c, \ a \oplus (b \oplus c), \ \cdots$

# Dekker's Theorem

- <u>Theorem</u> [Dekker, 1971]

$$|a| \geq |b| \ \text{ and } \ r = a \oplus b \ominus a \ominus b$$
$$\implies \ r = (a \oplus b) - (a + b)$$

- Example: $a + b + c$, where $a \geq b \geq c > 0$

  - Naïve ways: $(a \oplus b) \oplus c, \ a \oplus (b \oplus c), \ \cdots$

  - More accurate way:

    1. Compute $r = a \oplus b \ominus a \ominus b$
    2. Return $(a \oplus b) \oplus (c \ominus r)$

# Dekker's Theorem

- <u>Theorem</u> [Dekker, 1971]

$$|a| \geq |b| \quad \text{and} \quad r = a \oplus b \ominus a \ominus b$$
$$\Longrightarrow \quad r = (a \oplus b) - (a + b)$$

- Example: $a + b + c$, where $\boxed{a \geq b} \geq c > 0$

  - Naïve ways: $(a \oplus b) \oplus c, \ a \oplus (b \oplus c), \ \cdots$

  - More accurate way:

    1. Compute $r = a \oplus b \ominus a \ominus b$
    2. Return $(a \oplus b) \oplus (c \ominus r)$

# Dekker's Theorem

- <u>Theorem</u> [Dekker, 1971]

$$|a| \geq |b| \ \text{ and } \ r = a \oplus b \ominus a \ominus b$$
$$\implies \ r = (a \oplus b) - (a + b)$$

- Example: $a + b + c$, where $\boxed{a \geq b} \geq c > 0$

  - Naïve ways: $(a \oplus b) \oplus c, \ a \oplus (b \oplus c), \ \cdots$

  - More accurate way:

    1. Compute $r = a \oplus b \ominus a \ominus b$
    2. Return $\underbrace{(a \oplus b)}_{\parallel} \oplus (c \ominus r)$
       $(a + b) + r$

# Dekker's Theorem

- <u>Theorem</u> [Dekker, 1971]

$$|a| \geq |b| \quad \text{and} \quad r = a \oplus b \ominus a \ominus b$$
$$\implies \quad r \ {\color{red}=} \ (a \oplus b) - (a + b)$$

- Example: $a + b + c$, where $\boxed{a \geq b} \geq c > 0$

    - Naïve ways: $(a \oplus b) \oplus c, \ a \oplus (b \oplus c), \ \cdots$

    - More accurate way:

        1. Compute $r = a \oplus b \ominus a \ominus b$
        2. Return $\underbrace{(a \oplus b)}_{=} \oplus (c \boxed{\ominus r})$

        $(a + b) \boxed{+ r}$

# Dekker's Theorem

- How to apply the theorem automatically?

# Dekker's Theorem

- How to apply the theorem automatically?

$$e \oplus e' \ominus e \ominus e' \quad \sqsubseteq \quad ?$$

# Dekker's Theorem

- How to apply the theorem automatically?

$$e \oplus e' \ominus e \ominus e' \sqsubseteq \quad ?$$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \; e' \sqsubseteq A'_{\vec{\delta}}, \; e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right) (1 + \delta)$

# Dekker's Theorem

- How to apply the theorem automatically?

$$e \oplus e' \ominus e \ominus e' \quad \sqsubseteq \quad ?$$

  - Construct: $e \sqsubseteq A_{\vec{\delta}}, \; e' \sqsubseteq A'_{\vec{\delta}}, \; e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right) (1 + \delta)$

$$\min_{x, \vec{\delta}} \left| A_{\vec{\delta}}(x) \right| \qquad \max_{x, \vec{\delta}} \left| A'_{\vec{\delta}}(x) \right|$$

# Dekker's Theorem

- How to apply the theorem automatically?

$$e \oplus e' \ominus e \ominus e' \quad \sqsubseteq \quad ?$$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}, \ e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right) (1 + \delta)$

- Check: $\min\limits_{x,\vec{\delta}} \left| A_{\vec{\delta}}(x) \right| \geq \max\limits_{x,\vec{\delta}} \left| A'_{\vec{\delta}}(x) \right|$

# Dekker's Theorem

- How to apply the theorem automatically?

$$e \oplus e' \ominus e \ominus e' \sqsubseteq \quad ?$$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \; e' \sqsubseteq A'_{\vec{\delta}}, \; e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right)(1 + \delta)$

- Check: $\min\limits_{x, \vec{\delta}} \left| A_{\vec{\delta}}(x) \right| \geq \max\limits_{x, \vec{\delta}} \left| A'_{\vec{\delta}}(x) \right|$ $\longleftarrow$ true

# Dekker's Theorem

- How to apply the theorem automatically?

$$e \oplus e' \ominus e \ominus e' \ \sqsubseteq \ ?$$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}, \ e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right)(1 + \delta)$

- Check: $\boxed{\min_{x,\vec{\delta}} \left| A_{\vec{\delta}}(x) \right| \geq \max_{x,\vec{\delta}} \left| A'_{\vec{\delta}}(x) \right|}$ ← <span style="color:red">true</span>

$\implies \quad \left| A_{\vec{\delta}}(x) \right| \geq \left| A'_{\vec{\delta}}(x) \right| \quad$ for all $x \in X$, all $|\delta_i| < \epsilon$

# Dekker's Theorem

- How to apply the theorem automatically?

$$e \oplus e' \ominus e \ominus e' \quad \sqsubseteq \quad ?$$

  - Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}, \ e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right)(1 + \delta)$

  - Check: $\quad \min\limits_{x, \vec{\delta}} \left| A_{\vec{\delta}}(x) \right| \geq \max\limits_{x, \vec{\delta}} \left| A'_{\vec{\delta}}(x) \right|$ $\longleftarrow$ <span style="color:red">true</span>

    $\Longrightarrow \quad \left| A_{\vec{\delta}}(x) \right| \geq \left| A'_{\vec{\delta}}(x) \right| \qquad$ for all $x \in X$, all $|\delta_i| < \epsilon$

  - Precondition of theorem:
    $\Longrightarrow \quad |e(x)| \geq |e'(x)| \qquad$ for all $x \in X$

# Dekker's Theorem

- How to apply the theorem automatically?

<span style="color:red">can apply the theorem</span> → $e \oplus e' \ominus e \ominus e' \sqsubseteq \ ?$

  - Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}, \ e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right) (1 + \delta)$

  - Check: $\min\limits_{x, \vec{\delta}} \left| A_{\vec{\delta}}(x) \right| \geq \max\limits_{x, \vec{\delta}} \left| A'_{\vec{\delta}}(x) \right|$ ← <span style="color:red">true</span>

    $\Longrightarrow \quad \left| A_{\vec{\delta}}(x) \right| \geq \left| A'_{\vec{\delta}}(x) \right| \quad$ for all $x \in X$, all $|\delta_i| < \epsilon$

  - Precondition of theorem:

    $\Longrightarrow \quad |e(x)| \geq |e'(x)| \quad$ for all $x \in X$

# Dekker's Theorem

- How to apply the theorem automatically?

can apply
the theorem $\longrightarrow$ $e \oplus e' \ominus e \ominus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right) \delta$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}, \ e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right) (1 + \delta)$

- Check: $\min\limits_{x, \vec{\delta}} \left| A_{\vec{\delta}}(x) \right| \geq \max\limits_{x, \vec{\delta}} \left| A'_{\vec{\delta}}(x) \right|$ $\longleftarrow$ true

$\Longrightarrow \quad \left| A_{\vec{\delta}}(x) \right| \geq \left| A'_{\vec{\delta}}(x) \right| \quad$ for all $x \in X$, all $|\delta_i| < \epsilon$

- Precondition of theorem:

$\Longrightarrow \quad |e(x)| \geq |e'(x)| \quad$ for all $x \in X$

# Dekker's Theorem

- How to apply the theorem automatically?

can apply the theorem $\longrightarrow$ $e \oplus e' \ominus e \ominus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right) \delta$

- Construct: $e \sqsubseteq A_{\vec{\delta}}, \ e' \sqsubseteq A'_{\vec{\delta}}, \ e \oplus e' \sqsubseteq \left( A_{\vec{\delta}}(x) + A'_{\vec{\delta}}(x) \right)(1 + \delta)$

- Check: $\min\limits_{x, \vec{\delta}} \left| A_{\vec{\delta}}(x) \right| \geq \max\limits_{x, \vec{\delta}} \left| A'_{\vec{\delta}}(x) \right|$ $\longleftarrow$ true

$\Longrightarrow \quad \left| A_{\vec{\delta}}(x) \right| \geq \left| A'_{\vec{\delta}}(x) \right| \quad$ for all $x \in X$, all $|\delta_i| < \epsilon$

- Precondition of theorem:

$\Longrightarrow \quad |e(x)| \geq |e'(x)| \quad$ for all $x \in X$

# More in Our Paper

- More complicated <span style="color:red">exactness results</span>

  - Their explicit statements
  - How to apply them automatically

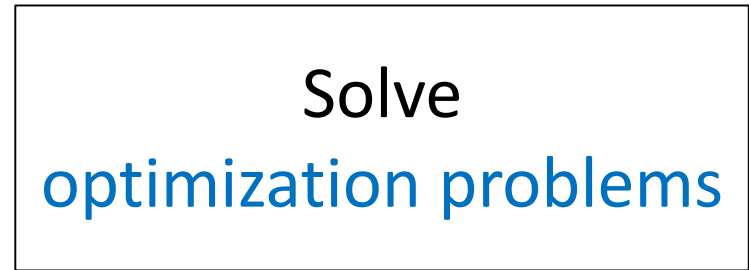# More in Our Paper

- More complicated <span style="color:red">exactness results</span>

  - Their explicit statements
  - How to apply them automatically

  - We check **preconditions** of exactness results

# More in Our Paper

- More complicated exactness results

  - Their explicit statements
  - How to apply them automatically

  - We check **preconditions** of exactness results

Solve
optimization problems

# More in Our Paper

- More complicated exactness results

  - Their explicit statements
  - How to apply them automatically

  - We check **preconditions** of exactness results

Solve
optimization problems

- Abstractions

  - In practice, we cannot use naïve abstractions
    because solving $\max\limits_{x,\,\vec{\delta}} \left( \cdots A_{\vec{\delta}}(x) \cdots \right)$ is difficult

# More in Our Paper

- More complicated exactness results

  - Their explicit statements
  - How to apply them automatically

  - We check **preconditions** of exactness results

- Abstractions

  - In practice, we cannot use naïve abstractions
    because solving $\max\limits_{x,\,\vec{\delta}} \left( \cdots A_{\vec{\delta}}(x) \cdots \right)$ is difficult

  - We apply our **"compress" operation** to abstractions
    to reduce # of $\delta$ variables without losing much precision

Solve
optimization problems

# More in Our Paper

- More complicated exactness results

  - Their explicit statements
  - How to apply them automatically

  - We check **preconditions** of exactness results

- Abstractions

  - In practice, we cannot use naïve abstractions
    because solving $\max\limits_{x,\,\vec{\delta}} \left( \cdots A_{\vec{\delta}}(x) \cdots \right)$ is difficult

  - We apply our **"compress" operation** to abstractions
    to reduce # of $\delta$ variables without losing much precision

Solve
optimization problems

# Case Studies

- Benchmarks
  - `sin`, `log`, `tan`: Intel's implementation of `math.h`
  - These are written directly in x86 assembly

# Case Studies

- Benchmarks
  - `sin`, `log`, `tan`: Intel's implementation of `math.h`
  - These are written directly in x86 assembly

- Use bit-level / integer arithmetic operations (bit-shift, int add, ⋯)

  - Mixed in with floating-point operations

# Case Studies

- Benchmarks
  - `sin, log, tan`: Intel's implementation of `math.h`
  - These are written directly in x86 assembly

- Use <span style="color:red">bit-level / integer arithmetic operations</span> (bit-shift, int add, ⋯)

  - Mixed in with floating-point operations
  - Eliminate them using the technique from [PLDI'16]
  - → Results in multiple independent subproblems (on subintervals)

# Case Studies

- Benchmarks
  - `sin`, `log`, `tan`: Intel's implementation of `math.h`
  - These are written directly in x86 assembly

- Use <span style="color:red">bit-level / integer arithmetic operations</span> (bit-shift, int add, ⋯)
  - Mixed in with floating-point operations
  - Eliminate them using the technique from [PLDI'16]
  - → Results in multiple independent subproblems (on subintervals)

- Apply our technique to the resulting floating-point expressions
  - Use Mathematica to solve optimization problems <span style="color:red">analytically</span>

ulp error
(log scale)

input value

| | error bounds from [PLDI'16] |
|---|---|
| ——— | our error bounds |
| —·—·— | 1 ulp |
| • | actual ulp errors |

# Results: `sin` on $[-\pi, \pi]$



ulp error
(log scale)

input value

| | |
|---|---|
| - - - - - | error bounds from [PLDI'16] |
| ——— | our error bounds |
| — · — | 1 ulp |
| • | actual ulp errors |

# Results: `sin` on $[-\pi, \pi]$



ulp error
(log scale)

input value

0.530 ulps

----- error bounds from [PLDI'16]
——— our error bounds
— · — 1 ulp
• actual ulp errors

# Results: sin on $[-\pi, \pi]$

# Results: `log` on $(2^{-1022}, 2^{1024})$



ulp error
(log scale)

input value
(log scale)

| | |
|---|---|
| - - - - - | error bounds from [PLDI'16] |
| ———— | our error bounds |
| — · — · — | 1 ulp |
| • | actual ulp errors |

error bounds from [PLDI'16]

our error bounds

1 ulp

actual ulp errors

# Results: $\log$ on $(2^{-1022}, 2^{1024})$



ulp error
(log scale)

0.583 ulps

input value
(log scale)

- - - - - error bounds from [PLDI'16]
———— our error bounds
— · — · 1 ulp
• actual ulp errors

• 461 hours on 16 cores



ulp error
(log scale)

1.0E+14

1.6E+01

1.0E+00

6.3E-02

2.2E-308   2.1E-154   2.0E+000   1.9E+154   1.8E+308

0.583 ulps

input value
(log scale)

- - - - - error bounds from [PLDI'16]
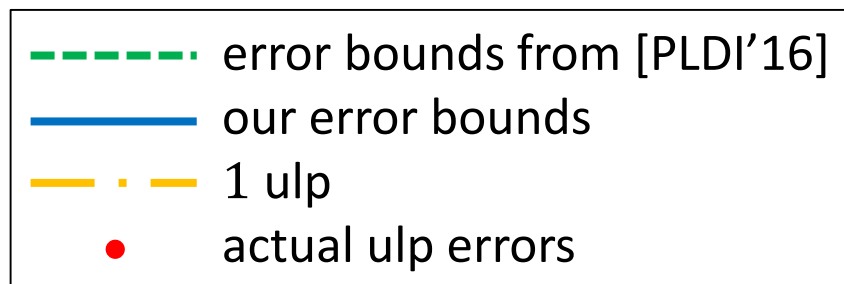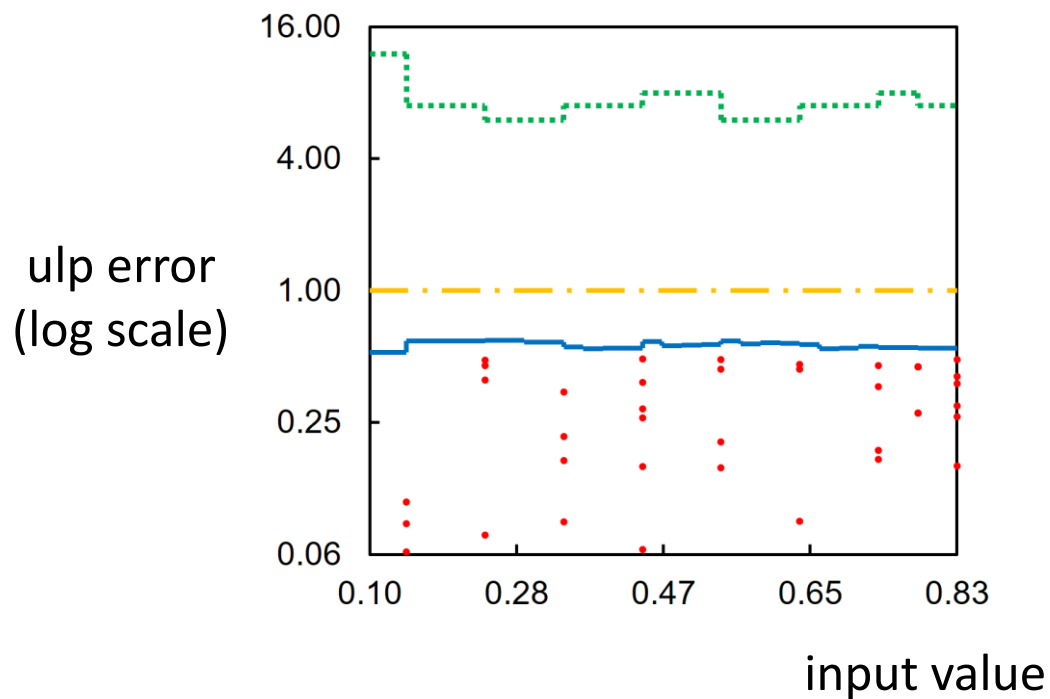———— our error bounds
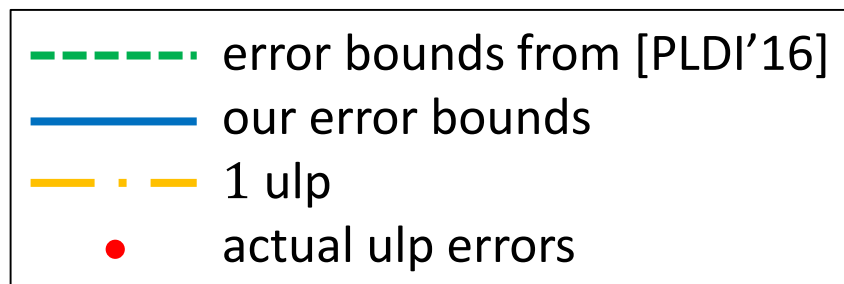— · — · — 1 ulp
●    actual ulp errors

# Results: `log o`

- 461 hours on 16 cores
- Highly parallelizable:
  - 4 million subintervals
  - < 16 sec on 1 core per subinterval

ulp error
(log scale)

1.0E+14

1.6E+01

1.0E+00

6.3E-02

0.583 ulps

2.2E-308    2.1E-154    2.0E+000    1.9E+154    1.8E+308

input value
(log scale)

- - - - - error bounds from [PLDI'16]
——— our error bounds
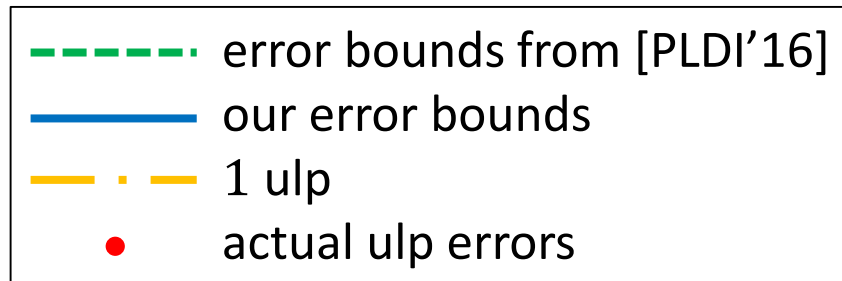— · — 1 ulp
• actual ulp errors

# Results: $\tan$ on $\left[\frac{13}{128}, \frac{17\pi}{64}\right), \left[\frac{17\pi}{64}, \frac{\pi}{2}\right)$



ulp error
(log scale)
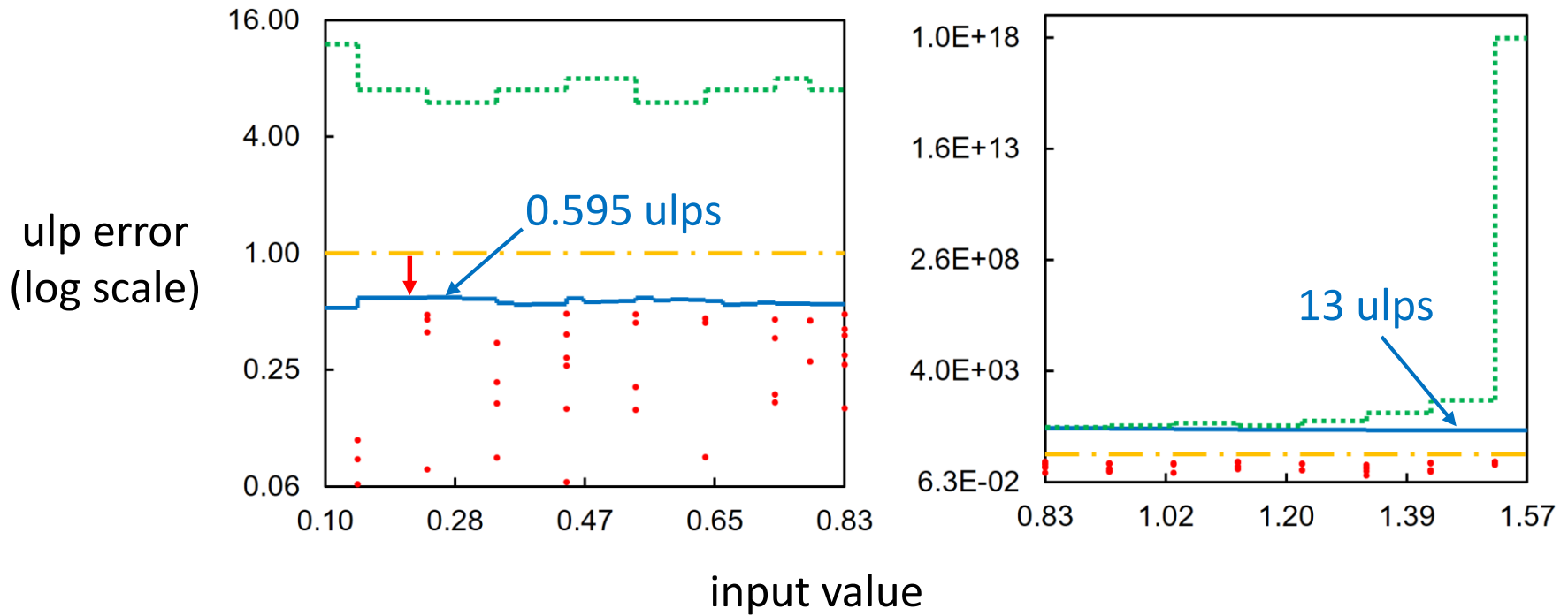
input value

error bounds from [PLDI'16]
our error bounds
1 ulp
actual ulp errors

# Results: $\tan$ on $\left[\frac{13}{128}, \frac{17\pi}{64}\right), \left[\frac{17\pi}{64}, \frac{\pi}{2}\right)$



ulp error
(log scale)

0.595 ulps

input value

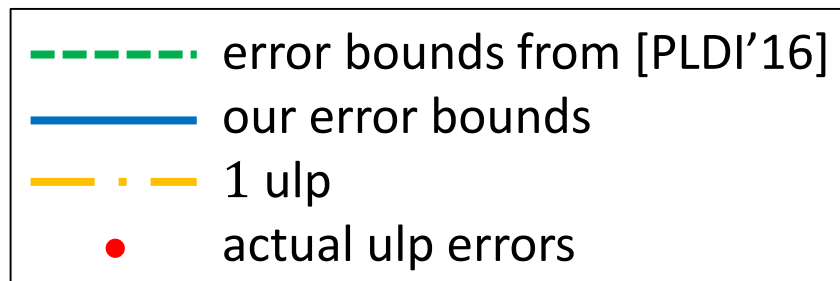| | |
|---|---|
| - - - - - | error bounds from [PLDI'16] |
| —— | our error bounds |
| — · — | 1 ulp |
| • | actual ulp errors |

# Results: $\tan$ on $\left[\frac{13}{128}, \frac{17\pi}{64}\right), \left[\frac{17\pi}{64}, \frac{\pi}{2}\right)$



ulp error
(log scale)

0.595 ulps

13 ulps

input value

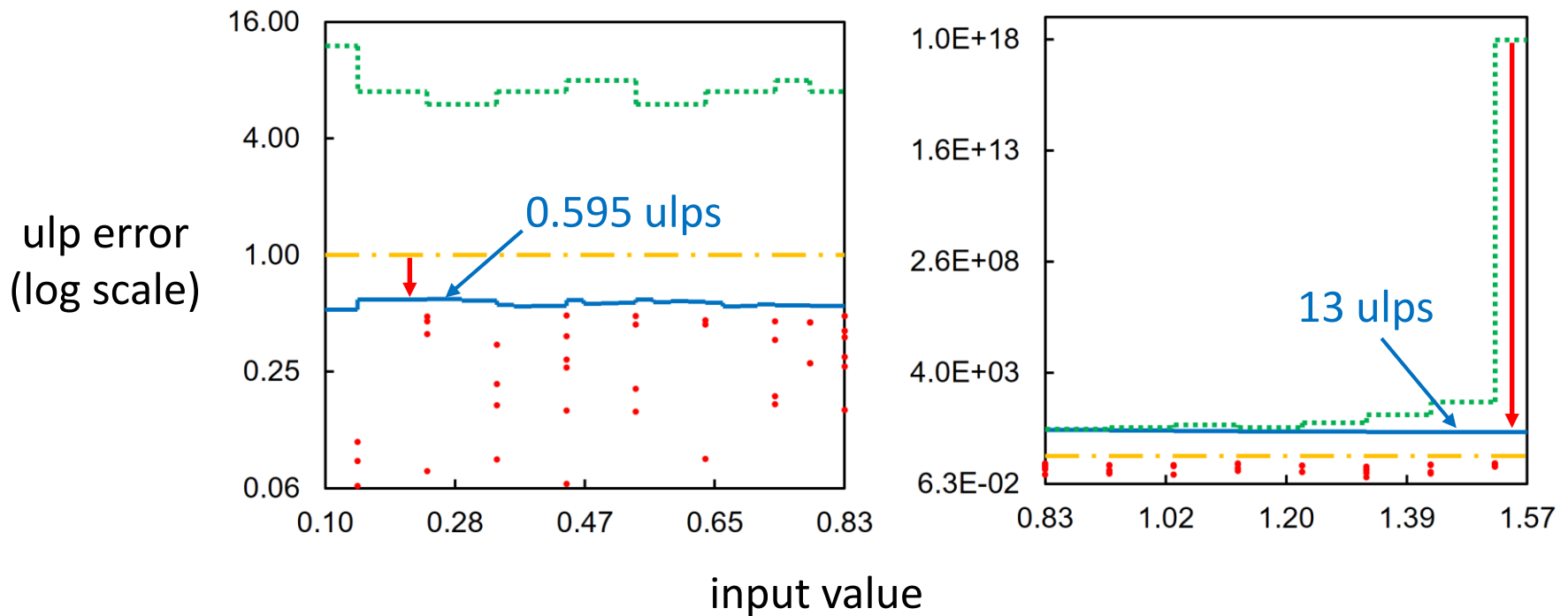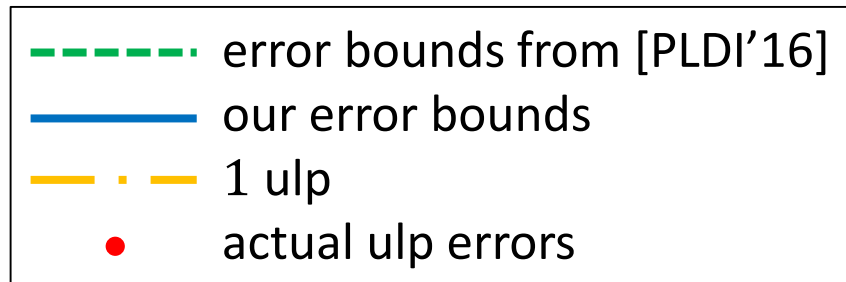| | error bounds from [PLDI'16] |
|---|---|
| | our error bounds |
| | 1 ulp |
| • | actual ulp errors |

# Results: $\tan$ on $\left[\frac{13}{128}, \frac{17\pi}{64}\right), \left[\frac{17\pi}{64}, \frac{\pi}{2}\right)$



ulp error (log scale)

0.595 ulps

13 ulps

input value

- - - - error bounds from [PLDI'16]
——— our error bounds
—·—·— 1 ulp
•  actual ulp errors

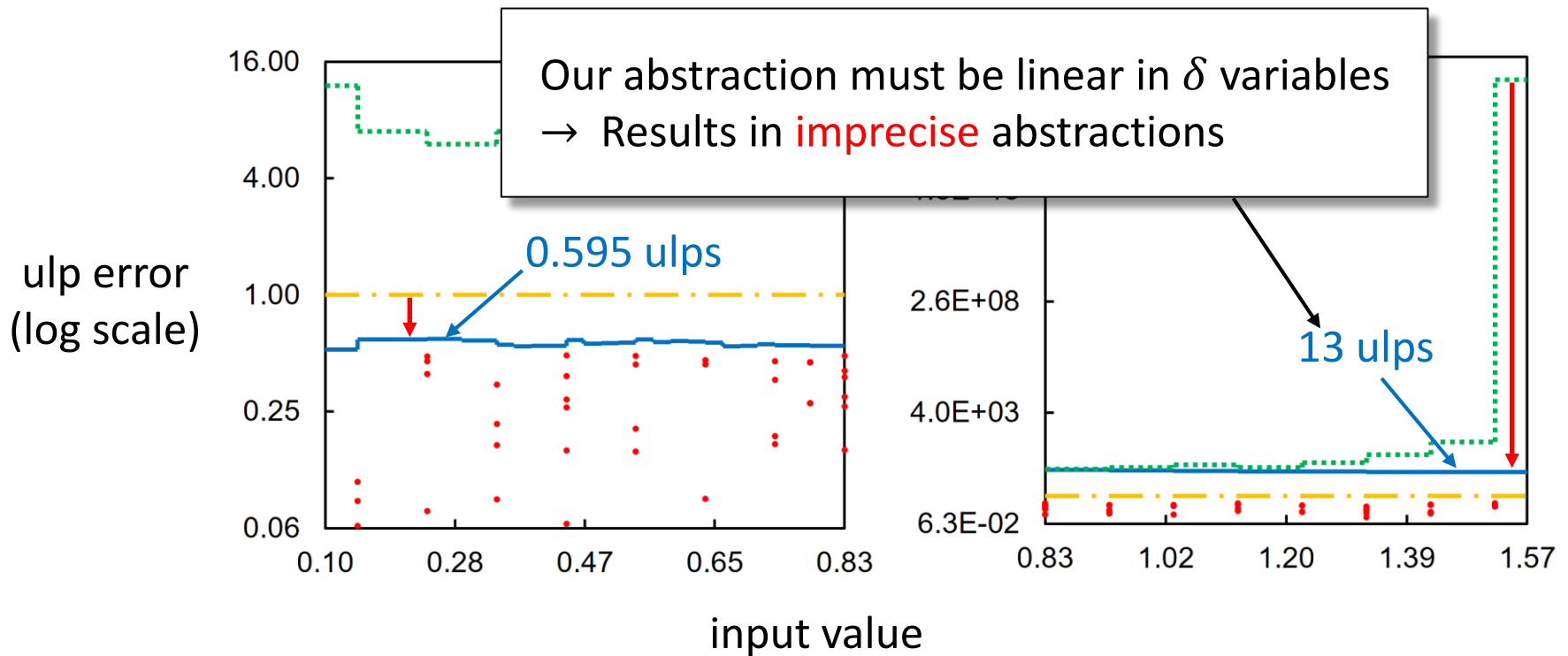# Results: $\tan$ on $\left[\frac{13}{128}, \frac{17\pi}{64}\right), \left[\frac{17\pi}{64}, \frac{\pi}{2}\right)$



Our abstraction must be linear in $\delta$ variables
→ Results in imprecise abstractions

0.595 ulps

13 ulps

ulp error
(log scale)

input value

- - - - - error bounds from [PLDI'16]
———— our error bounds
— · — · 1 ulp
• actual ulp errors

# Summary

- A novel static analysis for verifying `math.h` implementations automatically

- A reduction of this verification task to optimization problems

- Our technique verified some of Intel's `math.h` implementations automatically

# Summary

- A novel static analysis for verifying `math.h` implementations automatically

- A reduction of this verification task to optimization problems

- Our technique verified some of Intel's `math.h` implementations automatically

## Questions?