# Verifying Bit-Manipulations of Floating-Point

Wonyeol Lee    Rahul Sharma    Alex Aiken

Stanford University

# This Talk

- Example:

$$e^x$$

mathematical
specification

# This Talk

- Example:

$$e^x$$

mathematical
specification

```
...
vpslld  $20,    %xmm3, %xmm3
vpshufd $114,   %xmm3, %xmm3
vmulpd  C1,     %xmm2, %xmm1
vmulpd  C2,     %xmm2, %xmm2
...
```

floating-point implementation

# This Talk

- Example:

$$e^x$$

mathematical
specification

$$\neq$$

```
...
vpslld  $20,    %xmm3, %xmm3
vpshufd $114,   %xmm3, %xmm3
vmulpd  C1,     %xmm2, %xmm1
vmulpd  C2,     %xmm2, %xmm2
...
```

floating-point implementation

# This Talk

- Example:

$$e^x$$

mathematical
specification

← how different? →

```
...
vpslld  $20,    %xmm3, %xmm3
vpshufd $114,   %xmm3, %xmm3
vmulpd  C1,     %xmm2, %xmm1
vmulpd  C2,     %xmm2, %xmm2
...
```

floating-point implementation

# This Talk

- Example:

$$e^x$$

mathematical
specification

how different?

```
...
vpslld  $20,    %xmm3, %xmm3
vpshufd $114,   %xmm3, %xmm3
vmulpd  C1,     %xmm2, %xmm1
vmulpd  C2,     %xmm2, %xmm2
...
```

floating-point implementation

- Goal:  Bound the difference between spec and implementation

# This Talk

- Example:

$$e^x$$

mathematical
specification

how different?

```
...
vpslld  $20,    %xmm3,  %xmm3
vpshufd $114,   %xmm3,  %xmm3
vmulpd  C1,     %xmm2,  %xmm1
vmulpd  C2,     %xmm2,  %xmm2
...
```

floating-point implementation

- Goal:  Bound the difference between spec and implementation
- Key contribution: Verify binaries that mix floating-point and bit-level operations

# This Talk

- Example:
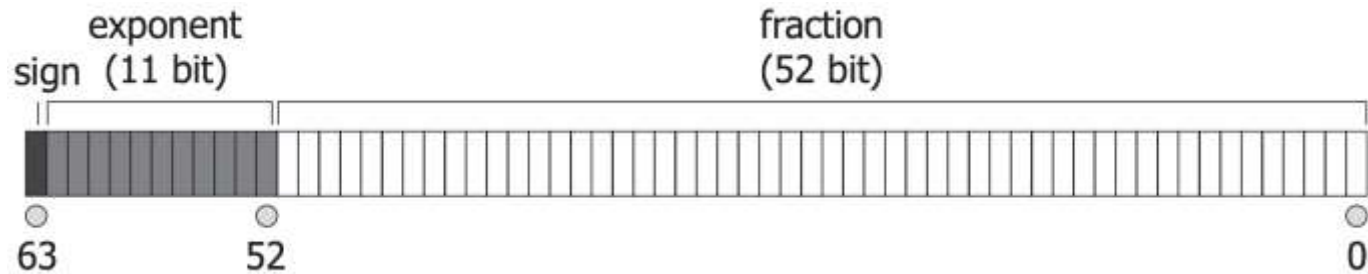
$$e^x$$

mathematical
specification

how different?

```
...
vpslld   $20,    %xmm3, %xmm3
vpshufd  $114,   %xmm3, %xmm3
vmulpd   C1,     %xmm2, %xmm1
vmulpd   C2,     %xmm2, %xmm2
...
```

floating-point implementation

- Goal: Bound the difference between spec and implementation
- Key contribution: Verify binaries that mix floating-point and bit-level operations
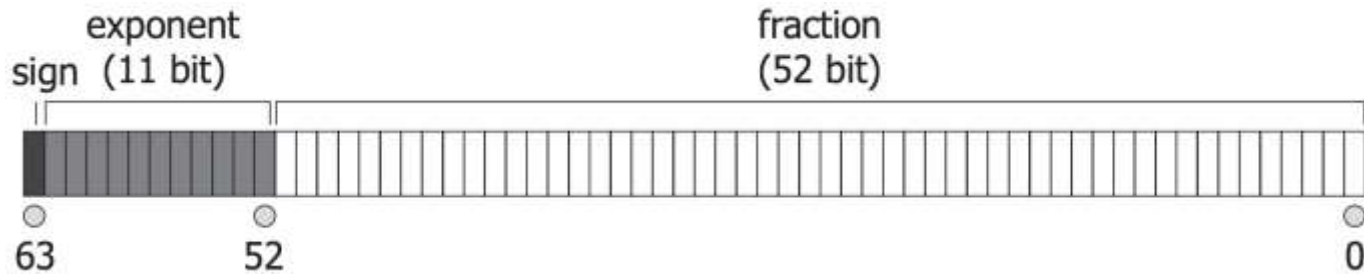  - Intel's implementations of transcendental functions

# Floating-Point Numbers



- Example:

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

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

# Floating-Point Numbers



- Example:

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

- Automatic reasoning about floating-point is not easy
  - have rounding errors
  - don't obey some algebraic rules of real numbers
  - Associativity: $1 + (10^{30} - 10^{30}) = 1 \neq 0 = (1 + 10^{30}) - 10^{30}$

# Floating-Point Numbers



- Example:  $\underset{\text{1}}{1}$  $\underset{\text{01111111111}}{}$  $\underset{\text{1100}\cdots\text{00}}{}$ $_{(2)}$

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

- Automatic reasoning about floating-point is not easy
  - have rounding errors
  - don't obey some algebraic rules of real numbers
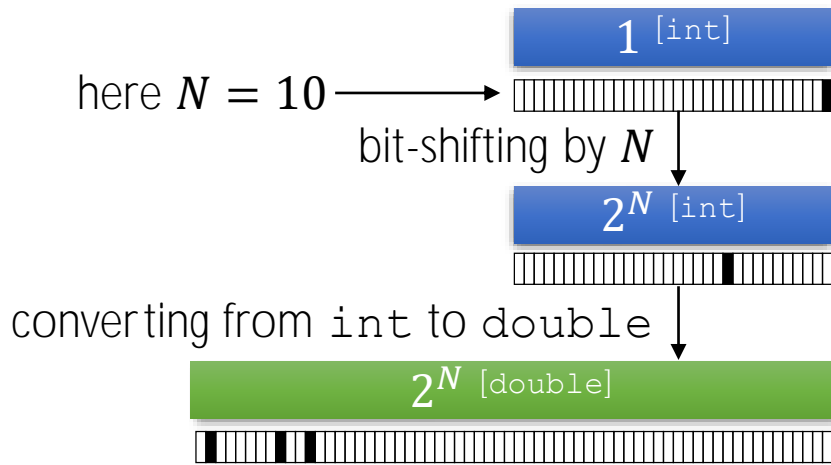  - Associativity: $1 + (10^{30} - 10^{30}) = 1 \neq 0 = (1 + 10^{30}) - 10^{30}$

- It becomes much harder if bit-level operations are used

# Bit-Level Operations

- Example: Given $N$ (in `int`), compute $2^N$ (in `double`)

# Bit-Level Operations

- Example: Given $N$ (in `int`), compute $2^N$ (in `double`)



here $N = 10$ → bit-shifting by $N$

$1$ [int]

$2^N$ [int]

converting from `int` to `double`

$2^N$ [double]

# Bit-Level Operations

- Example: Given $N$ (in `int`), compute $2^N$ (in `double`)



here $N = 10$ → 1 [int]

bit-shifting by $N$

$2^N$ [int]

expensive

converting from `int` to `double`

$2^N$ [double]

# Bit-Level Operations

- Example: Given $N$ (in `int`), compute $2^N$ (in `double`)



here $N = 10$

bit-shifting by $N$

1 [int]

$2^N$ [int]

expensive

converting from `int` to `double`

$2^N$ [double]

works only for $0 \leq N \leq 31$

# Bit-Level Operations

- Example: Given $N$ (in `int`), compute $2^N$ (in `double`)

here $N = 10$ ⟶ | $1$ [int] |

bit-shifting by $N$ ⟶ | $2^N$ [int] |

expensive

converting from `int` to `double` ⟶ | $2^N$ [double] |

works only for $0 \leq N \leq 31$

| $N$ [int] |

integer addition ⟶ | $N + 1023$ [int] |

bit-shifting by $52$

[12 bits] | $00 \cdots 0$ [52 bits]

# Bit-Level Operations

- Example: Given $N$ (in `int`), compute $2^N$ (in `double`)



here $N = 10$ → $1$ [int]

bit-shifting by $N$

$2^N$ [int]

expensive

converting from `int` to `double`

$2^N$ [double]

works only for $0 \leq N \leq 31$

$N$ [int]

integer addition

$N + 1023$ [int]

bit-shifting by $52$

$2^N$ [double]

# Bit-Level Operations

- Example: Given $N$ (in `int`), compute $2^N$ (in `double`)



here $N = 10$

$1$ [int]

bit-shifting by $N$

expensive

converting from `int` to `double`

$2^N$ [int]

$2^N$ [double]

works only for $0 \leq N \leq 31$

$N$ [int]

integer addition

$N + 1023$ [int]

bit-shifting by $52$

$2^N$ [double]

works for $-1022 \leq N \leq 1023$

# Bit-Level Operations

- Example: Given $N$ (in `int`), compute $2^N$ (in `double`)



here $N = 10$ → $1$ [int]

bit-shifting by $N$

$2^N$ [int]

expensive

converting from `int` to `double`

$2^N$ [double]

works only for $0 \leq N \leq 31$

$N$ [int]

integer addition

$N + 1023$ [int]

bit-shifting by $52$

$2^N$ [double]

works for $-1022 \leq N \leq 1023$

- Such bit-manipulations are ubiquitous in highly optimized floating-point implementations

- If a code mixes floating-point and bit-level operations, reasoning about the code is difficult

19

# Problem Statement

$$e^x$$

mathematical
specification
$f: \mathbb{R} \to \mathbb{R}$

# Problem Statement

$$e^x$$

mathematical
specification
$f \colon \mathbb{R} \to \mathbb{R}$

```
...
vpslld  $20,    %xmm3, %xmm3
vpshufd $114,   %xmm3, %xmm3
vmulpd  C1,     %xmm2, %xmm1
vmulpd  C2,     %xmm2, %xmm2
...
```

binary $P$ that mixes floating-point
and bit-level operations

# Problem Statement

$$e^x$$

mathematical
specification
$f : \mathbb{R} \rightarrow \mathbb{R}$

$[-1, 1]$
input range $X \subseteq \mathbb{R}$

⟷

```
...
vpslld  $20,     %xmm3, %xmm3
vpshufd $114,    %xmm3, %xmm3
vmulpd  C1,      %xmm2, %xmm1
vmulpd  C2,      %xmm2, %xmm2
...
```

binary $P$ that mixes floating-point
and bit-level operations

# Problem Statement

$$e^x$$

mathematical
specification
$f: \mathbb{R} \to \mathbb{R}$

$$[-1, 1]$$
input range $X \subseteq \mathbb{R}$

$\longleftrightarrow$

```
...
vpslld   $20,     %xmm3, %xmm3
vpshufd  $114,    %xmm3, %xmm3
vmulpd   C1,      %xmm2, %xmm1
vmulpd   C2,      %xmm2, %xmm2
...
```

binary $P$ that mixes floating-point
and bit-level operations

- Goal: Find a small $\Theta > 0$ such that
$$\left| \frac{f(x) - P(x)}{f(x)} \right| \leq \Theta \text{ for all } x \in X$$

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

# Possible Alternatives

- Exhaustive testing
    - feasible for 32-bit float: $\sim 30$ seconds      (with $1$ core for `sinf`)
    - infeasible for 64-bit double: $> 4000$ years   $(= 30$ seconds $\times\ 2^{32})$

# Possible Alternatives

- Exhaustive testing
  - feasible for 32-bit float: $\sim 30$ seconds        (with $1$ core for `sinf`)
  - infeasible for 64-bit double: $> 4000$ years  ($= 30$ seconds $\times 2^{32}$)
  - infeasible even for input range $X = [-1, 1]$
    $\because$ (# of doubles between $-1$ and $1$) $= \frac{1}{2}$ (# of all doubles)

# Possible Alternatives

- Exhaustive testing
  - feasible for 32-bit float: $\sim 30$ seconds      (with $1$ core for `sinf`)
  - infeasible for 64-bit double: $> 4000$ years  (= $30$ seconds $\times 2^{32}$)
  - infeasible even for input range $X = [-1, 1]$
    $\because$ (# of doubles between $-1$ and $1$) $= \frac{1}{2}$ (# of all doubles)

- Machine-checkable proofs
  - Harrison used HOL Light to prove Intel's transcendental functions are very accurate [FMCAD'00]

# Possible Alternatives

- Exhaustive testing
  - feasible for 32-bit float: $\sim 30$ seconds (with $1$ core for `sinf`)
  - infeasible for 64-bit double: $> 4000$ years $(= 30$ seconds $\times 2^{32})$
  - infeasible even for input range $X = [-1, 1]$
    $\because$ (# of doubles between $-1$ and $1$) $= \frac{1}{2}$ (# of all doubles)

- Machine-checkable proofs
  - Harrison used HOL Light to prove Intel's transcendental functions are very accurate [FMCAD'00]
  - "The construction of these proofs often requires considerable persistence." [FMSD'00]

# Possible Automatic Alternatives

- If only floating-point operations are used,
  various automatic techniques can be applied
  - e.g., Astree [PLDI'03], Fluctuat [FMICS'09], ROSA [POPL'14], FPTaylor [FM'15]

- Several commercial tools (e.g., Astree, Fluctuat) can handle
  certain bit-trick routines

# Possible Automatic Alternatives

- If only floating-point operations are used, various automatic techniques can be applied
  - e.g., Astree [PLDI'03], Fluctuat [FMICS'09], ROSA [POPL'14], FPTaylor [FM'15]

- Several commercial tools (e.g., Astree, Fluctuat) can handle certain bit-trick routines

- We are unaware of a general technique for verifying mixed floating-point and bit-level code

# Our Method

# $e^x$ Explained

```
 1  vmovddup    %xmm0,  %xmm0
 2  vmulpd      L2E,    %xmm0,  %xmm2
 3  vroundpd    $0,     %xmm2,  %xmm2
 4  vcvtpd2dqx  %xmm2,  %xmm3
 5  vpaddd      B,      %xmm3,  %xmm3
 6  vpslld      $20,    %xmm3,  %xmm3
 7  vpshufd     $114,   %xmm3,  %xmm3
 8  vmulpd      C1,     %xmm2,  %xmm1
 9  vmulpd      C2,     %xmm2,  %xmm2
10  vaddpd      %xmm1,  %xmm0,  %xmm1
11  vaddpd      %xmm2,  %xmm1,  %xmm1
12  vmovapd     T1,     %xmm0
13  vmulpd      T12,    %xmm1,  %xmm2
14  vaddpd      T11,    %xmm2,  %xmm2
    ...
36  vaddpd      %xmm0,  %xmm1,  %xmm0
37  vmulpd      %xmm3,  %xmm0,  %xmm0
38  retq
```

# $e^x$ Explained

```
 1 vmovddup    %xmm0,  %xmm0
 2 vmulpd      L2E,    %xmm0,  %xmm2
 3 vroundpd    $0,     %xmm2,  %xmm2
 4 vcvtpd2dqx  %xmm2,  %xmm3
 5 vpaddd      B,      %xmm3,  %xmm3
 6 vpslld      $20,    %xmm3,  %xmm3
 7 vpshufd     $114,   %xmm3,  %xmm3
 8 vmulpd      C1,     %xmm2,  %xmm1
 9 vmulpd      C2,     %xmm2,  %xmm2
10 vaddpd      %xmm1,  %xmm0,  %xmm1
11 vaddpd      %xmm2,  %xmm1,  %xmm1
12 vmovapd     T1,     %xmm0
13 vmulpd      T12,    %xmm1,  %xmm2
14 vaddpd      T11,    %xmm2,  %xmm2
   ...
36 vaddpd      %xmm0,  %xmm1,  %xmm0
37 vmulpd      %xmm3,  %xmm0,  %xmm0
38 retq
```

$x$

$N = \text{round}(x \cdot \log_2 e)$

# $e^x$ Explained

```
 1  vmovddup    %xmm0,  %xmm0
 2  vmulpd      L2E,    %xmm0,  %xmm2
 3  vroundpd    $0,     %xmm2,  %xmm2
 4  vcvtpd2dqx  %xmm2,  %xmm3
 5  vpaddd      B,      %xmm3,  %xmm3
 6  vpslld      $20,    %xmm3,  %xmm3
 7  vpshufd     $114,   %xmm3,  %xmm3
 8  vmulpd      C1,     %xmm2,  %xmm1
 9  vmulpd      C2,     %xmm2,  %xmm2
10  vaddpd      %xmm1,  %xmm0,  %xmm1
11  vaddpd      %xmm2,  %xmm1,  %xmm1
12  vmovapd     T1,     %xmm0
13  vmulpd      T12,    %xmm1,  %xmm2
14  vaddpd      T11,    %xmm2,  %xmm2
    ...
36  vaddpd      %xmm0,  %xmm1,  %xmm0
37  vmulpd      %xmm3,  %xmm0,  %xmm0
38  retq
```

$x$

$N = \text{round}(x \cdot \log_2 e)$

$2^N$

# $e^x$ Explained

```
 1  vmovddup   %xmm0,  %xmm0
 2  vmulpd     L2E,    %xmm0,  %xmm2
 3  vroundpd   $0,     %xmm2,  %xmm2
 4  vcvtpd2dqx %xmm2,  %xmm3
 5  vpaddd     B,      %xmm3,  %xmm3
 6  vpslld     $20,    %xmm3,  %xmm3
 7  vpshufd    $114,   %xmm3,  %xmm3
 8  vmulpd     C1,     %xmm2,  %xmm1
 9  vmulpd     C2,     %xmm2,  %xmm2
10  vaddpd     %xmm1,  %xmm0,  %xmm1
11  vaddpd     %xmm2,  %xmm1,  %xmm1
12  vmovapd    T1,     %xmm0
13  vmulpd     T12,    %xmm1,  %xmm2
14  vaddpd     T11,    %xmm2,  %xmm2
    ...
36  vaddpd     %xmm0,  %xmm1,  %xmm0
37  vmulpd     %xmm3,  %xmm0,  %xmm0
38  retq
```

$x$

$N = \text{round}(x \cdot \log_2 e)$

$2^N$

$r = x - N \cdot \ln 2$

$e^r \approx \sum_{i=0}^{12} \frac{r^i}{i!}$

$e^x = e^{N \cdot \ln 2} \cdot e^r \approx 2^N \cdot e^r$

34

# $e^x$ Explained

```
 1  vmovddup    %xmm0,  %xmm0
 2  vmulpd      L2E,    %xmm0,  %xmm2
 3  vroundpd    $0,     %xmm2,  %xmm2
 4  vcvtpd2dqx  %xmm2,  %xmm3
 5  vpaddd      B,      %xmm3,  %xmm3
 6  vpslld      $20,    %xmm3,  %xmm3
 7  vpshufd     $114,   %xmm3,  %xmm3
 8  vmulpd      C1,     %xmm2,  %xmm1
 9  vmulpd      C2,     %xmm2,  %xmm2
10  vaddpd      %xmm1,  %xmm0,  %xmm1
11  vaddpd      %xmm2,  %xmm1,  %xmm1
12  vmovapd     T1,     %xmm0
13  vmulpd      T12,    %xmm1,  %xmm2
14  vaddpd      T11,    %xmm2,  %xmm2
    ...
36  vaddpd      %xmm0,  %xmm1,  %xmm0
37  vmulpd      %xmm3,  %xmm0,  %xmm0
38  retq
```

$x$

$N = \mathrm{round}(x \cdot \log_2 e)$

$2^N$

$r = x - N \cdot \ln 2$

$e^r \approx \sum_{i=0}^{12} \frac{r^i}{i!}$

$e^x = e^{N \cdot \ln 2} \cdot e^r \approx 2^N \cdot e^r$

Goal: Find a small $\Theta > 0$ such that
$$\left| \frac{e^x - 2^N e^r}{e^x} \right| \leq \Theta \text{ for all } x \in X$$

35

# 1) Abstract Floating-Point Operations

- Assume only floating-point operations are used

# 1) Abstract Floating-Point Operations

- Assume only floating-point operations are used
- $(1 + \epsilon)$ property
  - A standard way to model rounding errors
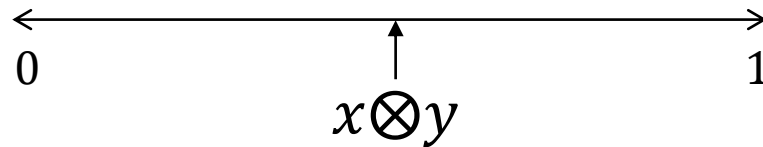
# 1) Abstract Floating-Point Operations

- Assume only floating-point operations are used
- $(1 + \epsilon)$ property
  - A standard way to model rounding errors

$$x \otimes_{\mathrm{f}} y \in \{(x \otimes y)(1 + \delta) : |\delta| < \epsilon\}$$

# 1) Abstract Floating-Point Operations

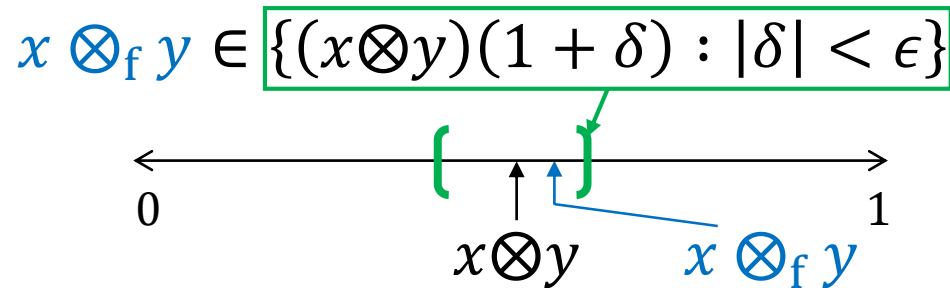- Assume only floating-point operations are used
- <span style="color:red">$(1 + \epsilon)$ property</span>
  - A standard way to model rounding errors

$$x \otimes_{\mathrm{f}} y \in \{(x \otimes y)(1 + \delta) : |\delta| < \epsilon\}$$



$0$        $x \otimes y$        $1$

# 1) Abstract Floating-Point Operations

- Assume only floating-point operations are used
- $(1 + \epsilon)$ property
  - A standard way to model rounding errors

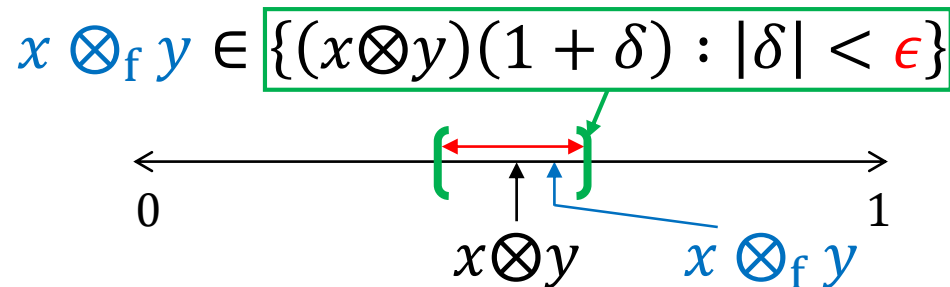$$x \otimes_{\mathrm{f}} y \in \{(x \otimes y)(1 + \delta) : |\delta| < \epsilon\}$$

# 1) Abstract Floating-Point Operations

- Assume only floating-point operations are used
- $(1 + \epsilon)$ property
  - A standard way to model rounding errors

$$x \otimes_{\mathrm{f}} y \in \{(x \otimes y)(1 + \delta) : |\delta| < \epsilon\}$$

# 1) Abstract Floating-Point Operations

- Assume only floating-point operations are used
- $(1 + \epsilon)$ property
  - A standard way to model rounding errors

$$x \otimes_{\mathrm{f}} y \in \{(x \otimes y)(1 + \delta) : |\delta| < \epsilon\}$$



$x \otimes y$    $x \otimes_{\mathrm{f}} y$

- For 64-bit doubles, $\epsilon = 2^{-53}$

# 1) Abstract Floating-Point Operations

- Assume only floating-point operations are used
- $(1 + \epsilon)$ property
  - A standard way to model rounding errors

$$x \otimes_{\mathrm{f}} y \in \{(x \otimes y)(1 + \delta) : |\delta| < \epsilon\}$$



- For 64-bit doubles, $\epsilon = 2^{-53}$
- This property has been used in previous automatic techniques (FPTaylor, ROSA, …) for verifying floating-point programs

# 1) Abstract Floating-Point Operations

- Compute a <span style="color:red">symbolic abstraction $A_{\vec{\delta}}(x)$</span> of a program $P$

# 1) Abstract Floating-Point Operations

- Compute a <span style="color:red">symbolic abstraction $A_{\vec{\delta}}(x)$</span> of a program $P$
  - Example:

$$P(x) = \big((2 \times_f x) \qquad\qquad +_f 3\big)$$

# 1) Abstract Floating-Point Operations

- Compute a <span style="color:red">symbolic abstraction $A_{\vec{\delta}}(x)$</span> of a program $P$
  - Example:

$$A_{\vec{\delta}}(x) = \big((2 \times_f x) \qquad +_f 3\big)$$

# 1) Abstract Floating-Point Operations

- Compute a <span style="color:red">symbolic abstraction $A_{\vec{\delta}}(x)$</span> of a program $P$
  - Example:

$$A_{\vec{\delta}}(x) = ((2 \times x) + 3)$$

# 1) Abstract Floating-Point Operations

- Compute a symbolic abstraction $A_{\vec{\delta}}(x)$ of a program $P$
  - Example:

    $$A_{\vec{\delta}}(x) = \big((2 \times x)(1 + \delta_1) + 3\big)(1 + \delta_2)$$

# 1) Abstract Floating-Point Operations

- Compute a symbolic abstraction $A_{\vec{\delta}}(x)$ of a program $P$
  - Example:

$$A_{\vec{\delta}}(x) = \big((2 \times x)(1 + \delta_1) + 3\big)(1 + \delta_2)$$

- From $(1 + \epsilon)$ property, $A_{\vec{\delta}}(x)$ satisfies

$$P(x) \in \big\{ A_{\vec{\delta}}(x) : |\delta_i| < \epsilon \big\} \text{ for all } x$$

# 1) Abstract Floating-Point Operations

- Compute a symbolic abstraction $A_{\vec{\delta}}(x)$ of a program $P$
  - Example:

  $$A_{\vec{\delta}}(x) = \big((2 \times x)(1 + \delta_1) + 3\big)(1 + \delta_2)$$

- From $(1 + \epsilon)$ property, $A_{\vec{\delta}}(x)$ satisfies

  $$P(x) \in \big\{A_{\vec{\delta}}(x) : |\delta_i| < \epsilon\big\} \text{ for all } x$$

  - Example:

  $$P(x) = \big((2 \times_f x) \qquad +_f 3\big)$$

# 1) Abstract Floating-Point Operations

- Compute a <span style="color:red">symbolic abstraction $A_{\vec{\delta}}(x)$</span> of a program $P$
  - Example:

  $$A_{\vec{\delta}}(x) = \big((2 \times x)(1 + \delta_1) + 3\big)(1 + \delta_2)$$

- From $(1 + \epsilon)$ property, $A_{\vec{\delta}}(x)$ satisfies

  $$P(x) \in \big\{ A_{\vec{\delta}}(x) : |\delta_i| < \epsilon \big\} \text{ for all } x$$

  - Example:

  $$P(x) \quad \big\{\big((2 \times x)(1 + \delta_1) + 3\big)(1 + \delta_2) : |\delta_1|, |\delta_2| < \epsilon \big\}$$

# 1) Abstract Floating-Point Operations

- Compute a symbolic abstraction $A_{\vec{\delta}}(x)$ of a program $P$
  - Example:

    $$A_{\vec{\delta}}(x) = \big((2 \times x)(1 + \delta_1) + 3\big)(1 + \delta_2)$$
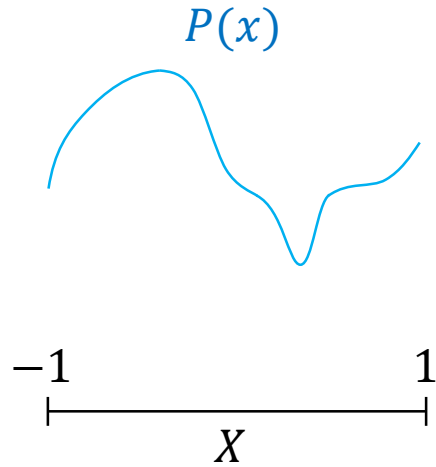
- From $(1 + \epsilon)$ property, $A_{\vec{\delta}}(x)$ satisfies

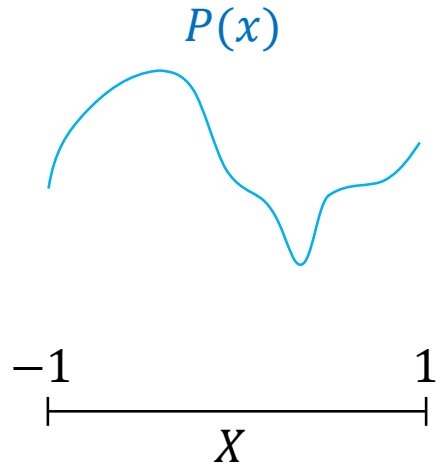  $$P(x) \in \big\{A_{\vec{\delta}}(x) : |\delta_i| < \epsilon\big\} \text{ for all } x$$

  - Example:

    $$P(x) \in \big\{\big((2 \times x)(1 + \delta_1) + 3\big)(1 + \delta_2) : |\delta_1|, |\delta_2| < \epsilon\big\}$$

# Our Method: Overview

$P(x)$



$-1$               $1$

$X$

```
...
vpslld  $20,    ...
vpshufd $114,   ...
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```

# Our Method: Overview

$P(x)$



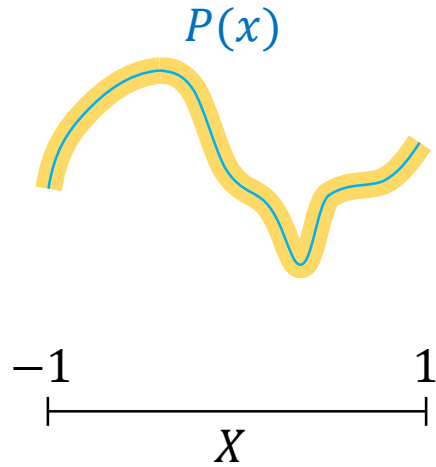$-1$                    $1$

$X$

```
...
vpslld  $20,    ...
vpshufd $114,   ...
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```
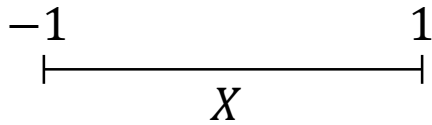
# Our Method: Overview

$P(x)$



$-1$                      $1$

$X$

```
...
vpslld  $20,    ...
vpshufd $114,   ...
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```

# Our Method: Overview

$P(x)$

hard to find

$-1$  ⊢————————⊣  $1$

$X$

```
. . .
vpslld  $20,    . . .
vpshufd $114,   . . .
vmulpd  C1,     . . .
vmulpd  C2,     . . .
. . .
```

# Our Method: Overview

$P(x)$

hard to find

$-1$        $1$
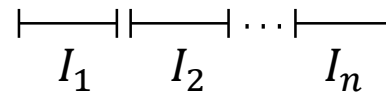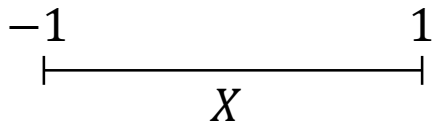
$X$

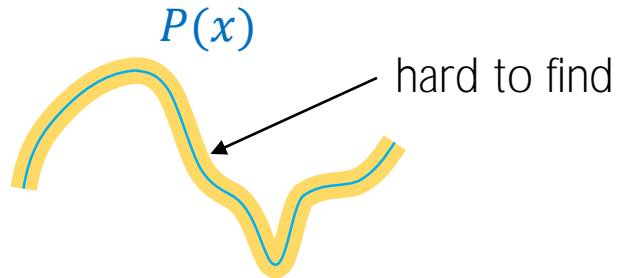not "smooth"

```
...
vpslld   $20,    ...
vpshufd  $114,   ...
vmulpd   C1,     ...
vmulpd   C2,     ...
...
```

abstract using
"smooth" functions

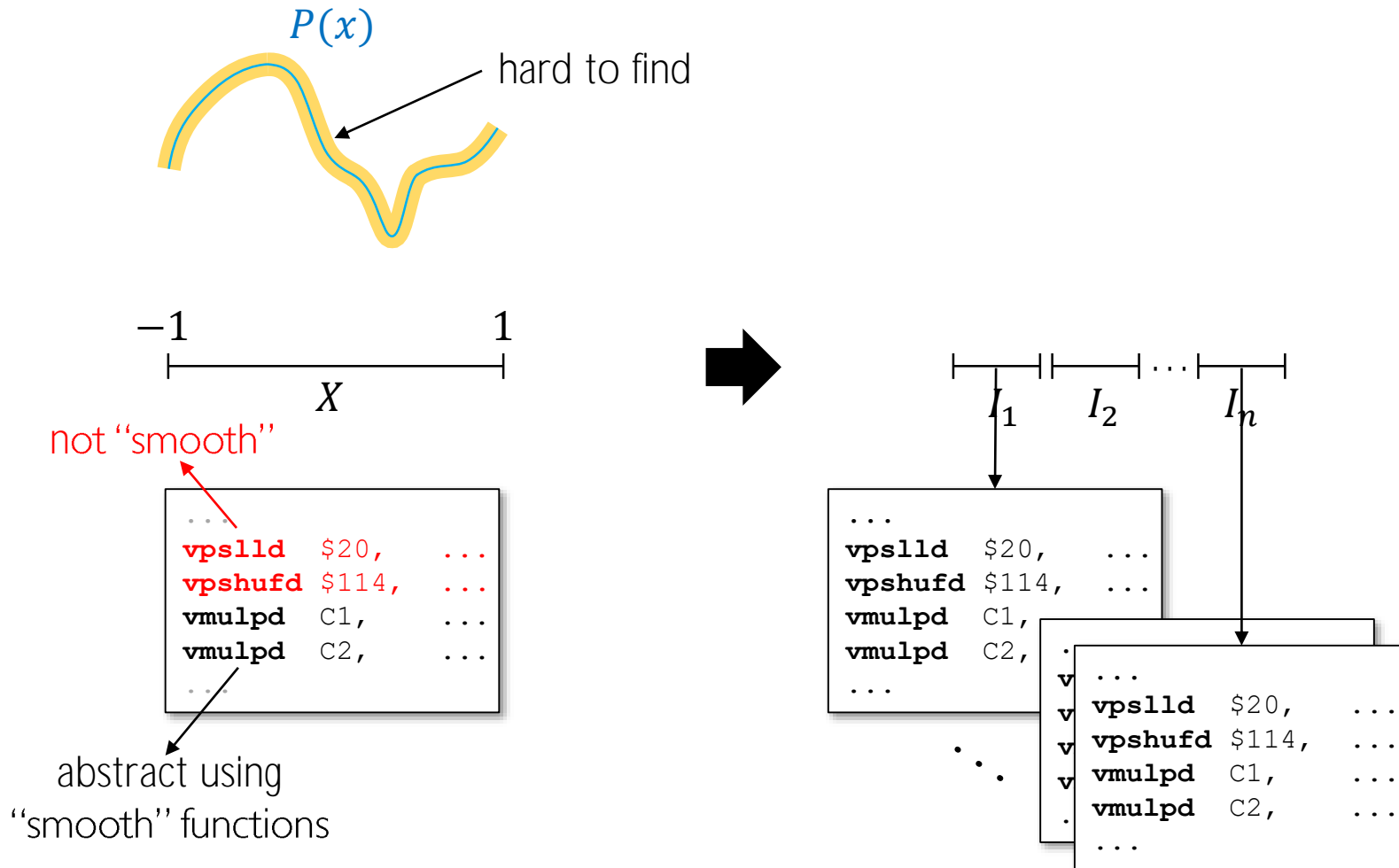# Our Method: Overview

$P(x)$

hard to find

$-1$       $1$

$X$

not "smooth"

```
...
vpslld   $20,    ...
vpshufd  $114,   ...
vmulpd   C1,     ...
vmulpd   C2,     ...
...
```

abstract using
"smooth" functions

$I_1$   $I_2$   $\cdots$   $I_n$

# Our Method: Overview

$P(x)$

hard to find

$-1$ ⊢————————————⊣ $1$

$X$

not "smooth"

```
...
vpslld  $20,    ...
vpshufd $114,   ...
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```

abstract using
"smooth" functions

➡

$I_1$ $I_2$ $\cdots$ $I_n$

```
...
vpslld  $20,    ...
vpshufd $114,   ...
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```

```
...
vpslld  $20,    ...
vpshufd $114,   ...
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```

$\cdots$

# Our Method: Overview

$P(x)$

hard to find

$-1$         $1$

$X$

not "smooth"

```
...
vpslld  $20,    ...
vpshufd $114,   ...
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```

abstract using "smooth" functions

$I_1$   $I_2$   $\cdots$   $I_n$
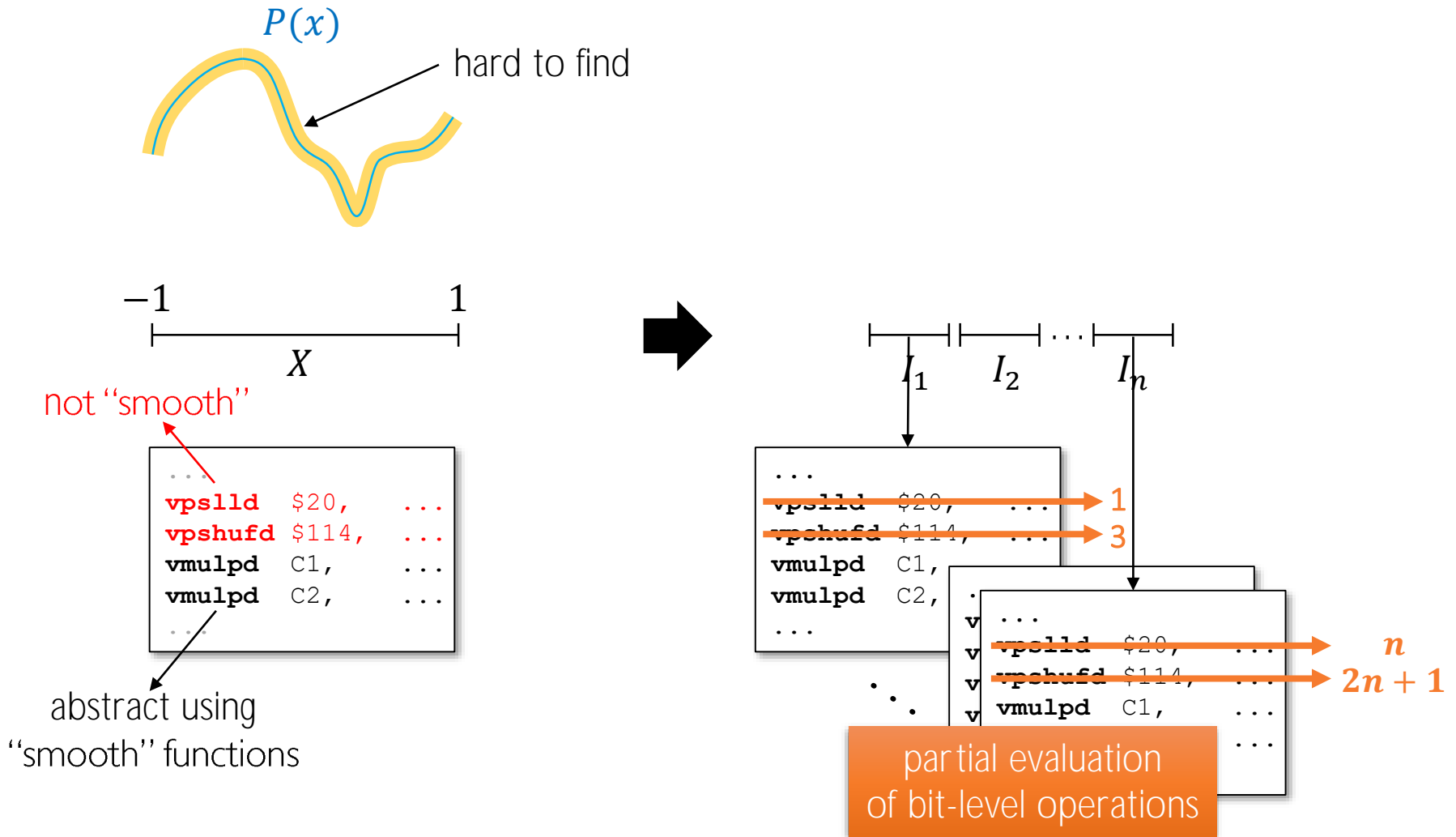
```
...
vpslld  $20,    ...    → 1
vpshufd $114,   ...    → 3
vmulpd  C1,     ...
vmulpd  C2,
...
```

```
...
vpslld  $20,    ...    → n
vpshufd $114,   ...    → 2n + 1
vmulpd  C1,     ...
...
```

partial evaluation of bit-level operations

# Our Method: Overview

$P(x)$

hard to find

$-1$         $1$

$X$

not "smooth"

```
...
vpslld  $20,    ...
vpshufd $114,   ...
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```

abstract using
"smooth" functions

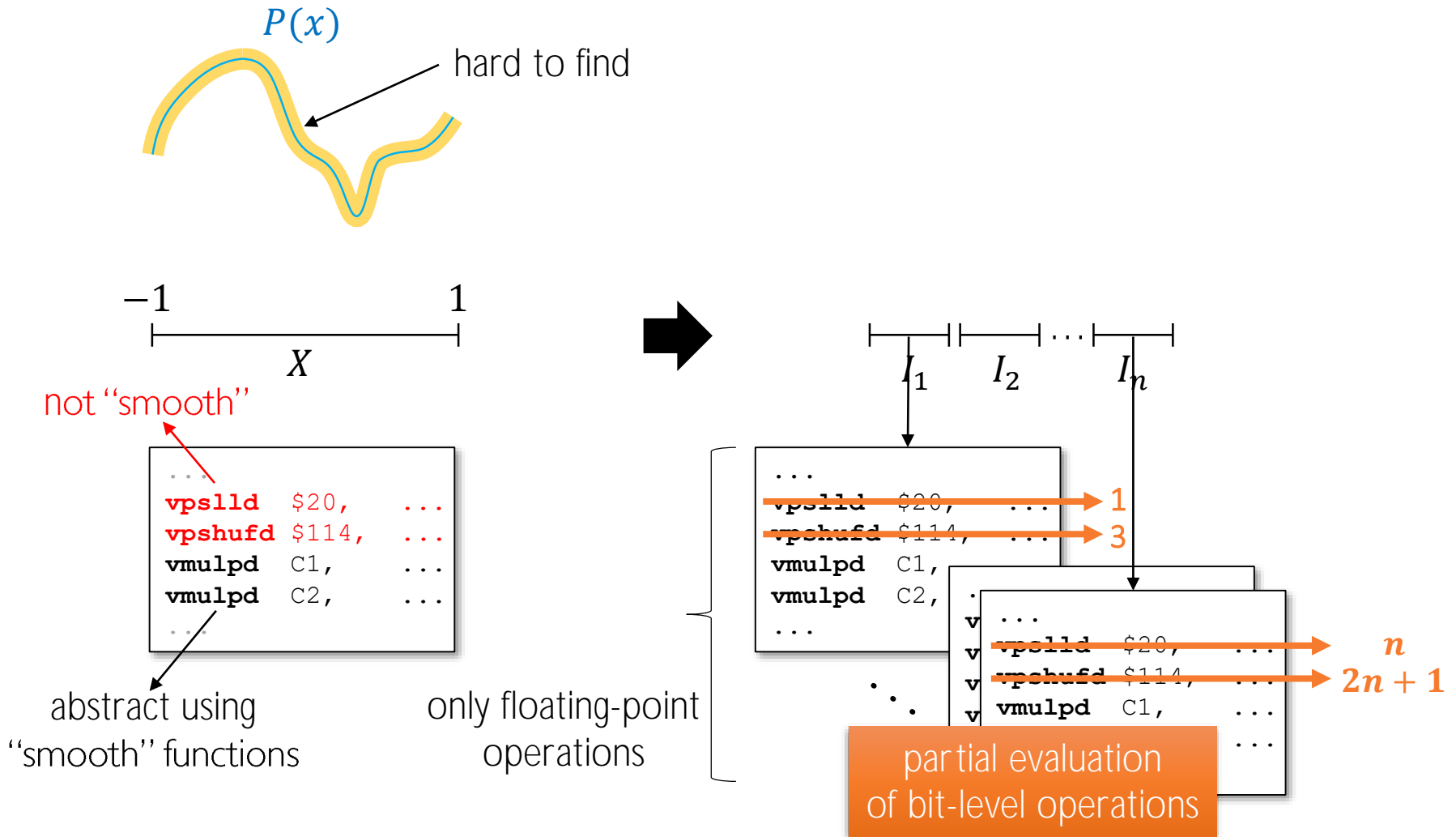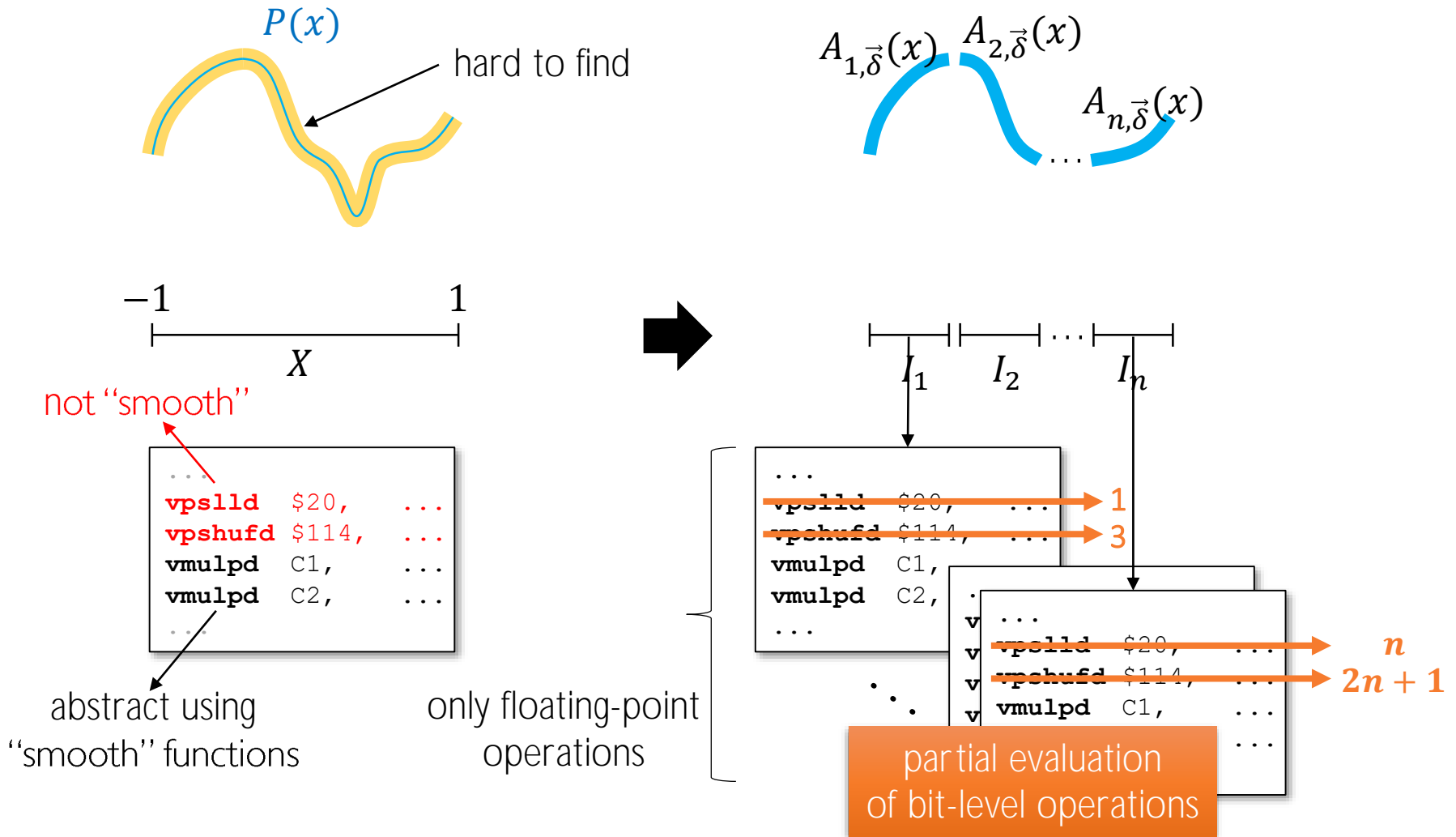only floating-point
operations

$I_1$    $I_2$   $\cdots$   $I_n$

```
...
vpslld  $20,    ...   → 1
vpshufd $114,   ...   → 3
vmulpd  C1,     ...
vmulpd  C2,     ...
...
```

```
...
vpslld  $20,    ...   → n
vpshufd $114,   ...   → 2n+1
vmulpd  C1,     ...
              ...
```

partial evaluation
of bit-level operations

# Our Method: Overview

$P(x)$

hard to find

$A_{1,\vec{\delta}}(x)$  $A_{2,\vec{\delta}}(x)$

$A_{n,\vec{\delta}}(x)$

$\cdots$

$-1$   $1$

$X$

$I_1$   $I_2$   $\cdots$   $I_n$

not "smooth"

```
...
vpslld   $20,    ...
vpshufd  $114,   ...
vmulpd   C1,     ...
vmulpd   C2,     ...
...
```

abstract using
"smooth" functions

only floating-point
operations

```
...
vpslld   $20,    ...   1
vpshufd  $114,   ...   3
vmulpd   C1,     ...
vmulpd   C2,     ...
...
```

```
...
vpslld   $20,    ...   n
vpshufd  $114,   ...   2n+1
vmulpd   C1,     ...
                       ...
```

partial evaluation
of bit-level operations

# Our Method: Overview

# Our Method: Overview

# Our Method: Overview

# Our Method: Overview



$A_{1,\vec{\delta}}(x)$  $A_{2,\vec{\delta}}(x)$

$A_{n,\vec{\delta}}(x)$

$\cdots$

$I_1$  $I_2$  $I_n$

```
...
vpslld   $20,  ...        1
vpshufd  $114, ...        3
vmulpd   C1,
vmulpd   C2,  .
...      v ...
         v vpslld   $20,    ...      n
         v vpshufd  $114,   ...    2n+1
         v vmulpd   C1,     ...
           ...
```

**partial evaluation of bit-level operations**

solve optimization problems

$$\max \left| \frac{f(x) - A_{1,\vec{\delta}}}{f(x)} \right| \quad \max \left| \frac{f(x) - A_{n,\vec{\delta}}(x)}{f(x)} \right|$$

$\cdots$

$I_1$  $I_2$  $I_n$

answer!

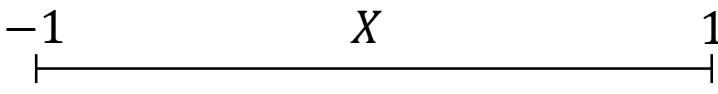# 2) Divide the Input Range

- Assume bit-level operations are used as well

# 2) Divide the Input Range

- Assume bit-level operations are used as well
- To handle bit-level operations, divide $X$ into intervals $I_k$,

  so that, on each $I_k$, we can statically know
  the result of each bit-level operation

# 2) Divide the Input Range

- Assume bit-level operations are used as well
- To handle bit-level operations, divide $X$ into intervals $I_k$,

  so that, on each $I_k$, we can statically know
  the result of each bit-level operation

- Example:

$$-1 \xrightarrow{\qquad\qquad X \qquad\qquad} 1$$
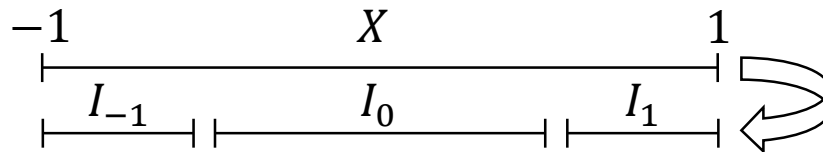
```
input x
y ← x ×f C
    (C=0x3ff71547652b82fe)
N ← round(y)
z ← int(N) +i 0x3ff
w ← z << 52
...
```

# 2) Divide the Input Range

- Assume bit-level operations are used as well
- To handle bit-level operations, divide $X$ into intervals $I_k$,

    so that, on each $I_k$, we can statically know
    the result of each bit-level operation

- Example:

$$-1 \qquad\qquad X \qquad\qquad 1$$

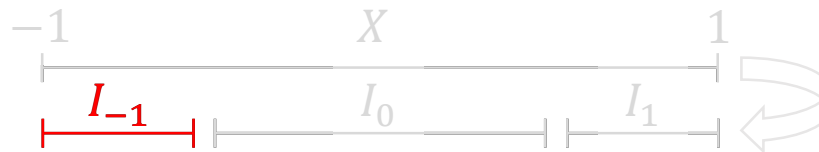$$I_{-1} \qquad\qquad I_0 \qquad\qquad I_1$$

```
input x
y ← x ×f C
     (C= 0x3ff71547652b82fe)
N ← round(y)
z ← int(N) +i 0x3ff
w ← z << 52
...
```

# 2) Divide the Input Range

- Assume bit-level operations are used as well
- To handle bit-level operations, divide $X$ into intervals $I_k$,

  so that, on each $I_k$, we can statically know
  the result of each bit-level operation

- Example:

$-1 \qquad\qquad X \qquad\qquad 1$

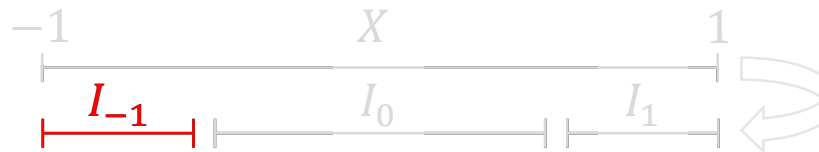$I_{-1}$ $\qquad\qquad I_0 \qquad\qquad I_1$

```
input x
y ← x ×f C
    (C= 0x3ff71547652b82fe)
N ← round(y)
z ← int(N) +i 0x3ff
w ← z << 52
...
```

# 2) Divide the Input Range

- Assume bit-level operations are used as well
- To handle bit-level operations, divide $X$ into intervals $I_k$,

  so that, on each $I_k$, we can statically know
  the result of each bit-level operation

- Example:

$$-1 \qquad\qquad X \qquad\qquad 1$$
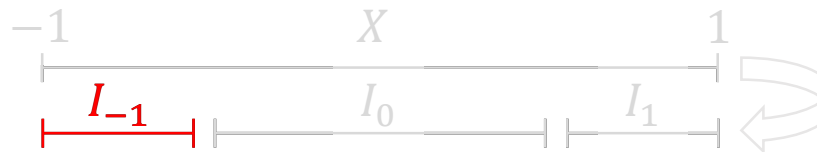
$$I_{-1} \qquad\qquad I_0 \qquad\qquad I_1$$

```
input x
y ← x ×f C
     (C=0x3ff71547652b82fe)
N ← round(y)      −1
z ← int(N) +i 0x3ff
w ← z << 52
...
```

# 2) Divide the Input Range

- Assume bit-level operations are used as well
- To handle bit-level operations, divide $X$ into intervals $I_k$,

  so that, on each $I_k$, we can statically know
  the result of each bit-level operation

- Example:



```
input x
y ← x ×f C
     (C = 0x3ff71547652b82fe)
N ← round(y)      −1
z ← int(N) +i 0x3ff
w ← z << 52
...
```
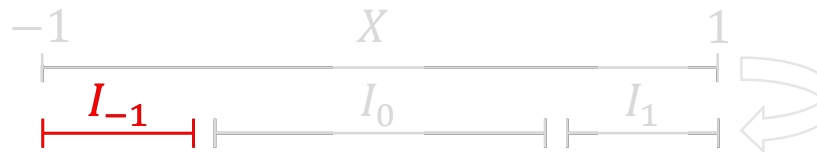
partial evaluation →

```
input x
y ← x ×f C
     (C = 0x3ff71547652b82fe)
N ← −1
z ← 1022
w ← 0.5
...
```

# 2) Divide the Input Range

- Assume bit-level operations are used as well
- To handle bit-level operations, divide $X$ into intervals $I_k$,

so that, on each $I_k$, we can statically know
the result of each bit-level operation

- Example:



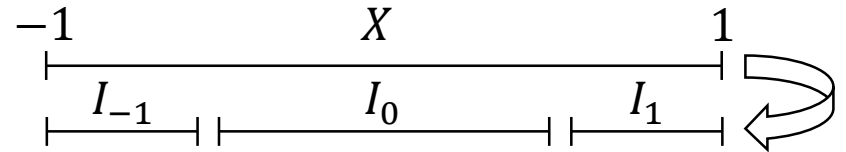| input x | | input x |
|---|---|---|
| y ← x ×f C | | y ← x ×f C |
| (C = 0x3ff71547652b82fe) | partial evaluation | (C = 0x3ff71547652b82fe) |
| N ← ~~round(y)~~ −1 | | N ← −1 |

Only floating-point operations are left
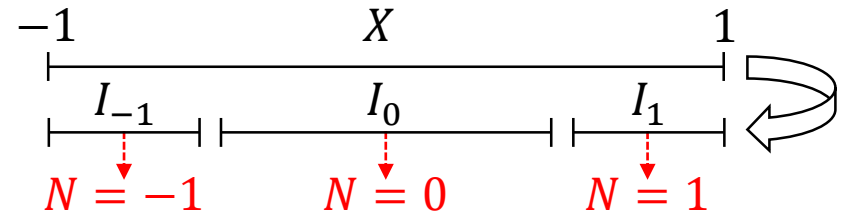→ Can compute $A_{\vec{\delta}}(x)$ on each $I_k$

# 2) Divide the Input Range

- How to find such intervals?

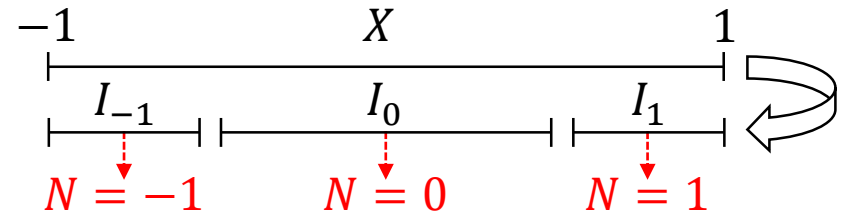# 2) Divide the Input Range

- How to find such intervals?

# 2) Divide the Input Range

- How to find such intervals?
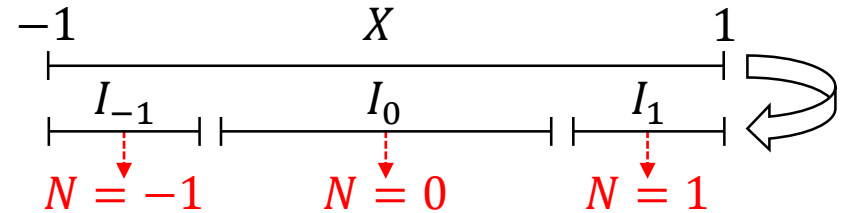
# 2) Divide the Input Range

- How to find such intervals?
  - Use symbolic abstractions

$-1 \quad\quad\quad\quad\quad X \quad\quad\quad\quad\quad 1$

$I_{-1} \quad\quad\quad I_0 \quad\quad\quad I_1$

$N = -1 \quad\quad N = 0 \quad\quad N = 1$

# 2) Divide the Input Range

- How to find such intervals?
  - Use symbolic abstractions

$$-1 \qquad\qquad X \qquad\qquad 1$$

$$I_{-1} \qquad\qquad I_0 \qquad\qquad I_1$$

$$N = -1 \qquad N = 0 \qquad N = 1$$

- Example:
  - $N = \text{round}(x \times_{\text{f}} \text{C})$

# 2) Divide the Input Range

- How to find such intervals?
  - Use symbolic abstractions

$$-1 \qquad\qquad X \qquad\qquad\qquad 1$$

$$I_{-1} \qquad\qquad I_0 \qquad\qquad I_1$$

$$N = -1 \qquad N = 0 \qquad N = 1$$

- Example:
  - $N = \text{round}(x \times_{\text{f}} \text{C})$
  - (symbolic abstraction of $x \times_{\text{f}} \text{C}$) = $(x \times \text{C})(1 + \delta)$

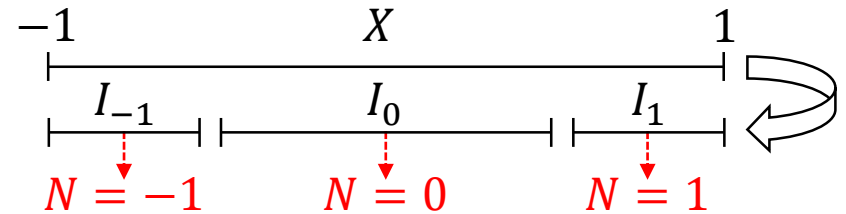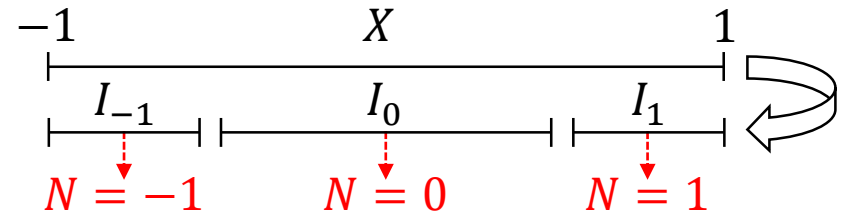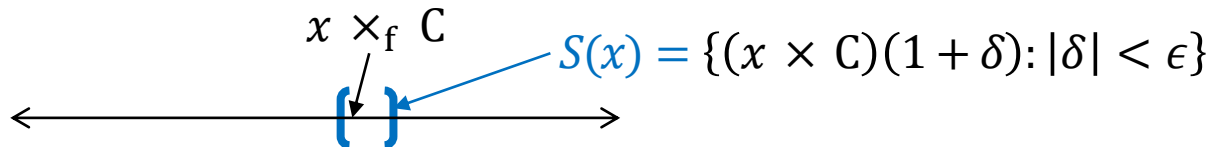# 2) Divide the Input Range

- How to find such intervals?
  - Use symbolic abstractions



- Example:
  - $N = \text{round}(x \times_f \text{C})$
  - (symbolic abstraction of $x \times_f \text{C}$) $= (x \times \text{C})(1 + \delta)$



$$S(x) = \{(x \times \text{C})(1 + \delta) : |\delta| < \epsilon\}$$

# 2) Divide the Input Range

- How to find such intervals?
  - Use symbolic abstractions



$-1$      $X$      $1$

$I_{-1}$    $I_0$    $I_1$

$N = -1$    $N = 0$    $N = 1$

- Example:
  - $N = \mathrm{round}(x \times_{\mathrm{f}} \mathrm{C})$
  - (symbolic abstraction of $x \times_{\mathrm{f}} \mathrm{C}$) $= (x \times \mathrm{C})(1 + \delta)$

$x \times_{\mathrm{f}} \mathrm{C}$

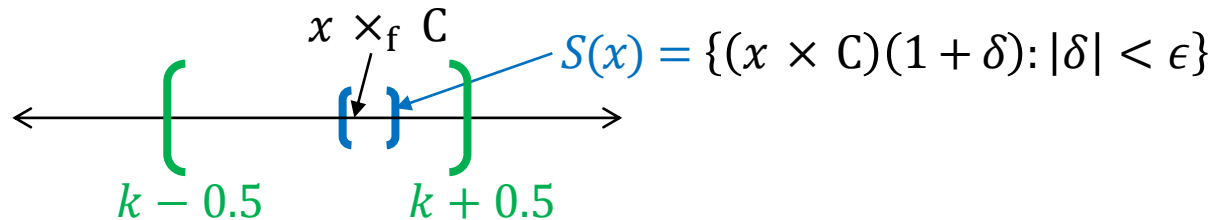$S(x) = \{(x \times \mathrm{C})(1 + \delta) : |\delta| < \epsilon\}$

$k - 0.5$      $k + 0.5$

# 2) Divide the Input Range

- How to find such intervals?
  - Use symbolic abstractions



- Example:
  - $N = \mathrm{round}(x \times_{\mathrm{f}} \mathrm{C})$
  - (symbolic abstraction of $x \times_{\mathrm{f}} \mathrm{C}$) $= (x \times \mathrm{C})(1 + \delta)$



$$x \times_{\mathrm{f}} \mathrm{C}$$

$$S(x) = \{(x \times \mathrm{C})(1 + \delta): |\delta| < \epsilon\}$$

$$k - 0.5 \qquad k + 0.5 \qquad \Rightarrow \quad N = k$$

# 2) Divide the Input Range

- How to find such intervals?
  - Use symbolic abstractions

$$-1 \qquad\qquad X \qquad\qquad 1$$
$$I_{-1} \qquad\qquad I_0 \qquad\qquad I_1$$
$$N = -1 \qquad N = 0 \qquad N = 1$$

- Example:
  - $N = \mathrm{round}(x \times_{\mathrm{f}} \mathrm{C})$
  - (symbolic abstraction of $x \times_{\mathrm{f}} \mathrm{C}$) = $(x \times \mathrm{C})(1 + \delta)$

$$x \times_{\mathrm{f}} \mathrm{C}$$
$$S(x) = \{(x \times \mathrm{C})(1 + \delta) : |\delta| < \epsilon\}$$
$$k - 0.5 \qquad k + 0.5 \qquad\qquad N = k$$

- Let $I_k = $ largest interval contained in
$$\{x \in X : S(x) \subset (k - 0.5, k + 0.5)\}$$

# 2) Divide the Input Range

- How to find such intervals?
  - Use symbolic abstractions

$-1$  $X$  $1$

$I_{-1}$  $I_0$  $I_1$

$N = -1$  $N = 0$  $N = 1$

- Example:
  - $N = \text{round}(x \times_{\text{f}} \text{C})$
  - (symbolic abstraction of $x \times_{\text{f}} \text{C}$) $= (x \times \text{C})(1 + \delta)$

$x \times_{\text{f}} \text{C}$

$S(x) = \{(x \times \text{C})(1 + \delta) : |\delta| < \epsilon\}$

$k - 0.5$  $k + 0.5$  $N = k$

  - Let $I_k = $ largest interval contained in
$$\{x \in X : S(x) \subset (k - 0.5, k + 0.5)\}$$
  - Then $N$ is evaluated to $k$ for every input in $I_k$

# 3) Compute a Bound on Precision Loss

- Precision loss on each interval $I_k$
  - Let $A_{\vec{\delta}}(x)$ be a symbolic abstraction on $I_k$
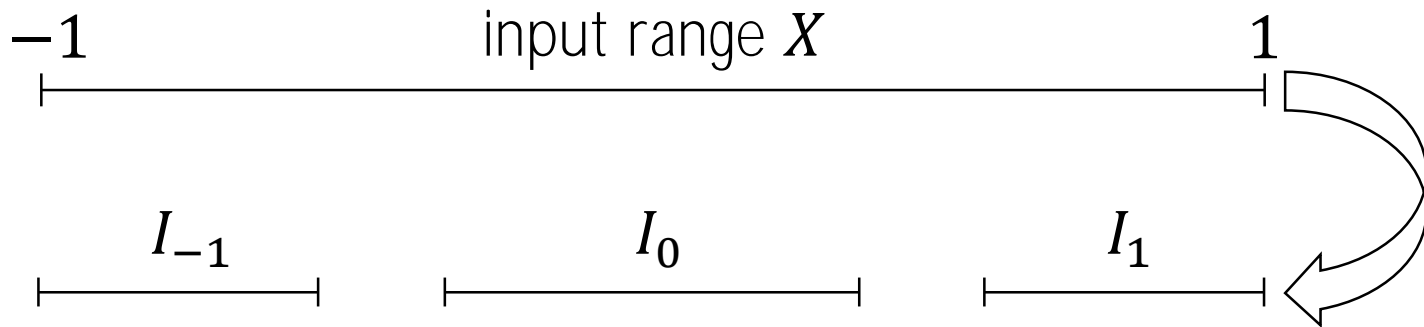
# 3) Compute a Bound on Precision Loss

- Precision loss on each interval $I_k$
  - Let $A_{\vec{\delta}}(x)$ be a symbolic abstraction on $I_k$
  - Analytical optimization:

$$\max_{x \in I_k,\ |\delta_i| < \epsilon} \left| \frac{e^x - A_{\vec{\delta}}(x)}{e^x} \right|$$

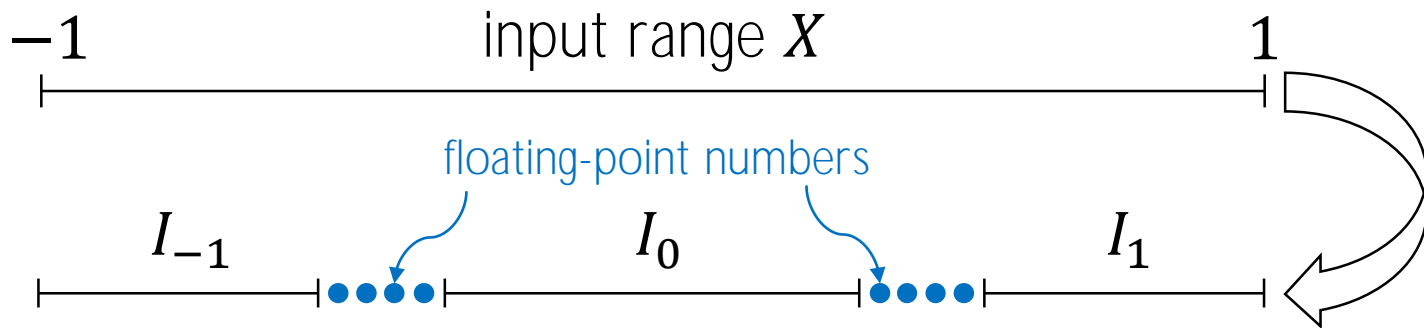  - Use Mathematica to solve optimization problems analytically

# Are We Done?

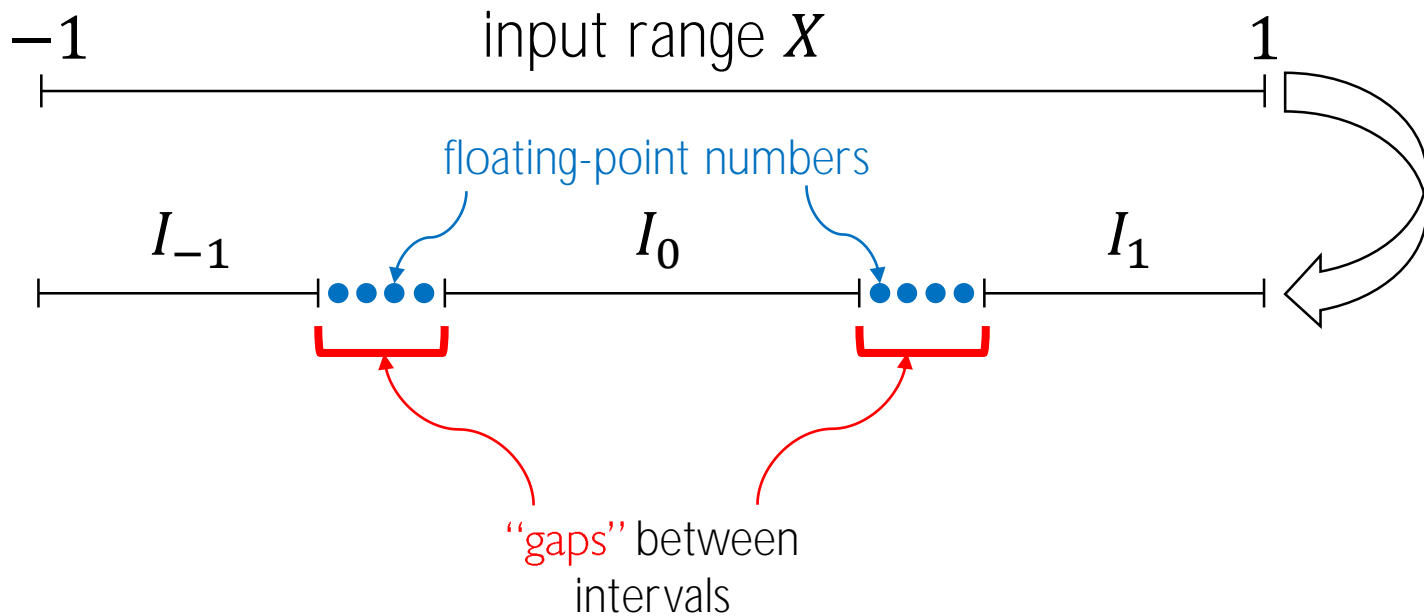- No. The constructed intervals <span style="color:red">do not</span> cover $X$ in general

$$-1 \qquad\qquad \text{input range } X \qquad\qquad 1$$

$$I_{-1} \qquad\qquad I_0 \qquad\qquad I_1$$

# Are We Done?

- No. The constructed intervals <span style="color:red">do not</span> cover $X$ in general

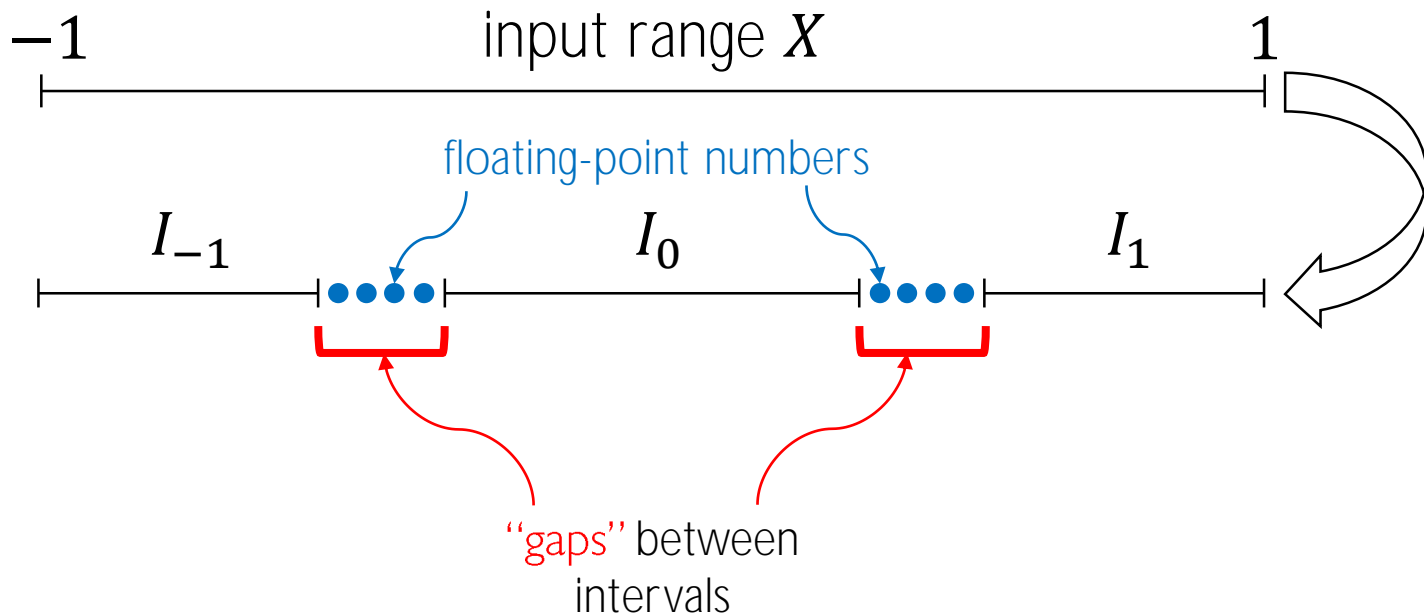# Are We Done?

- No. The constructed intervals <span style="color:red">do not</span> cover $X$ in general
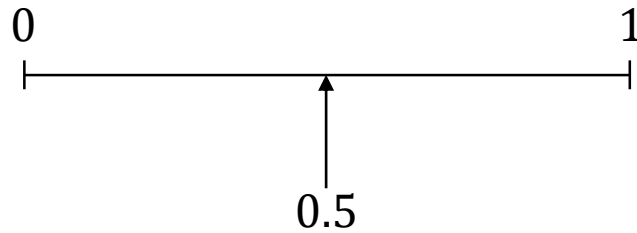
# Are We Done?

- No. The constructed intervals <span style="color:red">do not</span> cover $X$ in general
    - Because we made sound approximations

# Are We Done?

- Example: $N = \text{round}(x \times_\text{f} C)$

$[\ ]$ : abstraction of $x \times_\text{f} C$

$$
\begin{array}{l}
0 \qquad\qquad\qquad\qquad\qquad 1 \\
\vdash\!\!\!\!\!\!\!\!-\!\!\!-\!\!\!-\!\!\!-\!\!\!\uparrow\!\!\!-\!\!\!-\!\!\!-\!\!\!-\!\!\!\dashv \\
\qquad\qquad\qquad 0.5
\end{array}
$$

# Are We Done?

- Example: $N = \mathrm{round}(x \times_{\mathrm{f}} C)$

$[\ ]$ : abstraction of $x \times_{\mathrm{f}} C$

# Are We Done?

- Example: $N = \text{round}(x \times_f C)$

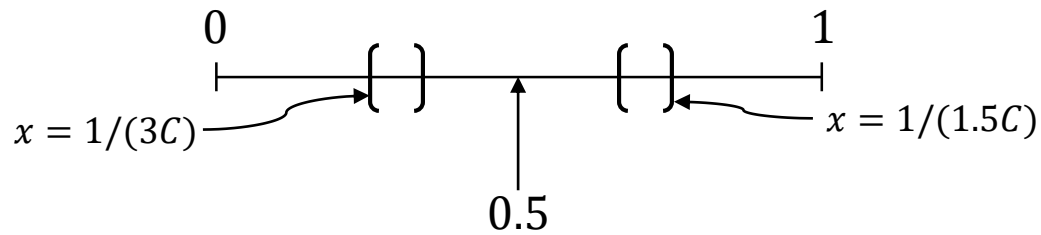$\big[\ \big]$ : abstraction of $x \times_f C$
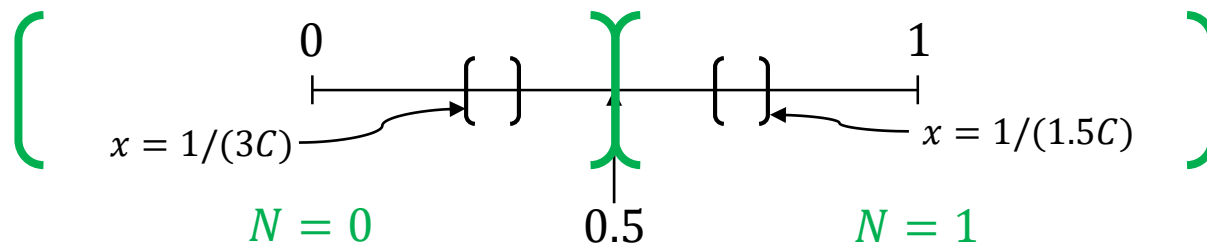
# Are We Done?
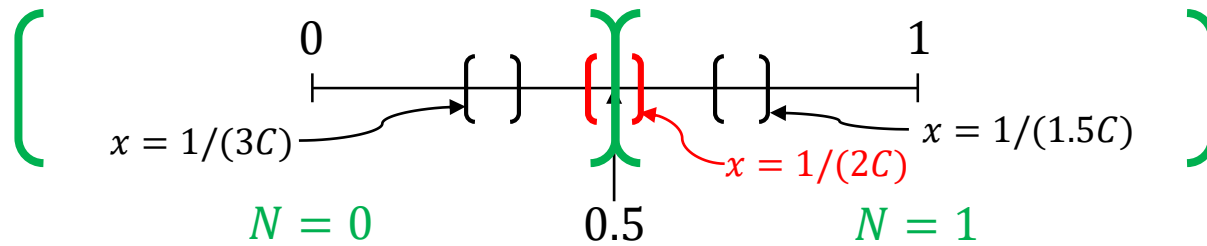
- Example: $N = \text{round}(x \times_f C)$

$\Big[\ \Big]$ : abstraction of $x \times_f C$



$x = 1/(3C)$    $N = 0$

$x = 1/(2C)$

$x = 1/(1.5C)$    $N = 1$

$0$    $0.5$    $1$

# Are We Done?

- Example: $N = \mathrm{round}(x \times_{\mathrm{f}} C)$



$[\ ]$: abstraction of $x \times_{\mathrm{f}} C$

$x \times_{\mathrm{f}} C$

???

0

1

$x = 1/(3C)$

$x = 1/(1.5C)$

$x = 1/(2C)$

$N = 0$

0.5

$N = 1$

# Are We Done?

- Example: $N = \text{round}(x \times_{\text{f}} C)$



$x \times_{\text{f}} C$

$[\ ]$: abstraction of $x \times_{\text{f}} C$

???

0             1

$x = 1/(3C)$     $x = 1/(1.5C)$

$x = 1/(2C)$

$N = 0$      0.5      $N = 1$

For $x = \dfrac{1}{2C}$, we can't statically know if $N$ would be 0 or 1

97

# Are We Done?

- Example: $N = \mathrm{round}(x \times_{\mathrm{f}} C)$



$[\ ]$ : abstraction of $x \times_{\mathrm{f}} C$

$x \times_{\mathrm{f}} C$

???

0

1

$x = 1/(3C)$

$x = 1/(2C)$

$x = 1/(1.5C)$

$N = 0$

0.5

$N = 1$

For $x = \frac{1}{2C}$, we can't statically know if $N$ would be 0 or 1

- Let $H = \{$floating-point numbers in the "gaps"$\}$
  - We observe that $|H|$ is small in experiment

# 3) Compute a Bound on Precision Loss

- Precision loss on each interval $I_k$
  - Let $A_{\vec{\delta}}(x)$ be a symbolic abstraction on $I_k$
  - Analytical optimization:

$$\max_{x \in I_k,\, |\delta_i| < \epsilon} \left| \frac{e^x - A_{\vec{\delta}}(x)}{e^x} \right|$$

  - Use Mathematica to solve optimization problems analytically

- Precision loss on $H$
  - For each $x \in H$, obtain $P(x)$ by executing the binary
  - Brute force:

$$\max_{x \in H} \left| \frac{e^x - P(x)}{e^x} \right|$$

  - Use Mathematica to compute $e^x$ and precision loss exactly

# 3) Compute a Bound on Precision Loss

- Precision loss on each interval $I_k$
  - Let $A_{\vec{\delta}}(x)$ be a symbolic abstraction on $I_k$
  - Analytical optimization:

$$\max_{x \in I_k,\ |\delta_i| < \epsilon} \left| \frac{e^x - A_{\vec{\delta}}(x)}{e^x} \right|$$

take maximum
$\rightarrow$ answer!

  - Use Mathematica to solve optimization problems analytically

- Precision loss on $H$
  - For each $x \in H$, obtain $P(x)$ by executing the binary
  - Brute force:

$$\max_{x \in H} \left| \frac{e^x - P(x)}{e^x} \right|$$

  - Use Mathematica to compute $e^x$ and precision loss exactly

# Case Studies

# Settings

- Benchmarks
  - `exp`: from S3D (a combustion simulation engine)
  - `sin, log`: from Intel's `<math.h>`

- Measures of precision loss
  - Relative error: $\mathbf{RelErr}(a, b) = \left| \frac{a-b}{a} \right|$
  - ULP error:
    - Rounding errors of numeric libraries are typically measured by ULPs

# Settings

- Benchmarks
  - `exp`: from S3D (a combustion simulation engine)
  - `sin, log`: from Intel's `<math.h>`

- Measures of precision loss
  - Relative error: $\mathbf{RelErr}(a, b) = \left| \frac{a-b}{a} \right|$
  - <span style="color:red">ULP error:</span>
    - Rounding errors of numeric libraries are typically measured by ULPs
    - $\mathbf{ULPErr}(a, b) = (\# \text{ of floating-point numbers between } a \text{ and } b)$
    - Example:



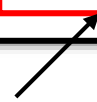  - $\mathbf{ULPErr}(a, b) \leq 2 \cdot \mathbf{RelErr}(a, b)/\epsilon$

# Results

| | Interval | Bound on ULP error | # of intervals | # of $\delta$'s | Size of "gaps" |
|---|---|---|---|---|---|
| exp | $[-4, 4]$ | 14 | 13 | 29 | 36 |
| sin | $\left[-\dfrac{\pi}{2}, \dfrac{\pi}{2}\right]$ | 9 | 33 | 53 | 110 |
| log | $(0,4) \setminus \left[\dfrac{4095}{4096}, 1\right)$ | 21 | $2^{21}$ | 25 | 0 |
| | $\left[\dfrac{4095}{4096}, 1\right)$ | $1 \times 10^{14}$ | 1 | 25 | 0 |

# Results

| | Interval | Bound on ULP error | # of intervals | # of $\delta$'s | Size of "gaps" |
|---|---|---|---|---|---|
| exp | $[-4, 4]$ | 14 | 13 | 29 | 36 |
| sin | $\left[-\dfrac{\pi}{2}, \dfrac{\pi}{2}\right]$ | 9 | 33 | 53 | 110 |
| log | $(0,4) \setminus \left[\dfrac{4095}{4096}, 1\right)$ | 21 | $2^{21}$ | 25 | 0 |
| | $\left[\dfrac{4095}{4096}, 1\right)$ | $1 \times 10^{14}$ | 1 | 25 | 0 |

best illustrates
the power of our method

# Results: `sin, log`

sin        log

$10^{14}$

y-axis:
ULP error

x-axis: input value

bounds on the intervals
errors on the "gaps"

# Results: `sin, log`



sin

log

$10^{14}$

y-axis: ULP error

x-axis: input value

- bounds on the intervals
- errors on the "gaps"
- Intel's bound

# Limitations of Our Method

- May construct a <span style="color:red">large</span> number of intervals
  - Example: `0x5fe6eb50c7b537a9 - (x >> 1)`
  - For this example, our method constructs $2^{63}$ intervals

# Limitations of Our Method

- May construct a large number of intervals
  - Example: `0x5fe6eb50c7b537a9 – (x >> 1)`
  - For this example, our method constructs $2^{63}$ intervals


- May produce loose error bounds
  - Example: $10^{14}$ ULPs for `log` on $\left[\frac{4095}{4096}, 1\right)$
  - Sometimes $(1 + \epsilon)$ property provides an imprecise abstraction
  - Proving a precise error bound requires more sophisticated error analysis beyond $(1 + \epsilon)$ property
  - Our recent result: 6 ULPs for for `log` on $(0,4)$

# Summary

- First systematic method for verifying binaries
  that mix floating-point and bit-level operations

- Use abstraction, analytical optimization, and testing

- Directly applicable to highly optimized binaries
  of transcendental functions

# Questions?