

University of Perpignan Via Domitia
LAMPS Laboratory

Fast and Efficient Bit-Level Precision Tuning

Ph.D. Defense

Dorra BEN KHALIFA

Jury Members

Directors: Matthieu MARTEL
Assalé ADJE

Reviewers: Ganesh GOPALAKRISHNAN
Eric GOUBAULT
David MONNIAUX

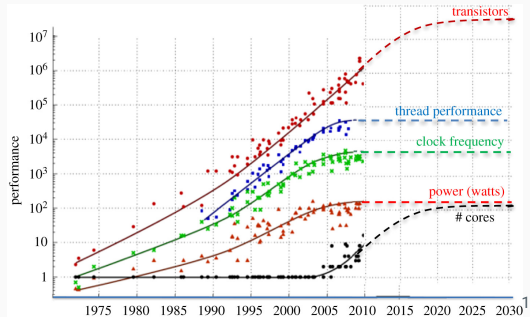
Examiners: Eva DARULOVA
Philippe LANGLOIS
Laura TITOLO

November 29, 2021



Introduction

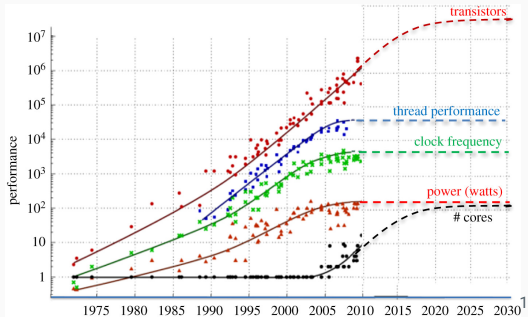
- **SW** used to solve more & more complex tasks
- Thanks to **HW** whose performance double every 18 months \implies **Moore's law in 1965**



¹Photo credits from <https://royalsocietypublishing.org/doi/10.1098/rsta.2019.0061>

Introduction

- SW used to solve more & more complex tasks
- Thanks to HW whose performance double every 18 months \implies **Moore's law in 1965**



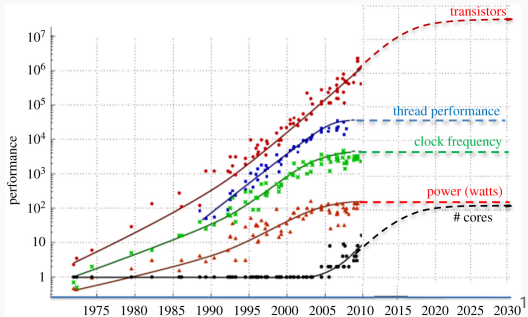
Issue

- ☞ Less perspectives to further performance improvements after 2025
- ☞ End of Moore's law

¹Photo credits from <https://royalsocietypublishing.org/doi/10.1098/rsta.2019.0061>

Introduction

- SW used to solve more & more complex tasks
- Thanks to HW whose performance double every 18 months \implies **Moore's law in 1965**



Issue

- ☞ Less perspectives to further performance improvements after 2025
- ☞ End of Moore's law

☞ What is the future of computing?

¹Photo credits from <https://royalsocietypublishing.org/doi/10.1098/rsta.2019.0061>

Motivation

- Developers use **highest precision** available (e.g. 64 bits, 128 bits) → increases running time, memory footprint, energy consumption



Motivation

- Developers use **highest precision** available (e.g. 64 bits, 128 bits) → increases running time, memory footprint, energy consumption



- Sometimes variables in **lower precision** (e.g. 16 bits, 32 bits) is enough

Motivation

- Developers use **highest precision** available (e.g. 64 bits, 128 bits) → increases running time, memory footprint, energy consumption



- Sometimes variables in **lower precision** (e.g. 16 bits, 32 bits) is enough
- Manually changing for optimized precision is fastidious

Motivation

- Developers use **highest precision** available (e.g. 64 bits, 128 bits) → increases running time, memory footprint, energy consumption



- Sometimes variables in **lower precision** (e.g. 16 bits, 32 bits) is enough
- Manually changing for optimized precision is fastidious
- **Precision Tuning technique**
 - ☞ Not a Find-and-Replace button like in a text editor
 - ☞ More complex task with program semantics analysis

Motivation

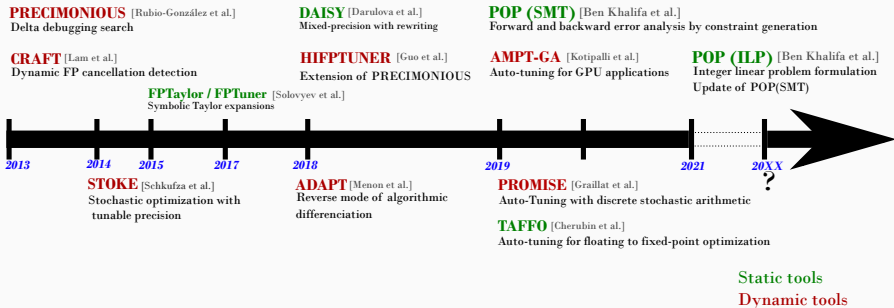
- Developers use **highest precision** available (e.g. 64 bits, 128 bits) → increases running time, memory footprint, energy consumption



- Sometimes variables in **lower precision** (e.g. 16 bits, 32 bits) is enough
- Manually changing for optimized precision is fastidious
- **Precision Tuning technique**
 - ↳ Not a Find-and-Replace button like in a text editor
 - ↳ More complex task with program semantics analysis

Goal

- ↳ Develop automated techniques for precision tuning
- ↳ Best trade-offs between performance & user requirements



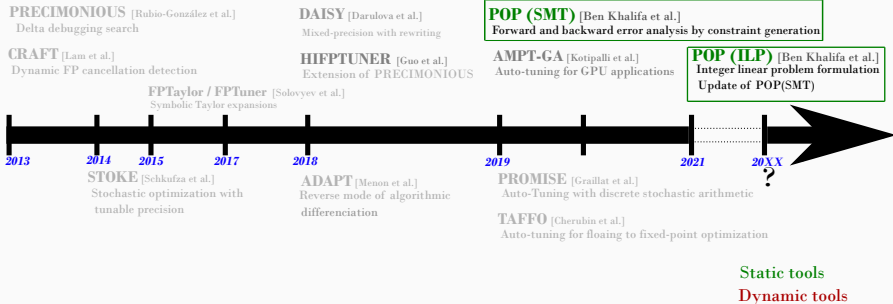
Weaknesses:

- Trial-and-Error strategy
- Search algorithm-based tools (**CRAFT**, **PRECIMONIOUS**, ...)
- Static tools limited to straight-line programs (**Daisy**, **FPTUNER**, ...) or uniform precision (**Rosa**,...)

👉 How to go beyond?

²A more exhaustive survey is given in Chapter 3 of the dissertation

Precision Tuning A Survey²

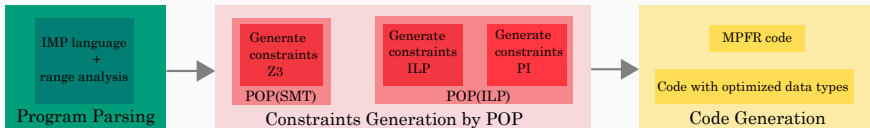


POP strengths:

- ✓ No trial-and-Error strategy
- ✓ Finds directly the minimal number of bits needed
- ✓ Complexity does not increase as number of data types increases
- ✓ Loops and conditionals

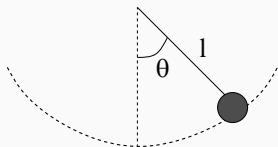
²A more exhaustive survey is given in Chapter 3 of the dissertation

POP in One Slide



- Open-source tool
- Implementation in Java
- \approx **10 000** lines of code
- ANTLR, Z3, GLPK libraries
- Bit-level, IEEE754, fixed and MPFR datatypes
- Publications [FTSCS'19, IINTEC'19, IoTaIS'20, ICICT'21, ICCSA'21, SAS'21]



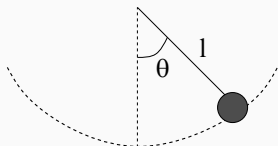


Pendulum ($\theta = \frac{\pi}{4}$)

Second order differential equation

$$(E1): m \cdot l \cdot \frac{d^2\theta}{dt^2} = -m \cdot g \cdot \sin \theta$$

$$(E1) \Leftrightarrow (E2): \frac{dy_1}{dt} = y_2 \quad \text{and} \quad \frac{dy_2}{dt} = -\frac{g}{l} \cdot \sin y_1$$



Pendulum ($\theta = \frac{\pi}{4}$)

Second order differential equation

$$(E1): m \cdot l \cdot \frac{d^2\theta}{dt^2} = -m \cdot g \cdot \sin \theta$$

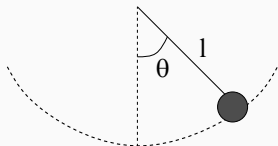
$$(E1) \Leftrightarrow (E2): \frac{dy_1}{dt} = y_2 \quad \text{and} \quad \frac{dy_2}{dt} = -\frac{g}{l} \cdot \sin y_1$$

```

1  gl1 = 9.81l0; ll3 = 0.5l2;
2  y1l5 = 0.785398l4;
3  y2l7 = 0.785398l6;
4  hl9 = 0.1l8; tl11 = 0.0l10;
5  while (tl13 <l15 10.0l14)l59 {
6    y1newl24 = y1l17 +l23 y2l19 *l22 hl21;
7    aux1l28 = sin(y1l26)l27;
8    aux2l40 = aux1l30 *l39 hl32
9      *l38 gl34 /l37 ll36;
10   y2newl46 = y2l42 -l45 aux2l44;
11   tl52 = tl48 +l51 hl50;
12   y1l55 = y1newl54;
13   y2l58 = y2newl57; };
14  require_nsb(y2, 20)l61;
  
```

POP Label File

Running Example ... to Tuning



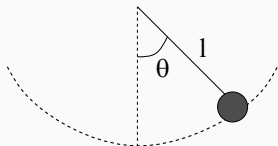
Pendulum ($\theta = \frac{\pi}{4}$)

Second order differential equation

$$(E1): m \cdot l \cdot \frac{d^2\theta}{dt^2} = -m \cdot g \cdot \sin \theta$$

$$(E1) \Leftrightarrow (E2): \frac{dy_1}{dt} = y_2 \quad \text{and} \quad \frac{dy_2}{dt} = -\frac{g}{l} \cdot \sin y_1$$

Running Example ... to Tuning



Pendulum ($\theta = \frac{\pi}{4}$)

Second order differential equation

$$(E1): m \cdot l \cdot \frac{d^2\theta}{dt^2} = -m \cdot g \cdot \sin \theta$$

$$(E1) \Leftrightarrow (E2): \frac{dy_1}{dt} = y_2 \quad \text{and} \quad \frac{dy_2}{dt} = -\frac{g}{l} \cdot \sin y_1$$

```
1 g|20| = 9.81|20|; l|20| = 1.5|20|;
2 y1|29| = 0.785398|29|;
3 y2|21| = 0.0|21|;
4 h|21| = 0.1|21|; t|21| = 0.0|21|;
5 while (t<1.0) {
6   y1new|20| = y1|21| +|20| y2|21|*|22| h|21|;
7   aux1|20| = sin(y1|29|)|20|;
8   aux2|20| = aux1|19| *|20| h|18|*|19|g|17| /
9     |18||17|;
10  y2new|20| = y2|21| -|20| aux2|18|;
11  t|20| = t|21| +|20| h|17|;
12  y1|20| = y1new|20|;
13  y2|20| = y2new|20|;
14  };
require_nsb(y2, 20);
```

POP Output File

- 1 Preliminary Elements**

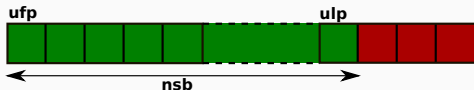
- 2 **Constraint Generation**

 -
 -
- 3 **Code Generation**

 -
- 4 **Conclusion and Perspectives**

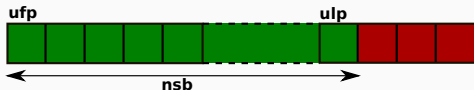
Preliminary Elements ufp, nsb, ulp





- The **unit in the first place** of a real number x

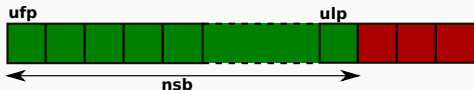
$$\text{ufp}(x) = \begin{cases} \min\{i \in \mathbb{Z} : 2^{i+1} > |x|\} = \lfloor \log_2(|x|) \rfloor & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}$$



- The **unit in the first place** of a real number x

$$\text{ufp}(x) = \begin{cases} \min\{i \in \mathbb{Z} : 2^{i+1} > |x|\} = \lfloor \log_2(|x|) \rfloor & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}$$

- The **number of significant bits** nsb of x
 - \hat{x} : approximation of x in finite precision
 - **Absolute error**: $\varepsilon_x \leq 2^{\text{ufp}(x) - \text{nsb}(x) + 1}$



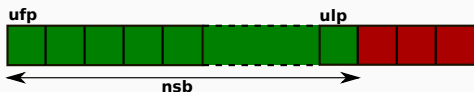
- The **unit in the first place** of a real number x

$$\text{ufp}(x) = \begin{cases} \min\{i \in \mathbb{Z} : 2^{i+1} > |x|\} = \lfloor \log_2(|x|) \rfloor & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}$$

- The **number of significant bits** nsb of x
 - \hat{x} : approximation of x in finite precision
 - **Absolute error**: $\varepsilon_x \leq 2^{\text{ufp}(x) - \text{nsb}(x) + 1}$

- The **unit in the last place** of x

$$\text{ulp}(x) = \text{ufp}(x) - \text{nsb}(x) + 1$$



- The **unit in the first place** of a real number x

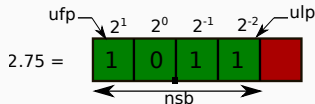
$$\text{ufp}(x) = \begin{cases} \min\{i \in \mathbb{Z} : 2^{i+1} > |x|\} = \lfloor \log_2(|x|) \rfloor & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}$$

- The **number of significant bits** nsb of x

- \hat{x} : approximation of x in finite precision
- Absolute error**: $\varepsilon_x \leq 2^{\text{ufp}(x) - \text{nsb}(x) + 1}$

- The **unit in the last place** of x

$$\text{ulp}(x) = \text{ufp}(x) - \text{nsb}(x) + 1$$



- $\text{ufp}(2.75) = 1$
- $\text{nsb}(2.75) = 4$
- $\text{ulp}(2.75) = -2$

```
1 [...]
2  $z^{\ell_2} = x^{\ell_0} + y^{\ell_1};$ 
3 require_nsb(z,18) $\ell_3$ ;
```

```
1 [...]
2  $z^{\ell_2} = x^{\ell_0} + y^{\ell_1};$ 
3  $\text{require\_nsb}(z, 18)^{\ell_3};$ 
```

- Preliminary range analysis for all the program variables

example $x \in [0.2, 0.7], y \in [0.6, 0.8], z \in [0.8, 1.5]$


```
1 [...]
2  $z^{\ell_2} = x^{\ell_0} + y^{\ell_1};$ 
3  $\text{require\_nsb}(z, 18)^{\ell_3};$ 
```

- Preliminary range analysis for all the program variables

example $x \in [0.2, 0.7], y \in [0.6, 0.8], z \in [0.8, 1.5]$

- Pre-computation of the ufp of each range

example $\text{ufp}(0.2) = -3, \text{ufp}(0.7) = -1 \implies \text{ufp}(x) = -1, \text{ufp}(y) = -1, \text{ufp}(z) = 0$

```

1 [...]
2  $z^{\ell_2} = x^{\ell_0} + y^{\ell_1};$ 
3 require_nsb(z, 18) $\ell_3$ ;

```

- Preliminary range analysis for all the program variables

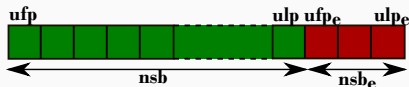
example $x \in [0.2, 0.7], y \in [0.6, 0.8], z \in [0.8, 1.5]$

- Pre-computation of the ufp of each range

example $ufp(0.2) = -3, ufp(0.7) = -1 \implies ufp(x) = -1, ufp(y) = -1, ufp(z) = 0$

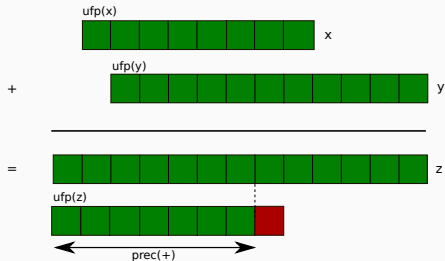
- **Goal: Computation of the minimal nsb**

- $nsb(z) = 18$ // thanks to the annotation
- $nsb(x) = ufp(x) - ufp(\varepsilon(x))$
- $nsb(x) = -1 - ufp(\varepsilon(x))$?



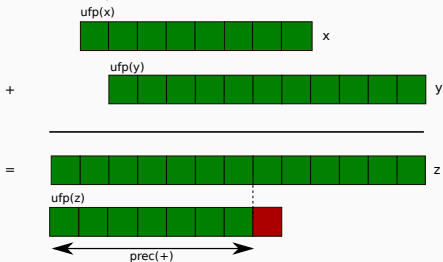
👉 How to compute errors?

Case 1: Initially values x and y are exact



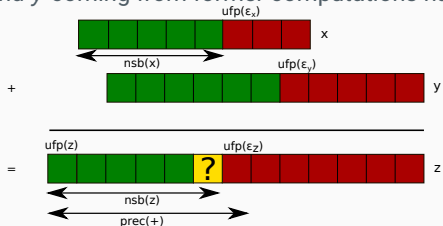
- **Truncation error:** $\varepsilon_+ \leq 2^{u_{fp}(z) - prec(+)}$

Case 1: Initially values x and y are exact



- **Truncation error:** $\varepsilon_+ \leq 2^{ufp(z) - prec(+)}$

Case 2: Values x and y coming from former computations have errors

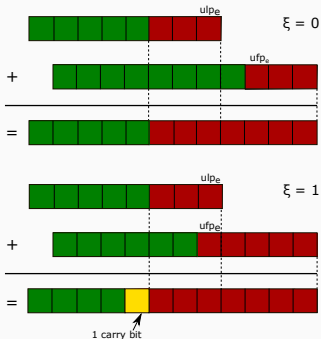


- $ufp(\varepsilon_z)$ depends on carry

✗ Before this work, very costly over-approximation in [NFM'17]

✗ Before this work, very costly over-approximation in [NFM'17]

✓ Optimized carry bit function



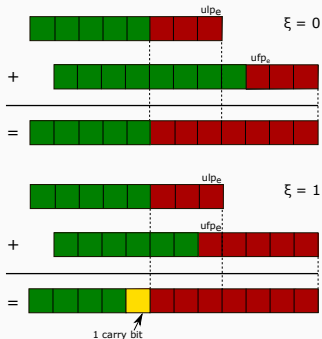
Optimized Carry Bit

Let $x^l = c_1^{l_1} \odot c_2^{l_2}$, $l \in Lab$ and $\odot \in \{+, -, \times, \div\}$

$$carry(l, l_1, l_2) = \begin{cases} 0 & \text{if no overlap errors} \\ 1 & \text{otherwise.} \end{cases}$$

✗ Before this work, very costly over-approximation in [NFM'17]

✓ Optimized carry bit function



Optimized Carry Bit

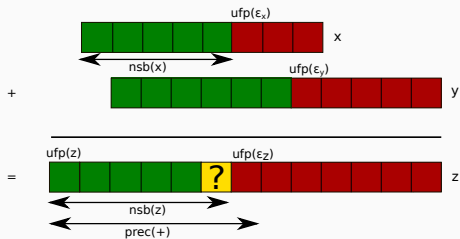
Let $x^l = c_1^{l_1} \odot c_2^{l_2}$, $l \in Lab$ and $\odot \in \{+, -, \times, \div\}$

$$carry(l, l_1, l_2) = \begin{cases} 0 & \text{if no overlap errors} \\ 1 & \text{otherwise.} \end{cases}$$

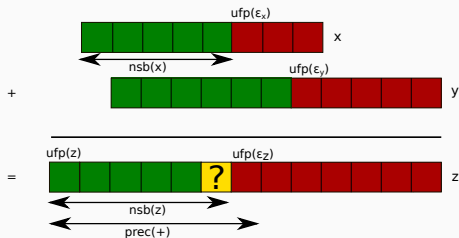
$$carry(l, l_1, l_2) = \min \left(\begin{array}{l} \max \left(\text{ufp}(l_2) - \text{ufp}(l_1) + \text{nsb}(l_1) - \text{nsb}(l_2) - \text{nsb}_e(l_2), 0 \right), \\ \max \left(\text{ufp}(l_1) - \text{ufp}(l_2) + \text{nsb}(l_2) - \text{nsb}(l_1) - \text{nsb}_e(l_1), 0 \right), 1 \end{array} \right)$$

Non-linear computations (min and max operators)

Preliminary Elements Transfer Functions



³Chapter 4 of the dissertation discusses the arithmetic expressions



Forward addition

$$\vec{\oplus}(x, y) = z$$

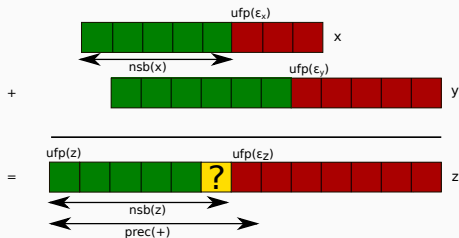
$$nsb(z) = ufp(x + y) - ufp(\epsilon(x) + \epsilon(y) + \epsilon_+)$$

Backward addition

$$\overleftarrow{\oplus}(z, y) = (z - y)$$

$$nsb(x) = ufp(z - y) - ufp(\epsilon(z) - \epsilon(y) - \epsilon_+)$$

³Chapter 4 of the dissertation discusses the arithmetic expressions



Forward addition

$$\vec{\oplus}(x, y) = z$$

$$\text{nsb}(z) = \text{ufp}(x + y) - \text{ufp}(\varepsilon(x) + \varepsilon(y) + \varepsilon_+)$$

Backward addition

$$\overleftarrow{\oplus}(z, y) = (z - y)$$

$$\text{nsb}(x) = \text{ufp}(z - y) - \text{ufp}(\varepsilon(z) - \varepsilon(y) - \varepsilon_+)$$

- 👉 Same for the multiplication, subtraction and division³
- 👉 Semantic of the commands in the dissertation
- 👉 Technique generalizable to sets of values

³Chapter 4 of the dissertation discusses the arithmetic expressions

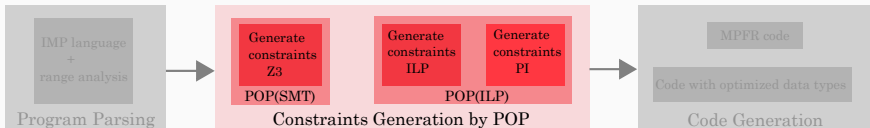
- 1 Preliminary Elements

- 2 **Constraint Generation**

- ⋮
- 3 Code Generation

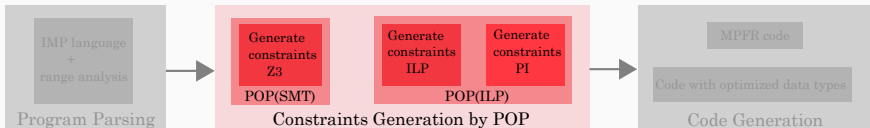
- ⋮
- 4 Conclusion and Perspectives

Constraints Generation by POP



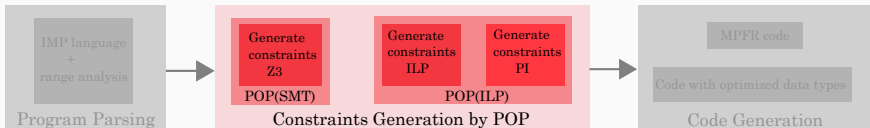
- **Approach:** static analysis of the error propagation

Constraints Generation by POP



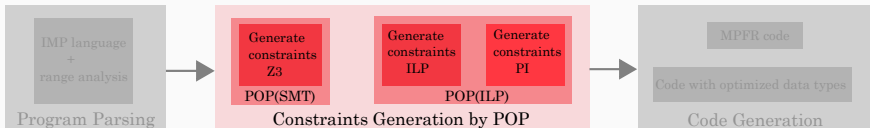
- **Approach:** static analysis of the error propagation
- **Formulation:** three systems of constraints

Constraints Generation by POP



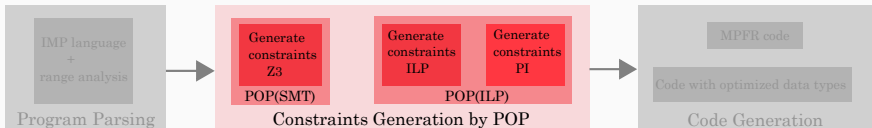
- **Approach:** static analysis of the error propagation
- **Formulation:** three systems of constraints
 1. **SMT-Based Method** checked by SMT solver
 - POP(SMT) version [FTSCS'19, IINTEC'19, IoTaIS'20, ICICT'21]

Constraints Generation by POP



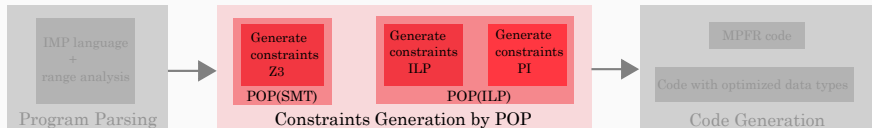
- **Approach:** static analysis of the error propagation
- **Formulation:** three systems of constraints
 1. **SMT-Based Method** checked by SMT solver
 - POP(SMT) version [FTSCS'19, IINTEC'19, IoTaIS'20, ICICT'21]
 2. **ILP-Based Method** solved by LP solver
 - Pessimistic carry bit
 - POP(ILP) version [ICCSA'21, SAS'21]

Constraints Generation by POP



- **Approach:** static analysis of the error propagation
- **Formulation:** three systems of constraints
 1. **SMT-Based Method** checked by SMT solver
 - POP(SMT) version [FTSCS'19, IINTEC'19, IoTaIS'20, ICICT'21]
 2. **ILP-Based Method** solved by LP solver
 - Pessimistic carry bit
 - POP(ILP) version [ICCSA'21, SAS'21]
 3. **PI-Based Method** solved with policy iteration algorithm
 - Optimized carry bit
 - POP(ILP) version [ICCSA'21, SAS'21]

Constraints Generation by POP



- **Approach:** static analysis of the error propagation
- **Formulation:** three systems of constraints
 1. **SMT-Based Method** checked by SMT solver
 - POP(SMT) version [FTSCS'19, IINTEC'19, IoTaIS'20, ICICT'21]
 2. **ILP-Based Method** solved by LP solver
 - Pessimistic carry bit
 - POP(ILP) version [ICCSA'21, SAS'21]
 3. **PI-Based Method** solved with policy iteration algorithm
 - Optimized carry bit
 - POP(ILP) version [ICCSA'21, SAS'21]
- **Output:** precision at bit-level translatable to IEEE754, fixed arithmetic, . . .

1 Preliminary Elements

2 **Constraint Generation**

• SMT-Based Method

•

•

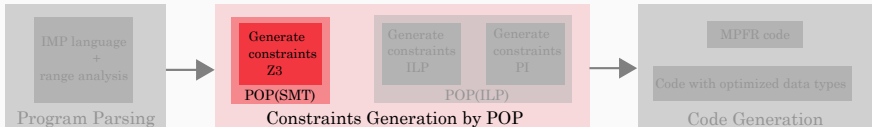
3 Code Generation

•

•

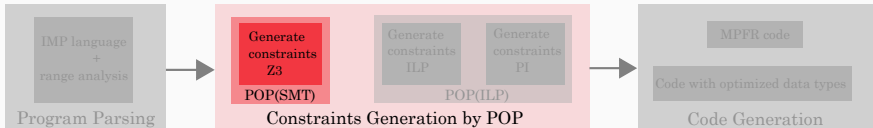
4 Conclusion and Perspectives

Formulation 1 SMT-Based Method



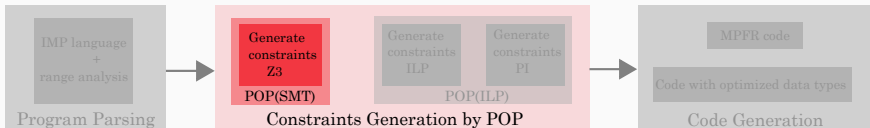
- First introduced in [NFM'17], improvement of carries in [FTSCS'19]

Formulation 1 SMT-Based Method



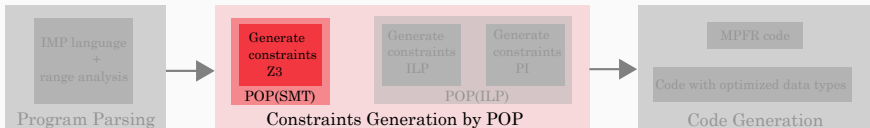
- First introduced in [NFM'17], improvement of carries in [FTSCS'19]
- **SMT-Based Method^{ab} = forward analysis + backward analysis**

Formulation 1 SMT-Based Method



- First introduced in [NFM'17], improvement of carries in [FTSCS'19]
- **SMT-Based Method^{ab} = forward analysis + backward analysis**
- Assign three parameters to each variable: nsb_F , nsb_B and nsb

Formulation 1 SMT-Based Method



- First introduced in [NFM'17], improvement of carries in [FTSCS'19]
- **SMT-Based Method^{ab} = forward analysis + backward analysis**
- Assign three parameters to each variable: nsb_F , nsb_B and nsb
- Constraints of **first order predicates** and with **linear integer relations** only

Formulation 1 SMT-Based Method



- First introduced in [NFM'17], improvement of carries in [FTSCS'19]
- **SMT-Based Method^{ab} = forward analysis + backward analysis**
- Assign three parameters to each variable: nsb_F , nsb_B and nsb
- Constraints of **first order predicates** and with **linear integer relations** only
- Easy to solve for an SMT solver (Z3)

^aM. Martel. "Floating-Point Format Inference in Mixed-Precision". NFM'17

^bD. Ben Khalifa, M. Martel and A. Adjé. "POP: A Tuning Assistant for Mixed-Precision Floating-Point Computations". FTSCS'19

- Recall that for the forward addition we have

$$\varepsilon(z) = \varepsilon(x) + \varepsilon(y) + \varepsilon_+$$

- Recall that for the forward addition we have

$$\varepsilon(z) = \varepsilon(x) + \varepsilon(y) + \varepsilon_+$$

- First over-approximation

$$\varepsilon(z) \leq 2^{\text{ufp}(x) - \text{nsb}_F(x) + 1} + 2^{\text{ufp}(y) - \text{nsb}_F(y) + 1} + 2^{\text{ufp}(z) - \text{prec}(+)}$$

- Recall that for the forward addition we have

$$\varepsilon(z) = \varepsilon(x) + \varepsilon(y) + \varepsilon_+$$

- First over-approximation

$$\varepsilon(z) \leq 2^{\text{ufp}(x) - \text{nsb}_F(x) + 1} + 2^{\text{ufp}(y) - \text{nsb}_F(y) + 1} + 2^{\text{ufp}(z) - \text{prec}(+)}$$

- Second over-approximation

$$\varepsilon(z) \leq 2^{\max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+) + \text{carry}(z, x, y))}$$

- Recall that for the forward addition we have

$$\varepsilon(z) = \varepsilon(x) + \varepsilon(y) + \varepsilon_+$$

- First over-approximation

$$\varepsilon(z) \leq 2^{\text{ufp}(x) - \text{nsb}_F(x) + 1} + 2^{\text{ufp}(y) - \text{nsb}_F(y) + 1} + 2^{\text{ufp}(z) - \text{prec}(+)}$$

- Second over-approximation

$$\varepsilon(z) \leq 2^{\max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+)) + \text{carry}(z, x, y)}$$

- Forward addition** $\rightarrow \oplus$

$$\text{nsb}_F(z) \geq \text{ufp}(x + y) - \max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+))$$

$$- \text{carry}(z, x, y)$$

- Recall that for the forward addition we have

$$\varepsilon(z) = \varepsilon(x) + \varepsilon(y) + \varepsilon_+$$

- First over-approximation

$$\varepsilon(z) \leq 2^{\text{ufp}(x) - \text{nsb}_F(x) + 1} + 2^{\text{ufp}(y) - \text{nsb}_F(y) + 1} + 2^{\text{ufp}(z) - \text{prec}(+)}$$

- Second over-approximation

$$\varepsilon(z) \leq 2^{\max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+) + \text{carry}(z, x, y))}$$

- Forward addition** $\rightarrow \oplus$

$$\text{nsb}_F(z) \geq \text{ufp}(x + y) - \max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+))$$

$$- \text{carry}(z, x, y)$$

- Backward addition** $\leftarrow \oplus$

$$\text{nsb}_B(x) \leq \text{ufp}(z - y) - \text{ufp}(z) + \text{nsb}(z)$$

- Recall that for the forward addition we have

$$\varepsilon(z) = \varepsilon(x) + \varepsilon(y) + \varepsilon_+$$

- First over-approximation

$$\varepsilon(z) \leq 2^{\text{ufp}(x) - \text{nsb}_F(x) + 1} + 2^{\text{ufp}(y) - \text{nsb}_F(y) + 1} + 2^{\text{ufp}(z) - \text{prec}(+)}$$

- Second over-approximation

$$\varepsilon(z) \leq 2^{\max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+) + \text{carry}(z, x, y))}$$

- Forward addition** $\rightarrow \oplus$

$$\text{nsb}_F(z) \geq \text{ufp}(x + y) - \max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+))$$

$$- \text{carry}(z, x, y)$$

- Backward addition** $\leftarrow \oplus$

$$\text{nsb}_B(x) \leq \text{ufp}(z - y) - \text{ufp}(z) + \text{nsb}(z)$$

- Final precision**

$$0 \leq \text{nsb}_B(x) \leq \text{nsb}(x) \leq \text{nsb}_F(x)$$

Elementary functions (sin, cos, tan, arcsin, log, . . .)

- Loss of precision of φ bits, where $\varphi \in \mathbb{N}$ a parameter of the analysis

example $x = 3.0$, $\varphi = 9$ and $require_nsb(\sin(x), 26) \implies \vec{\sin}(x) = \sin(3.0|35|)|26|$

⁴Discussed in Chapter 5 of the dissertation

Elementary functions (sin, cos, tan, arcsin, log, . . .)

- Loss of precision of φ bits, where $\varphi \in \mathbb{N}$ a parameter of the analysis

example $x = 3.0, \varphi = 9$ and $require_nsb(\sin(x), 26) \implies \vec{\sin}(x) = \sin(3.0|35|)|26|$

Loops

- Managed correctly thanks to range analysis
- Relate `nsb` at the end of the `body` to `nsb` of the same variables and the beginning of the loop

⁴Discussed in Chapter 5 of the dissertation

Elementary functions (sin, cos, tan, arcsin, log, . . .)

- Loss of precision of φ bits, where $\varphi \in \mathbb{N}$ a parameter of the analysis

example $x = 3.0, \varphi = 9$ and $require_nsb(\sin(x), 26) \implies \vec{\sin}(x) = \sin(3.0|35|)|26|$

Loops

- Managed correctly thanks to range analysis
- Relate `nsb` at the end of the `body` to `nsb` of the same variables and the beginning of the loop

Conditionals

- Analyze both `then` and `else` branches of the `if` statement without reducing the ranges of the variables
- ✗ Do not take care of the guards yet

⁴Discussed in Chapter 5 of the dissertation


```

1  gℓ1 = 9.81ℓ0; iℓ3 = 0.5ℓ2;
2  y1ℓ5 = 0.785398ℓ4;
3  y2ℓ7 = 0.785398ℓ6;
4  hℓ9 = 0.1ℓ8; tℓ11 = 0.0ℓ10;
5  while (tℓ13 <ℓ15 10.0ℓ14)ℓ59 {
6    y1newℓ24 = y1ℓ17 +ℓ23 y2ℓ19 *ℓ22 hℓ21;
7    aux1ℓ28 = sin(y1ℓ26)ℓ27;
8    aux2ℓ40 = aux1ℓ30 *ℓ39 hℓ32
9      *ℓ38 gℓ34 /ℓ37 iℓ36;
10   y2newℓ46 = y2ℓ42 -ℓ45 aux2ℓ44;
11   tℓ52 = tℓ48 +ℓ51 hℓ50;
12   y1ℓ55 = y1newℓ54;
13   y2ℓ58 = y2newℓ57; };
14   require_nsb(y2, 20)ℓ61;

```

POP Label File

⁵D. Ben Khalifa, M. Martel and A. Adjé. "POP: A Tuning Assistant for Mixed-Precision Floating-Point Computations". FTSCS'19

Running Example SMT Constraints⁵

```

1  gℓ1 = 9.81ℓ0; iℓ3 = 0.5ℓ2;
2  y1ℓ5 = 0.785398ℓ4;
3  y2ℓ7 = 0.785398ℓ6;
4  hℓ9 = 0.1ℓ8; tℓ11 = 0.0ℓ10;
5  while (tℓ13 <ℓ15 10.0ℓ14)ℓ59 {
6  y1newℓ24 = y1ℓ17 +ℓ23 y2ℓ19 *ℓ22 hℓ21;
7  aux1ℓ28 = sin(y1ℓ26)ℓ27;
8  aux2ℓ40 = aux1ℓ30 *ℓ39 hℓ32
9  *ℓ38 gℓ34 /ℓ37 iℓ36;
10 y2newℓ46 = y2ℓ42 -ℓ45 aux2ℓ44;
11 tℓ52 = tℓ48 +ℓ51 hℓ50;
12 y1ℓ55 = y1newℓ54;
13 y2ℓ58 = y2newℓ57; };
14 require_nsb(y2, 20)ℓ61;

```


POP Label File

$C_{Z3} =$

```

//ASSIGN
(assert(<= nsb(ℓ24) nsbF(ℓ24)))
(assert(<= nsbF(ℓ56) nsbF(ℓ24)))
(assert(<= nsbB(ℓ24) nsbB(ℓ19)))
(assert(>= nsbB(ℓ24) nsbB(ℓ13)))
(assert(= nsbB(ℓ56) nsbB(ℓ24)))
...
//ADDITION
(assert(<= nsb(ℓ17) nsbF(ℓ17)))
(assert(<= nsbF(ℓ17) nsbF(ℓ5)))
...
//MULTIPLICATION
(assert(or(and(<= nsb(ℓ17) nsbF(ℓ17))(>= nsb(ℓ22)
nsbB(ℓ22)))(and(<= nsb(ℓ22)
nsbF(ℓ22))(>= nsb(ℓ17) nsbB(ℓ17)))))
...
//CARRY BIT
(assert(= carry(ℓ23, ℓ17, ℓ22)(ite(> (ite(> ulpe(ℓ0)(-0 56)) 0 1)
(ite(> ulpe(ℓ0) (- 0 52)) 0 1))(ite(> ulpe(ℓ0)(- 0 56)) 0 1)
(ite(> ulpe(ℓ0)(- 0 52)) 0 1)))
...
...
...

```

 Linear number of constraints / variables in the size of the analyzed program

⁵D. Ben Khalifa, M. Martel and A. Adjé. "POP: A Tuning Assistant for Mixed-Precision Floating-Point Computations". FTSCS'19

```
1   g|20| = 9.81|20|; l|20| = 1.5|20|;
2   y1|29| = 0.785398|29|;
3   y2|21| = 0.0|21|;
4   h|21| = 0.1|21|; t|21| = 0.0|21|;
5   while (t<1.0) {
6     y1new|20| = y1|21| +|20| y2|21|*|22| h|21|;
7     aux1|20| = sin(y1|29|)|20|;
8     aux2|20| = aux1|19| *|20| h|18|*|19|g|17| /|18|l|17|;
9     y2new|20| = y2|21| -|20| aux2|18|;
10    t|20| = t|21| +|20| h|17|;
11    y1|20| = y1new|20|;
12    y2|20| = y2new|20|;
13  };
14  require_nsb(y2, 20);
```

File POP_output_Z3

Running Example Results with POP(SMT)

```
1  g|20| = 9.81|20|; l|20| = 1.5|20|;
2  y1|29| = 0.785398|29|;
3  y2|21| = 0.0|21|;
4  h|21| = 0.1|21|; t|21| = 0.0|21|;
5  while (t<1.0) {
6    y1new|20| = y1|21| +|20| y2|21|*|22| h|21|;
7    aux1|20| = sin(y1|29|)|20|;
8    aux2|20| = aux1|19| *|20| h|18|*|19|g|17| /|18|l|17|;
9    y2new|20| = y2|21| -|20| aux2|18|;
10   t|20| = t|21| +|20| h|17|;
11   y1|20| = y1new|20|;
12   y2|20| = y2new|20|;
13 };
14 require_nsb(y2, 20);
```

File POP_output_Z3

Tool	#Bits	#Var Solver	#Cstr Solver	Time (s)
POP(SMT)	960	314	431	13.1

Running Example Results with POP(SMT)

```
1      g|20| = 9.81|20|; l|20| = 1.5|20|;
2      y1|29| = 0.785398|29|;
3      y2|21| = 0.0|21|;
4      h|21| = 0.1|21|; t|21| = 0.0|21|;
5      while (t<1.0) {
6          y1new|20| = y1|21| +|20| y2|21|*|22| h|21|;
7          aux1|20| = sin(y1|29|)|20|;
8          aux2|20| = aux1|19| *|20| h|18|*|19|g|17| /|18||17|;
9          y2new|20| = y2|21| -|20| aux2|18|;
10         t|20| = t|21| +|20| h|17|;
11         y1|20| = y1new|20|;
12         y2|20| = y2new|20|;
13     };
14     require_nsb(y2, 20);
```

File POP_output_Z3

Tool	#Bits	#Var Solver	#Cstr Solver	Time (s)
POP(SMT)	960	314	431	13.1

Cons

- ✗ Too complex system of constraints
- ✗ Non-optimizing solver coupled with binary search
- ✗ Is the forward analysis really useful?

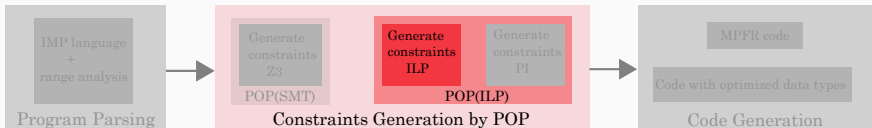
- 1 Preliminary Elements

- 2 **Constraint Generation**

 - ILP-Based Method
- 3 Code Generation

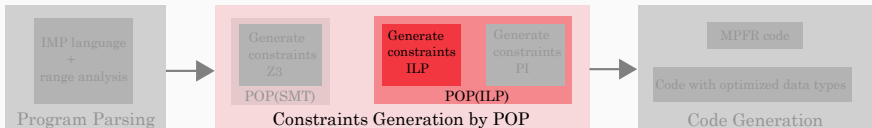
 -
- 4 Conclusion and Perspectives

Formulation 2 ILP-Based Method



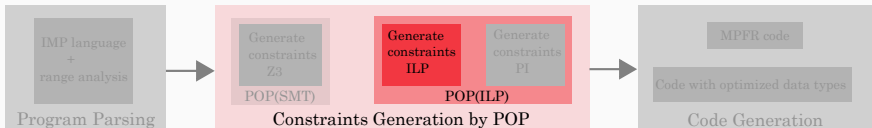
- Relaxation of the SMT-based method \implies backward analysis only

Formulation 2 ILP-Based Method



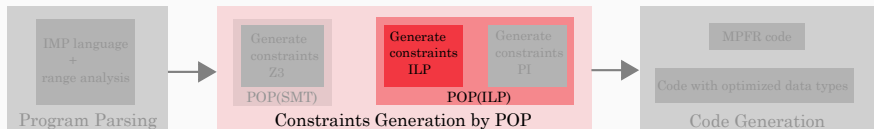
- Relaxation of the SMT-based method \implies backward analysis only
- Generate an ILP from the program source code

Formulation 2 ILP-Based Method



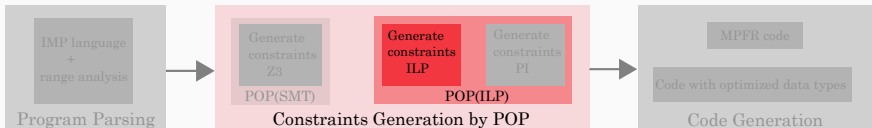
- Relaxation of the SMT-based method \implies backward analysis only
- Generate an ILP from the program source code
- Over-approximated carry bit function \implies $carry = 1$

Formulation 2 ILP-Based Method



- Relaxation of the SMT-based method \implies backward analysis only
- Generate an ILP from the program source code
- Over-approximated carry bit function $\implies carry = 1$
- Former constraints used to generate an ILP

Formulation 2 ILP-Based Method



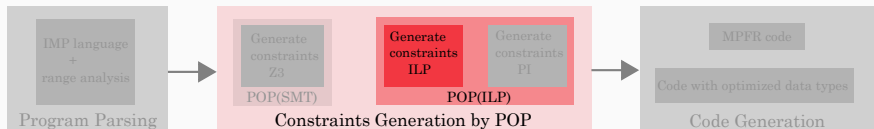
- Relaxation of the SMT-based method \implies backward analysis only
- Generate an ILP from the program source code
- Over-approximated carry bit function $\implies carry = 1$
- Former constraints used to generate an ILP

Case of Addition

☞ Since

$$nsb(z) \leq ufp(z) - \max(ufp(x) - nsb(x), ufp(y) - nsb(y)) - carry(z, y, x)$$

Formulation 2 ILP-Based Method



- Relaxation of the SMT-based method \implies backward analysis only
- Generate an ILP from the program source code
- Over-approximated carry bit function $\implies carry = 1$
- Former constraints used to generate an ILP

Case of Addition

☞ Since

$$nsb(z) \leq ufp(z) - \max(ufp(x) - nsb(x), ufp(y) - nsb(y)) - carry(z, y, x)$$

☞ We generate an ILP (with relaxation)

$$(S) = \begin{cases} nsb(x) \geq nsb(z) + ufp(x) - ufp(z) + 1 \\ nsb(y) \geq nsb(z) + ufp(y) - ufp(z) + 1 \end{cases}$$

```

1 gℓ1 = 9.81ℓ0; |ℓ3 = 0.5ℓ2;
2 y1ℓ5 = 0.785398ℓ4;
3 y2ℓ7 = 0.785398ℓ6;
4 hℓ9 = 0.1ℓ8; tℓ11 = 0.0ℓ10;
5 while (tℓ13 <ℓ15 10.0ℓ14)ℓ59 {
6   y1newℓ24 = y1ℓ17 +ℓ23 y2ℓ19 *ℓ22 hℓ21;
7   aux1ℓ28 = sin(y1ℓ26)ℓ27;
8   aux2ℓ40 = aux1ℓ30 *ℓ39 hℓ32
9     *ℓ38 gℓ34 /ℓ37 |ℓ36;
10  y2newℓ46 = y2ℓ42 -ℓ45 aux2ℓ44;
11  tℓ52 = tℓ48 +ℓ51 hℓ50;
12  y1ℓ55 = y1newℓ54;
13  y2ℓ58 = y2newℓ57; };
14 require_nsb(y2, 20)ℓ61;

```

POP Label File

⁶ILP nature of the problem presented in Chapter 5 Theorem 5.2

⁷A. Adjé, D. Ben Khalifa and M. Martel. "Fast and Efficient Bit-Level Precision Tuning". SAS'21

```

1 gℓ1 = 9.81ℓ0; lℓ3 = 0.5ℓ2;
2 y1ℓ5 = 0.785398ℓ4;
3 y2ℓ7 = 0.785398ℓ6;
4 hℓ9 = 0.1ℓ8; tℓ11 = 0.0ℓ10;
5 while (tℓ13 <ℓ15 10.0ℓ14)ℓ59 {
6   y1newℓ24 = y1ℓ17 +ℓ23 y2ℓ19 *ℓ22 hℓ21;
7   aux1ℓ28 = sin(y1ℓ26)ℓ27;
8   aux2ℓ40 = aux1ℓ30 *ℓ39 hℓ32
9     *ℓ38 gℓ34 /ℓ37 lℓ36;
10  y2newℓ46 = y2ℓ42 -ℓ45 aux2ℓ44;
11  tℓ52 = tℓ48 +ℓ51 hℓ50;
12  y1ℓ55 = y1newℓ54;
13  y2ℓ58 = y2newℓ57; }
14 require_nsb(y2, 20)ℓ61;

```

$$C_{GLPK} = \left\{ \begin{array}{l} \text{nsb}(\ell_{17}) \geq \text{nsb}(\ell_{23}) + (-1) + \text{carry}(\ell_{23})(\ell_{17}, \ell_{22}) - (-1) // \text{ADD} \\ \text{nsb}(\ell_{22}) \geq \text{nsb}(\ell_{23}) + 0 + \text{carry}(\ell_{23})(\ell_{17}, \ell_{22}) - (1) // \text{ADD} \\ \text{nsb}(\ell_{19}) \geq \text{nsb}(\ell_{22}) + \text{carry}(\ell_{22})(\ell_{19}, \ell_{21}) - 1 // \text{MULT} \\ \text{nsb}(\ell_{21}) \geq \text{nsb}(\ell_{22}) + \text{carry}(\ell_{22})(\ell_{19}, \ell_{21}) - 1 // \text{MULT} \\ \text{nsb}(\ell_{23}) \geq \text{nsb}(\ell_{24}) // \text{ASSIGN} \\ \text{carry}(\ell_{23}, \ell_{17}, \ell_{22}) = 1, \text{carry}(\ell_{22}, \ell_{19}, \ell_{21}) = 1 // \text{CARRY BIT} \end{array} \right.$$

POP Label File

⁶ILP nature of the problem presented in Chapter 5 Theorem 5.2

⁷A. Adjé, D. Ben Khalifa and M. Martel. "Fast and Efficient Bit-Level Precision Tuning". SAS'21

Running Example ILP Constraints⁷

```
1 gℓ1 = 9.81ℓ0; lℓ3 = 0.5ℓ2;
2 y1ℓ5 = 0.785398ℓ4;
3 y2ℓ7 = 0.785398ℓ6;
4 hℓ9 = 0.1ℓ8; tℓ11 = 0.0ℓ10;
5 while (tℓ13 <ℓ15 10.0ℓ14)ℓ59 {
6   y1newℓ24 = y1ℓ17 +ℓ23 y2ℓ19 *ℓ22 hℓ21;
7   aux1ℓ28 = sin(y1ℓ26)ℓ27;
8   aux2ℓ40 = aux1ℓ30 *ℓ39 hℓ32
9     *ℓ38 gℓ34 /ℓ37 lℓ36;
10  y2newℓ46 = y2ℓ42 -ℓ45 aux2ℓ44;
11  tℓ52 = tℓ48 +ℓ51 hℓ50;
12  y1ℓ55 = y1newℓ54;
13  y2ℓ58 = y2newℓ57; };
14 require_nsb(y2, 20)ℓ61;
```

POP Label File

$$C_{GLPK} = \left\{ \begin{array}{l} \text{nsb}(\ell_{17}) \geq \text{nsb}(\ell_{23}) + (-1) + \text{carry}(\ell_{23})(\ell_{17}, \ell_{22}) - (-1) // \text{ADD} \\ \text{nsb}(\ell_{22}) \geq \text{nsb}(\ell_{23}) + 0 + \text{carry}(\ell_{23})(\ell_{17}, \ell_{22}) - (1) // \text{ADD} \\ \text{nsb}(\ell_{19}) \geq \text{nsb}(\ell_{22}) + \text{carry}(\ell_{22})(\ell_{19}, \ell_{21}) - 1 // \text{MULT} \\ \text{nsb}(\ell_{21}) \geq \text{nsb}(\ell_{22}) + \text{carry}(\ell_{22})(\ell_{19}, \ell_{21}) - 1 // \text{MULT} \\ \text{nsb}(\ell_{23}) \geq \text{nsb}(\ell_{24}) // \text{ASSIGN} \\ \text{carry}(\ell_{23}, \ell_{17}, \ell_{22}) = 1, \text{carry}(\ell_{22}, \ell_{19}, \ell_{21}) = 1 // \text{CARRY BIT} \end{array} \right.$$

POP(ILP)⁶

- Linear number of constraints / variables in the size of the analyzed program
- Integer solution computed in **polynomial time**
- LP solver among reals

⁶ILP nature of the problem presented in Chapter 5 Theorem 5.2

⁷A. Adjé, D. Ben Khalifa and M. Martel. "Fast and Efficient Bit-Level Precision Tuning". SAS'21

Running Example Results with POP(ILP)

```
1 g|20| = 9.81|20|; l|20| = 1.5|20|;
2 y1|29| = 0.785398|29|;
3 y2|22| = 0.0|22|;
4 h|22| = 0.1|22|; t|21| = 0.0|21|;
5 while (t < 1.0) {
6   y1new|20| = y1|21| +|20| y2|22|*|22| h|22|;
7   aux1|20| = sin(y1|29|)|20|;
8   aux2|20| = aux1|20| *|20| h|20|*|20| g|20|/|20|l|20|;
9   y2new|20| = y2|21| -|20| aux2|18|;
10  t|20| = t|21| +|20| h|17|;
11  y1|20| = y1new|20|;
12  y2|20| = y2new|20|;
13 };
14 require_nsb(y2, 20);
```

File POP_output_GLPK

Tool	#Bits	#Var Solver	#Cstr Solver	Time (s)
POP(SMT)	960	314	431	13.1
POP(ILP)	861	52	69	3.5

Running Example Results with POP(ILP)

```
1 g|20| = 9.81|20|; l|20| = 1.5|20|;
2 y1|29| = 0.785398|29|;
3 y2|22| = 0.0|22|;
4 h|22| = 0.1|22|; t|21| = 0.0|21|;
5 while (t < 1.0) {
6   y1new|20| = y1|21| +|20| y2|22|*|22| h|22|;
7   aux1|20| = sin(y1|29|)|20|;
8   aux2|20| = aux1|20| *|20| h|20|*|20| g|20|/|20|l|20|;
9   y2new|20| = y2|21| -|20| aux2|18|;
10  t|20| = t|21| +|20| h|17|;
11  y1|20| = y1new|20|;
12  y2|20| = y2new|20|;
13 };
14 require_nsb(y2, 20);
```

File POP_output_GLPK

Tool	#Bits	#Var Solver	#Cstr Solver	Time (s)
POP(SMT)	960	314	431	13.1
POP(ILP)	861	52	69	3.5

👉 How to avoid the pessimistic carry bit?

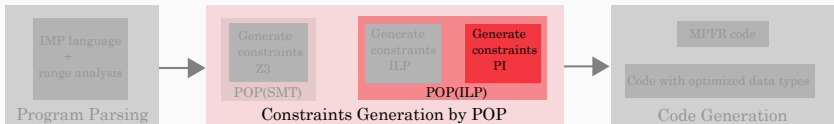
- 1 Preliminary Elements

- 2 **Constraint Generation**

 - PI-Based Method
- 3 Code Generation

 -
- 4 Conclusion and Perspectives

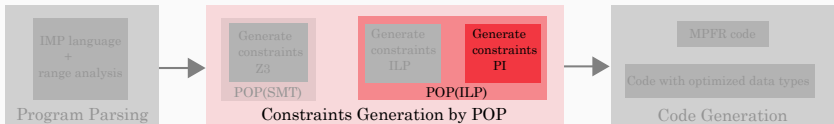
Formulation 3 PI-Based Method



- Policy iteration algorithm⁸ to solve min – max equations

⁸Algorithm 2 presented in Chapter 5 page 86

Formulation 3 PI-Based Method

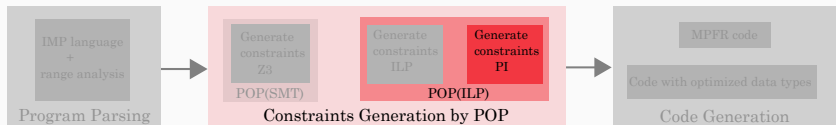


- Policy iteration algorithm⁸ to solve min – max equations
- Back to Line 6 of the pendulum example

$$\text{carry}(\ell, \ell_1, \ell_2) = \min \left(\begin{array}{l} \max(\text{ufp}(\ell_2) - \text{ufp}(\ell_1) + \text{nsb}(\ell_1) - \text{nsb}(\ell_2) - \text{nsb}_e(\ell_2), 0), \\ \max(\text{ufp}(\ell_1) - \text{ufp}(\ell_2) + \text{nsb}(\ell_2) - \text{nsb}(\ell_1) - \text{nsb}_e(\ell_1), 0), 1 \end{array} \right)$$

⁸Algorithm 2 presented in Chapter 5 page 86

Formulation 3 PI-Based Method



- Policy iteration algorithm⁸ to solve min – max equations
- Back to Line 6 of the pendulum example

$$\text{carry}(\ell, \ell_1, \ell_2) = \min \left(\begin{array}{l} \max(\text{ufp}(\ell_2) - \text{ufp}(\ell_1) + \text{nsb}(\ell_1) - \text{nsb}(\ell_2) - \text{nsb}_e(\ell_2), 0), \\ \max(\text{ufp}(\ell_1) - \text{ufp}(\ell_2) + \text{nsb}(\ell_2) - \text{nsb}(\ell_1) - \text{nsb}_e(\ell_1), 0), 1 \end{array} \right)$$

Principle

- Choose a policy $\pi_0 \in \Pi$ by breaking the min operator
- Associate policy map f^{π_0} to π_0
- example** $f^{\pi_0} = \max(\text{ufp}(\ell_2) - \text{ufp}(\ell_1) + \text{nsb}(\ell_1) - \text{nsb}(\ell_2) - \text{nsb}_e(\ell_2), 0)$
- Solve corresponding ILP problem
- If a more precise solution is found, PI algorithm iterates

⁸Algorithm 2 presented in Chapter 5 page 86

```
1  gℓ1 = 9.81ℓ0; lℓ3 = 0.5ℓ2;
2  y1ℓ5 = 0.785398ℓ4;
3  y2ℓ7 = 0.785398ℓ6;
4  hℓ9 = 0.1ℓ8; tℓ11 = 0.0ℓ10;
5  while (tℓ13 <ℓ15 10.0ℓ14)ℓ59 {
6    y1newℓ24 = y1ℓ17 +ℓ23 y2ℓ19 *ℓ22 hℓ21;
7    aux1ℓ28 = sin(y1ℓ26)ℓ27;
8    aux2ℓ40 = aux1ℓ30 *ℓ39 hℓ32
9      *ℓ38 gℓ34 /ℓ37 lℓ36;
10   y2newℓ46 = y2ℓ42 -ℓ45 aux2ℓ44;
11   tℓ52 = tℓ48 +ℓ51 hℓ50;
12   y1ℓ55 = y1newℓ54;
13   y2ℓ58 = y2newℓ57; };
14   require_nsb(y2, 20)ℓ61;
```

POP Label File

⁹A. Adjé, D. Ben Khalifa and M. Martel. "Fast and Efficient Bit-Level Precision Tuning". SAS'21

Running Example PI Constraints⁹

```

1  gℓ1 = 9.81ℓ0; lℓ3 = 0.5ℓ2;
2  y1ℓ5 = 0.785398ℓ4;
3  y2ℓ7 = 0.785398ℓ6;
4  hℓ9 = 0.1ℓ8; tℓ11 = 0.0ℓ10;
5  while (tℓ13 <ℓ15 10.0ℓ14)ℓ59 {
6    y1newℓ24 = y1ℓ17 +ℓ23 y2ℓ19 *ℓ22 hℓ21;
7    aux1ℓ28 = sin(y1ℓ26)ℓ27;
8    aux2ℓ40 = aux1ℓ30 *ℓ39 hℓ32
9      *ℓ38 gℓ34 /ℓ37 lℓ36;
10   y2newℓ46 = y2ℓ42 -ℓ45 aux2ℓ44;
11   tℓ52 = tℓ48 +ℓ51 hℓ50;
12   y1ℓ55 = y1newℓ54;
13   y2ℓ58 = y2newℓ57; };
14   require_nsb(y2, 20)ℓ61;

```

POP Label File

$$C_{PI} = \left\{ \begin{array}{l} \text{nsb}_e(\ell_{23}) \geq \text{nsb}_e(\ell_{17}), //\text{ADD} \\ \text{nsb}_e(\ell_{23}) \geq \text{nsb}_e(\ell_{22}), //\text{ADD} \\ \text{nsb}_e(\ell_{23}) \geq -1 - 0 + \text{nsb}(\ell_{22}) - \text{nsb}(\ell_{17}) + \text{nsb}_e(\ell_{22}) + \\ \text{carry}(\ell_{23}, \ell_{17}, \ell_{22}), //\text{ADD} \\ \text{nsb}_e(\ell_{23}) \geq 0 - (-1) + \text{nsb}(\ell_{17}) - \text{nsb}(\ell_{22}) + \text{nsb}_e(\ell_{17}) + \\ \text{carry}(\ell_{23}, \ell_{17}, \ell_{22}), //\text{ADD} \\ \text{nsb}_e(\ell_{23}) \geq \text{nsb}_e(\ell_{24}), //\text{ADD} \\ \text{nsb}_e(\ell_{22}) \geq \text{nsb}(\ell_{19}) + \text{nsb}_e(\ell_{19}) + \text{nsb}_e(\ell_{21}) - 2, //\text{MULT} \\ \text{nsb}_e(\ell_{22}) \geq \text{nsb}(\ell_{21}) + \text{nsb}_e(\ell_{21}) + \text{nsb}_e(\ell_{19}) - 2, //\text{MULT} \\ \text{carry}(\ell_{23}, \ell_{17}, \ell_{22}) = \min \left(\begin{array}{l} \max(0 - 6 + \text{nsb}(\ell_{17}) - \text{nsb}(\ell_{22}) - \\ \text{nsb}_e(\ell_{17}), 0), \\ \max(6 - 0 + \text{nsb}(\ell_{22}) - \text{nsb}(\ell_{17}) - \\ \text{nsb}_e(\ell_{22}), 0), 1 \end{array} \right) \\ //\text{CARRY} \end{array} \right.$$

⁹A. Adjé, D. Ben Khalifa and M. Martel. "Fast and Efficient Bit-Level Precision Tuning". SAS'21

```

1  gℓ1 = 9.81ℓ0; lℓ3 = 0.5ℓ2;
2  y1ℓ5 = 0.785398ℓ4;
3  y2ℓ7 = 0.785398ℓ6;
4  hℓ9 = 0.1ℓ8; tℓ11 = 0.0ℓ10;
5  while (tℓ13 <ℓ15 10.0ℓ14)ℓ59 {
6    y1newℓ24 = y1ℓ17 +ℓ23 y2ℓ19 *ℓ22 hℓ21;
7    aux1ℓ28 = sin(y1ℓ26)ℓ27;
8    aux2ℓ40 = aux1ℓ30 *ℓ39 hℓ32
9      *ℓ38 gℓ34 /ℓ37 lℓ36;
10   y2newℓ46 = y2ℓ42 -ℓ45 aux2ℓ44;
11   tℓ52 = tℓ48 +ℓ51 hℓ50;
12   y1ℓ55 = y1newℓ54;
13   y2ℓ58 = y2newℓ57; };
14   require_nsb(y2, 20)ℓ61;

```

POP Label File

$$C_{PI} = \left\{ \begin{array}{l} \text{nsb}_e(\ell_{23}) \geq \text{nsb}_e(\ell_{17}), // \text{ADD} \\ \text{nsb}_e(\ell_{23}) \geq \text{nsb}_e(\ell_{22}), // \text{ADD} \\ \text{nsb}_e(\ell_{23}) \geq -1 - 0 + \text{nsb}(\ell_{22}) - \text{nsb}(\ell_{17}) + \text{nsb}_e(\ell_{22}) + \\ \text{carry}(\ell_{23}, \ell_{17}, \ell_{22}), // \text{ADD} \\ \text{nsb}_e(\ell_{23}) \geq 0 - (-1) + \text{nsb}(\ell_{17}) - \text{nsb}(\ell_{22}) + \text{nsb}_e(\ell_{17}) + \\ \text{carry}(\ell_{23}, \ell_{17}, \ell_{22}), // \text{ADD} \\ \text{nsb}_e(\ell_{23}) \geq \text{nsb}_e(\ell_{24}), // \text{ADD} \\ \text{nsb}_e(\ell_{22}) \geq \text{nsb}(\ell_{19}) + \text{nsb}_e(\ell_{19}) + \text{nsb}_e(\ell_{21}) - 2, // \text{MULT} \\ \text{nsb}_e(\ell_{22}) \geq \text{nsb}(\ell_{21}) + \text{nsb}_e(\ell_{21}) + \text{nsb}_e(\ell_{19}) - 2, // \text{MULT} \\ \\ \text{carry}(\ell_{23}, \ell_{17}, \ell_{22}) = \min \left(\begin{array}{l} \max(0 - 6 + \text{nsb}(\ell_{17}) - \text{nsb}(\ell_{22}) - \\ \text{nsb}_e(\ell_{17}), 0), \\ \max(6 - 0 + \text{nsb}(\ell_{22}) - \text{nsb}(\ell_{17}) - \\ \text{nsb}_e(\ell_{22}), 0), 1 \end{array} \right) \\ \\ // \text{CARRY} \end{array} \right.$$

POP(ILP) with PI

- $C = C_{GLPK} \cup C_{PI}$: final set of constraints
- Activate optimized *carry*
- Feasible solution fastly computed

⁹A. Adjé, D. Ben Khalifa and M. Martel. "Fast and Efficient Bit-Level Precision Tuning". SAS'21


```

1 g|20| = 9.81|20|; l|20| = 1.5|20|;
2 y1|29| = 0.785398|29|;
3 y2|21| = 0.0|21|;
4 h|21| = 0.1|21|; t|21| = 0.0|21|;
5 while (t<1.0) {
6   y1new|20| = y1|21| +|20| y2|21|*|22| h|21|;
7   aux1|20| = sin(y1|29|)|20|;
8   aux2|20| = aux1|19| *|20| h|18|*|19|g|17| /|18|l|17|;
9   y2new|20| = y2|21| -|20| aux2|18|;
10  t|20| = t|21| +|20| h|17|;
11  y1|20| = y1new|20|;
12  y2|20| = y2new|20|;
13 };
14 require_nsb(y2, 20);

```

File POP_output_PI

Tool	#Bits	#Var Solver	#Cstr Solver	Time (s)
POP(SMT)	960	314	431	13.1
POP(ILP)	861	52	69	3.5
POP(ILP)/ PI	843	97	204 (1 policy)	4.1

```

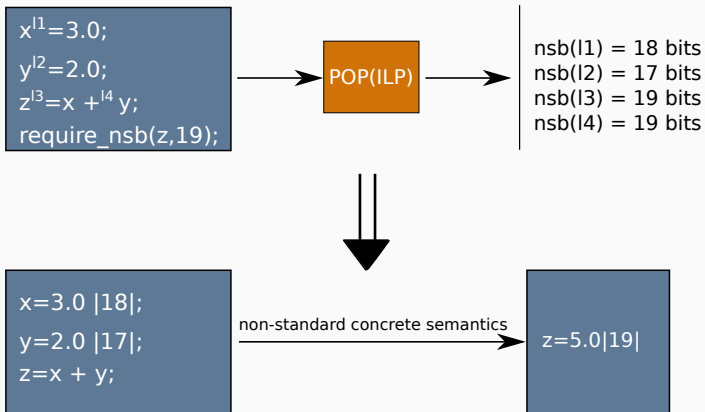
1 g|20| = 9.81|20|; l|20| = 1.5|20|;
2 y1|29| = 0.785398|29|;
3 y2|21| = 0.0|21|;
4 h|21| = 0.1|21|; t|21| = 0.0|21|;
5 while (t<1.0) {
6   y1new|20| = y1|21| +|20| y2|21|*|22| h|21|;
7   aux1|20| = sin(y1|29|)|20|;
8   aux2|20| = aux1|19| *|20| h|18|*|19|g|17| /|18|l|17|;
9   y2new|20| = y2|21| -|20| aux2|18|;
10  t|20| = t|21| +|20| h|17|;
11  y1|20| = y1new|20|;
12  y2|20| = y2new|20|;
13 };
14 require_nsb(y2, 20);

```

File POP_output_PI

Tool	#Bits	#Var Solver	#Cstr Solver	Time (s)
POP(SMT)	960	314	431	13.1
POP(ILP)	861	52	69	3.5
POP(ILP)/ PI	843	97	204 (1 policy)	4.1

- Soundness of the constraint system of the ILP formulation¹⁰



¹⁰Presented and proved in Chapter 5 Theorem 5.1

- 1 Preliminary Elements

- 2 Constraint Generation

- ⋮
- 3 Code Generation**

- ⋮
- 4 Conclusion and Perspectives



Code with mixed-precision at bit-level

- Translation of bit-level precision to the upper number of bits corresponding to an IEEE-754 format
- Fixed-point formats \implies soon!

¹¹<https://www.mpfr.org/>



Code with mixed-precision at bit-level

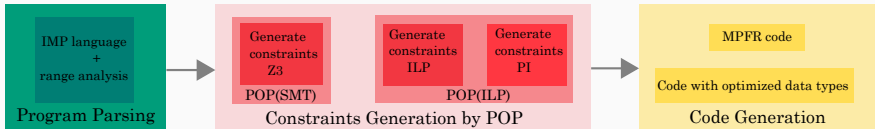
- Translation of bit-level precision to the upper number of bits corresponding to an IEEE-754 format
- Fixed-point formats \implies soon!

MPFR¹¹ code generation

- Generate MPFR program by assuming that the exact computations are performed in higher precision (**300 bits**)
- Generate MPFR program with precision returned by POP
- **Goal:** measure the difference between the two programs in fonction of the theoretical error given by the user

¹¹<https://www.mpfr.org/>

Evaluation Plan



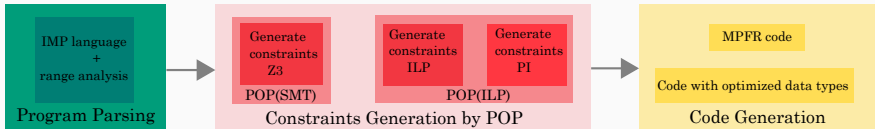
Evaluation¹² is meant to answer questions

¹²Intel Core i5-8350U CPU cadenced at 1.7GHz on a Linux machine with 8 GB RAM

¹³<https://www.gnu.org/software/gsl/>

¹⁴<https://fpbench.org/>

Evaluation Plan



Evaluation¹² is meant to answer questions

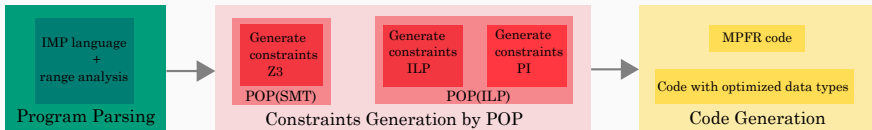
- Which version of POP is more efficient?

¹²Intel Core i5-8350U CPU cadenced at 1.7GHz on a Linux machine with 8 GB RAM

¹³<https://www.gnu.org/software/gsl/>

¹⁴<https://fpbench.org/>

Evaluation Plan



Evaluation¹² is meant to answer questions

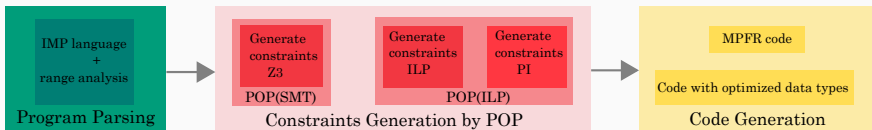
- Which version of POP is more efficient?
- How the returned precision are validated?

¹²Intel Core i5-8350U CPU cadenced at 1.7GHz on a Linux machine with 8 GB RAM

¹³<https://www.gnu.org/software/gsl/>

¹⁴<https://fpbench.org/>

Evaluation Plan



Evaluation¹² is meant to answer questions

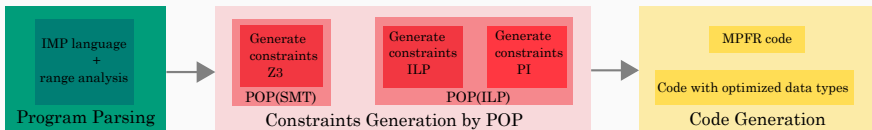
- Which version of POP is more efficient?
- How the returned precision are validated?
- Which impact has the policy iteration (PI) method on POP(ILP)?

¹²Intel Core i5-8350U CPU cadenced at 1.7GHz on a Linux machine with 8 GB RAM

¹³<https://www.gnu.org/software/gsl/>

¹⁴<https://fpbench.org/>

Evaluation Plan



Evaluation¹² is meant to answer questions

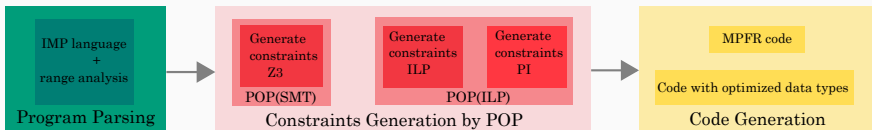
- Which version of POP is more efficient?
- How the returned precision are validated?
- Which impact has the policy iteration (PI) method on POP(ILP)?
- POP(SMT) vs POP(ILP) vs Precimonious?

¹²Intel Core i5-8350U CPU cadenced at 1.7GHz on a Linux machine with 8 GB RAM

¹³<https://www.gnu.org/software/gsl/>

¹⁴<https://fpbench.org/>

Evaluation Plan



Evaluation¹² is meant to answer questions

- Which version of POP is more efficient?
- How the returned precision are validated?
- Which impact has the policy iteration (PI) method on POP(ILP)?
- POP(SMT) vs POP(ILP) vs Precimonious?

Target applications

- IoT, GNU scientific library¹³ (arclength, simpson, . . .), physics (N-Body program \approx 400 LOCs), FPBench¹⁴ (Trapeze, Runge Kutta, PID, . . .), . . .

¹²Intel Core i5-8350U CPU cadenced at 1.7GHz on a Linux machine with 8 GB RAM

¹³<https://www.gnu.org/software/gsl/>

¹⁴<https://fpbench.org/>

Optimization parameters

- 100 \equiv all control points (CP) in **FP64** precision (53 mantissa size)

example If we have 8 |CP| $\implies 100 = 8 \times 53 = 424$ bits

Optimization parameters

- 100 \equiv all control points (CP) in **FP64** precision (53 mantissa size)

example If we have 8 |CP| \implies 100 = 8 \times 53 = 424 bits

Optimization in Bit-Level (BL)

$$BL = \frac{\sum_{I \in CP} nsb(\ell) \times 100}{|CP| \times 53}$$

Optimization parameters

- 100 \equiv all control points (CP) in **FP64** precision (53 mantissa size)

example If we have 8 |CP| $\implies 100 = 8 \times 53 = 424$ bits

Optimization in Bit-Level (BL)

$$BL = \frac{\sum_{l \in CP} nsb(l) \times 100}{|CP| \times 53}$$

From precision at bit-level to IEEE mode

- Round to the upper IEEE-754 format

example $nsb = 8$ bits \longrightarrow 10 (FP16)

example $nsb = 20$ bits \longrightarrow 23 (FP32)

example $nsb = 25$ bits \longrightarrow 53 (FP64)

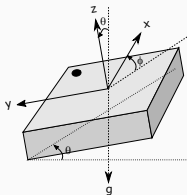
- 1 Preliminary Elements

- 2 Constraint Generation

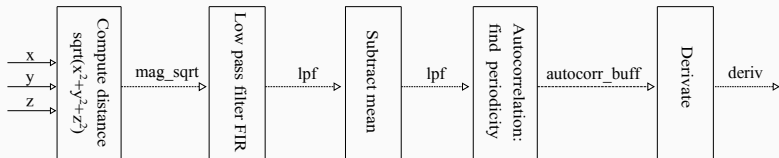
- ⋮
- 3 Code Generation**

 - Evaluation POP(SMT)
 -
- 4 Conclusion and Perspectives

- Application 1: Tilt angle measurements by a 3 axis accelerometer¹⁵



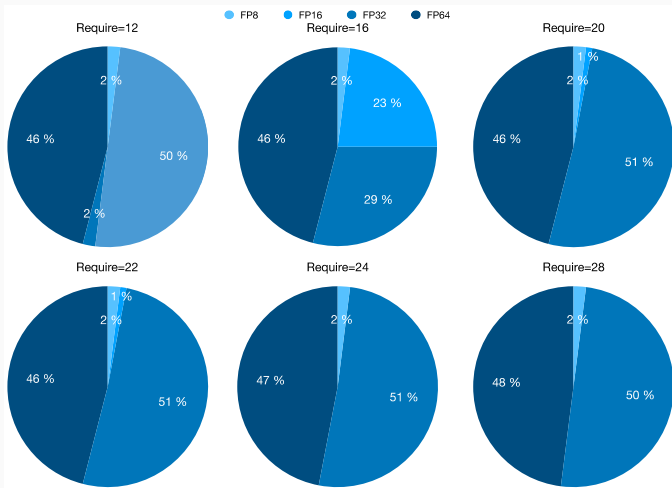
- Application 2: Footstep counter program¹⁶



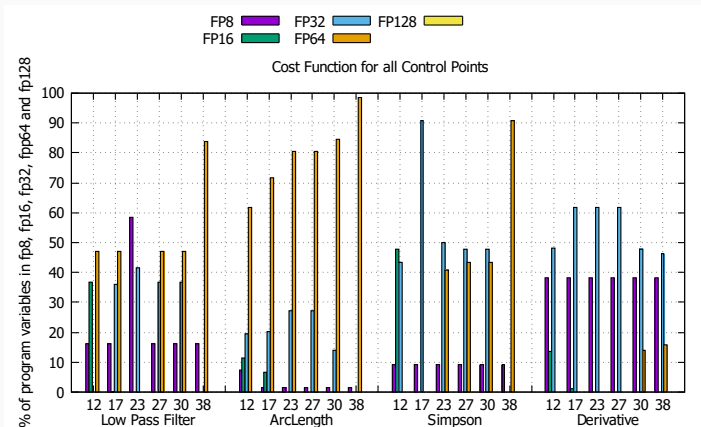
¹⁵D. Ben Khalifa and M.Martel. "Precision Tuning and Internet of Things".IINTEC'19.

¹⁶D.Ben Khalifa and M. Martel. "Precision Tuning of an Accelerometer-Based Pedometer Algorithm for IoT Devices".IoTalS'20

Pedometer Program Mixed-Precision in FP8, FP16, FP32 and FP64

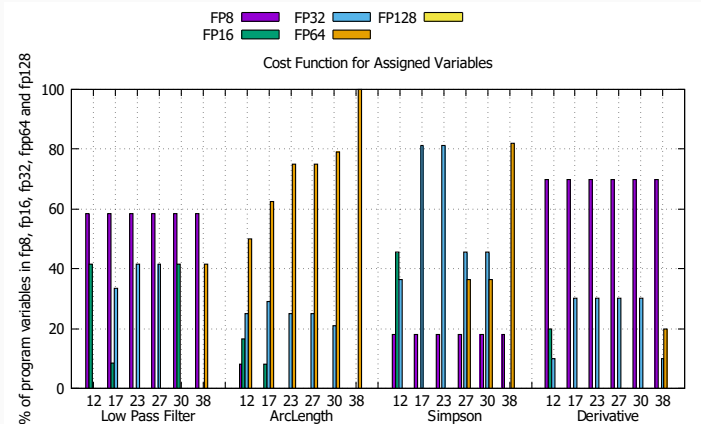


- **Goal:** measure the percentage of program variables in FP8, FP16, FP32 and FP64 after POP(SMT) analysis



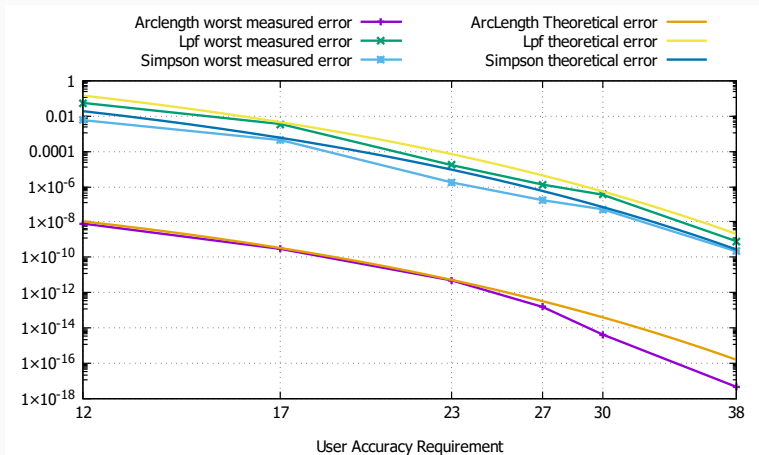
- Cost function **with all control points**
- Results in mixed precision for nsb = 12, 17, 23, 27, 30 and 38 bits

¹⁷D. Ben Khalifa and M. Martel. "An Evaluation of POP Performance for Tuning Numerical Programs in Floating-point Arithmetic". ICICT'21.



- More comparable cost function **with only assigned variables** \implies same as in Precimonious
- 📦 Better optimization of the number of bits of programs

¹⁷D. Ben Khalifa and M. Martel. "An Evaluation of POP Performance for Tuning Numerical Programs in Floating-point Arithmetic". ICICT'21.



- **Goal:** Curve of the difference between the exact results (300 bits) and the precision of POP(SMT) **sticks** to the theoretical curve by below

- 1 Preliminary Elements

- 2 Constraint Generation

- ⋮
- 3 Code Generation**

 - Evaluation POP(ILP)
- 4 Conclusion and Perspectives

Precision Tuning Results for ILP and PI Methods¹⁸

Program	TH	BL	IEEE	ILP-time	BL	IEEE	PI-time	H	S	D	LD
accelerometer	10 ⁻⁴	73%	61%	0.2s	76%	62%	1.0s	53	69	0	0
	10 ⁻⁶	62%	55%	0.2s	65%	55%	1.0s	2	102	0	0
	10 ⁻⁸	49%	15%	0.2s	52%	18%	1.0s	2	33	69	0
	10 ⁻¹⁰	36%	1%	0.2s	39%	1%	1.0s	2	0	102	0
	10 ⁻¹²	25%	1%	0.2s	28%	1%	1.0s	2	0	102	0
lowPassFilter	10 ⁻⁴	68%	46%	1.8s	69%	46%	10.7s	260	581	0	0
	10 ⁻⁶	57%	38%	1.8s	58%	45%	11.0s	258	580	3	0
	10 ⁻⁸	44%	-7%	2.0s	45%	-7%	11.4s	258	2	581	0
	10 ⁻¹⁰	31%	-7%	1.7s	32%	-7%	10.9s	258	0	583	0
	10 ⁻¹²	20%	-7%	1.8s	21%	-7%	11.3s	258	0	583	0
2-Body	10 ⁻⁴	51%	41%	0.81s	51%	41%	0.82s	5	39	5	0
	10 ⁻⁶	49%	18%	0.78s	49%	18%	0.9s	0	44	5	0
	10 ⁻⁸	-7%	5%	0.8s	-7%	5%	0.78s	0	5	44	0
	10 ⁻¹⁰	-34%	-2%	0.8s	-34%	-2%	0.9	0	0	48	1
	10 ⁻¹²	-57%	-11%	0.9s	-57%	-11%	1.0s	0	0	44	0

Parameters

- **TH:** Error threshold
- **IEEE** Optimization in IEEE754
- **H, S:** FP16, FP32 precision
- **BL:** Optimization at Bit level
- **ILP-time/PI-time:** Analysis time
- **D, LD:** FP64, FP128 precision

¹⁸More examples given in Chapter 9 Table 9.1

Precision Tuning Results for ILP and PI Methods¹⁸

Program	TH	BL	IEEE	ILP-time	BL	IEEE	PI-time	H	S	D	LD
accelerometer	10^{-4}	73%	61%	0.2s	76%	62%	1.0s	53	69	0	0
	10^{-6}	62%	55%	0.2s	65%	55%	1.0s	2	102	0	0
	10^{-8}	49%	15%	0.2s	52%	18%	1.0s	2	33	69	0
	10^{-10}	36%	1%	0.2s	39%	1%	1.0s	2	0	102	0
	10^{-12}	25%	1%	0.2s	28%	1%	1.0s	2	0	102	0
lowPassFilter	10^{-4}	68%	46%	1.8s	69%	46%	10.7s	260	581	0	0
	10^{-6}	57%	38%	1.8s	58%	45%	11.0s	258	580	3	0
	10^{-8}	44%	-7%	2.0s	45%	-7%	11.4s	258	2	581	0
	10^{-10}	31%	-7%	1.7s	32%	-7%	10.9s	258	0	583	0
	10^{-12}	20%	-7%	1.8s	21%	-7%	11.3s	258	0	583	0
2-Body	10^{-4}	51%	41%	0.81s	51%	41%	0.82s	5	39	5	0
	10^{-6}	49%	18%	0.78s	49%	18%	0.9s	0	44	5	0
	10^{-8}	-7%	5%	0.8s	-7%	5%	0.78s	0	5	44	0
	10^{-10}	-34%	-2%	0.8s	-34%	-2%	0.9	0	0	48	1
	10^{-12}	-57%	-11%	0.9s	-57%	-11%	1.0s	0	0	44	0

Main observations

- More important BL when coupling POP(ILP) with PI technique
- PI-time > ILP-Time → more constraints to solve
- Ability of POP(ILP) to return new formats for any threshold

¹⁸More examples given in Chapter 9 Table 9.1

POP(ILP) vs POP(SMT) vs Precimonious²⁰

Program	Tool (LOC)	#Bits saved - Time in seconds			
		Threshold 10^{-4}	Threshold 10^{-6}	Threshold 10^{-8}	Threshold 10^{-10}
arclength	POP(ILP) (28)	2464b. - 1.8s.	2144b. - 1.5s.	1792b. - 1.7s.	1728b. - 1.8s.
	POP(SMT) (22)	1488b. - 4.7s.	1472b. - 3.04s.	864b. - 3.09s.	384b. - 2.9s.
	Precimonious (9)	576b. - 146.4s.	576b. - 156.0s.	576b. - 145.8s.	576b. - 215.0s.
simpson	POP(ILP) (14)	1344b. - 0.4s.	1152b. - 0.5s.	896b. - 0.4s.	896b. - 0.4s.
	POP(SMT) (11)	896b. - 2.9s.	896b. - 1.9s.	704b. - 1.7s.	704b. - 1.8s.
	Precimonious (10)	704b. - 208.1s.	704b. - 213.7s.	704b. - 207.5s.	704b. - 200.3s.
rotation	POP(ILP) (25)	2624b. - 0.47s.	2464b. - 0.47s.	2048b. - 0.54s.	1600b. - 0.48s.
	POP(SMT) (22)	1584b. - 1.85s.	2208b. - 1.7s.	1776b. - 1.6s.	1600b. - 1.7s.
	Precimonious (27)	2400b. - 9.53s.	2592b. - 12.2s.	2464b. - 10.7s.	2464b. - 7.4s.
accel.	POP(ILP) (18)	1776b. - 1.05s.	1728b. - 1.05s.	1248b. - 1.04s.	1152b. - 1.03s.
	POP(SMT) (15)	1488b. - 2.6s.	1440b. - 2.6s.	1056b. - 2.4s.	960b. - 2.4s.
	Precimonious (0)	-	-	-	-

Adjusting comparison criteria

- Consider **only** variables optimized by Precimonious
- Express all the error thresholds in base 10
- ...¹⁹

¹⁹ More criteria have been presented in Chapter 8 Section 8.3

²⁰ C. Rubio González, C. Nguyen, H. D. Nguyen, J. Demmel, W. Kahan, K. Sen, D. H. Bailey, C. Iancu, D. Hough. "Precimonious: tuning assistant for floating-point precision". SC'13.

POP(ILP) vs POP(SMT) vs Precimonious¹⁹

Program	Tool (LOC)	#Bits saved - Time in seconds			
		Threshold 10^{-4}	Threshold 10^{-6}	Threshold 10^{-8}	Threshold 10^{-10}
arclength	POP(ILP) (28)	2464b. - 1.8s.	2144b. - 1.5s.	1792b. - 1.7s.	1728b. - 1.8s.
	POP(SMT) (22)	1488b. - 4.7s.	1472b. - 3.04s.	864b. - 3.09s.	384b. - 2.9s.
	Precimonious (9)	576b. - 146.4s.	576b. - 156.0s.	576b. - 145.8s.	576b. - 215.0s.
simpson	POP(ILP) (14)	1344b. - 0.4s.	1152b. - 0.5s.	896b. - 0.4s.	896b. - 0.4s.
	POP(SMT) (11)	896b. - 2.9s.	896b. - 1.9s.	704b. - 1.7s.	704b. - 1.8s.
	Precimonious (10)	704b. - 208.1s.	704b. - 213.7s.	704b. - 207.5s.	704b. - 200.3s.
rotation	POP(ILP) (25)	2624b. - 0.47s.	2464b. - 0.47s.	2048b. - 0.54s.	1600b. - 0.48s.
	POP(SMT) (22)	1584b. - 1.85s.	2208b. - 1.7s.	1776b. - 1.6s.	1600b. - 1.7s.
	Precimonious (27)	2400b. - 9.53s.	2592b. - 12.2s.	2464b. - 10.7s.	2464b. - 7.4s.
accel.	POP(ILP) (18)	1776b. - 1.05s.	1728b. - 1.05s.	1248b. - 1.04s.	1152b. - 1.03s.
	POP(SMT) (15)	1488b. - 2.6s.	1440b. - 2.6s.	1056 - 2.4s.	960b. - 2.4s.
	Precimonious (0)	-	-	-	-

Main observations

- 👉 POP(ILP) saves more bits in fewer time
- 👉 Precimonious displays better results on the rotation example
- 👉 Precimonious fails to tune some benchmarks of POP

¹⁹C. Rubio González, C. Nguyen, H. D. Nguyen, J. Demmel, W. Kahan, K. Sen, D. H. Bailey, C. Iancu, D. Hough. "Precimonious: tuning assistant for floating-point precision". SC'13.

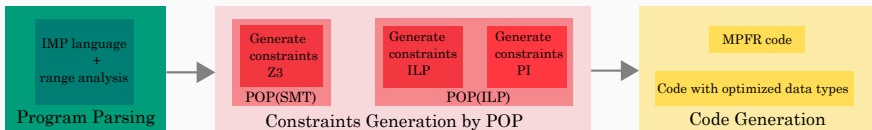
- 1 Preliminary Elements

- 2 Constraint Generation

- ⋮
- 3 Code Generation

- ⋮
- 4 Conclusion and Perspectives**

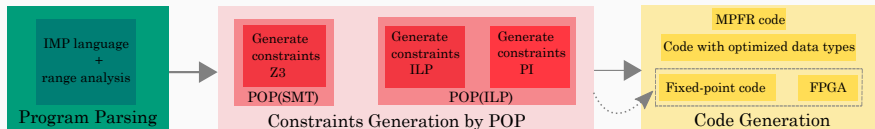
Takeaways



- A new tool for precision tuning with three variants of method
 - Forward and backward error analysis checked by SMT solver
 - Pure ILP with an over-approximation of the carry functions
 - PI technique for more precise carry bit function
- Fast and efficient bit-level precision tuning
- There are still challenges to apply precision tuning **at scale**

Source code available at

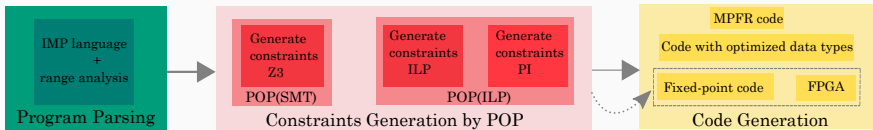
<https://github.com/benkhelifadorra/POP-v2.0>



First research direction

- Synthesis of fixed-point programs²⁰
- Satisfy the accuracy guarantee on the results
- POP computes at each point of this program the pair $|m, s|$
 - m : weight of the most significant bits
 - s : number of significants (nsb)
 - Shifts ?
- Synthesis of VHDL code for FPGA

²⁰Daniele Cattaneo et al. "Fixed point exploitation via compiler analyses and transformations: POSTER". CF'19



Scalability

- Reduce number of variables and constraints
- Explore commercial LP solvers for larger ILP problems

Extension

- Include functions in POP
- Combining POP with rewriting tools to improve accuracy²¹
- Precision tuning for DNN's and GPU applications

²¹N. Damouche, M. Martel. "Mixed Precision Tuning with Salsa". PECCS'18

Thank You...

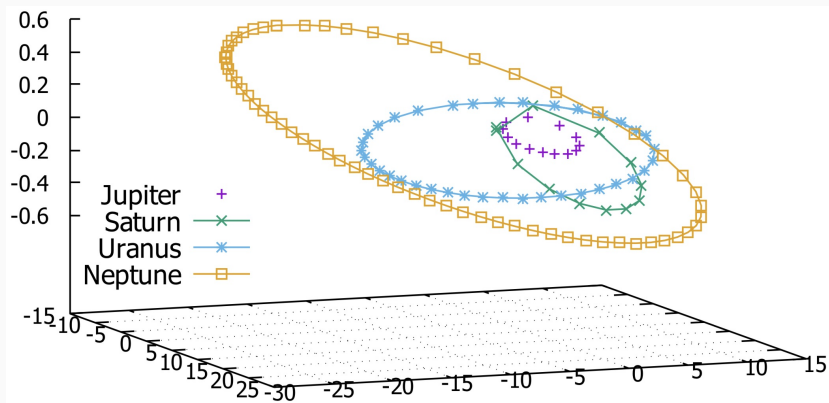
Merci...

شكراً...

List of Publications

- [1] Assalé Adjé, Dorra Ben Khalifa, and Matthieu Martel.
Fast and Efficient Bit-Level Precision Tuning.
In Static Analysis - 28th International Symposium, 2021.
- [2] Dorra Ben Khalifa and Matthieu Martel.
An Evaluation of POP Performance for Tuning Numerical Programs in Floating-Point Arithmetic.
In IEEE International Conference on Information and Computer Technologies, 2021.
- [3] Dorra Ben Khalifa and Matthieu Martel.
Precision Tuning and Internet of Things.
In IEEE International Conference on Internet of Things, Embedded Systems and Communications, 2019.
- [4] Dorra Ben Khalifa and Matthieu Martel.
Precision Tuning of an Accelerometer-Based Pedometer Algorithm for IoT Devices.
In IEEE International Conference on Internet of Things and Intelligence System, 2020.
- [5] Dorra Ben Khalifa and Matthieu Martel.
A Study of the Floating-Point Tuning Behaviour on the N-Body Problem.
In International Conference on Computational Science and Its Applications, 2021.
Springer.
- [6] Dorra Ben Khalifa, Matthieu Martel, and Assalé Adjé.
POP: A Tuning Assistant for Mixed-Precision Floating-Point Computations.
In Formal Techniques for Safety-Critical Systems - 7th International Workshop, 2019.

A Study on the N-Body Problem by POP(ILP)²²



- Simulation of the orbits of planets interacting with each other gravitationally

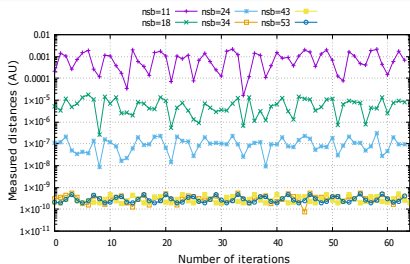
²²D. Ben Khalifa and M. Martel. "A study of the floating-point Tuning Behaviour on the N-body Problem." ICCSA'21

$dist(x, \hat{x})$ for a Single Position Planets

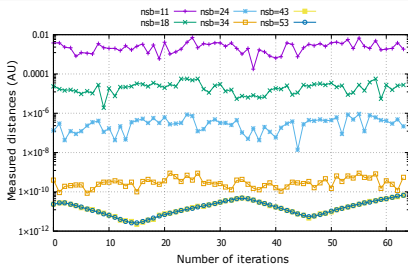
nsb	11	18	24	34	43	53
Simulation time: 10 years						
Jupiter	$5.542 \cdot 10^{-4}$	$1.650 \cdot 10^{-6}$	$1.577 \cdot 10^{-7}$	$4.998 \cdot 10^{-10}$	$5.077 \cdot 10^{-10}$	$5.076 \cdot 10^{-10}$
Saturn	$1.571 \cdot 10^{-3}$	$2.111 \cdot 10^{-5}$	$1.326 \cdot 10^{-7}$	$4.427 \cdot 10^{-10}$	$3.119 \cdot 10^{-10}$	$3.117 \cdot 10^{-10}$
Uranus	$2.952 \cdot 10^{-3}$	$2.364 \cdot 10^{-5}$	$1.140 \cdot 10^{-7}$	$3.072 \cdot 10^{-10}$	$7.212 \cdot 10^{-11}$	$7.236 \cdot 10^{-11}$
Neptune	$2.360 \cdot 10^{-3}$	$3.807 \cdot 10^{-5}$	$2.206 \cdot 10^{-7}$	$5.578 \cdot 10^{-10}$	$1.751 \cdot 10^{-10}$	$1.757 \cdot 10^{-10}$
Runtime	2'59	2'52	2'57	2'56	3'10	2'59
POP Time	25"	22"	22"	24"	23"	24"

- 👉 Satisfactory results which respect the user defined error
- 👉 POP(ILP) analysis time \approx 25 seconds

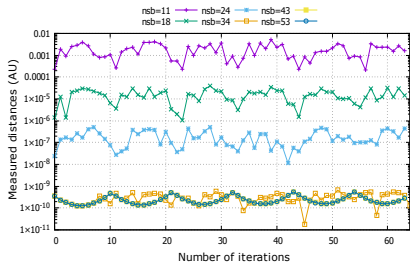
$dist(x, \hat{x})$ at each Instant of Simulation



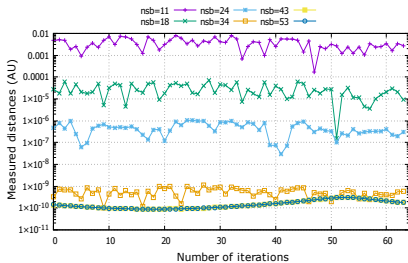
Jupiter



Uranus

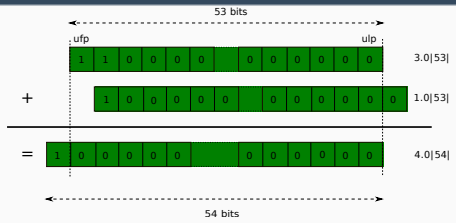


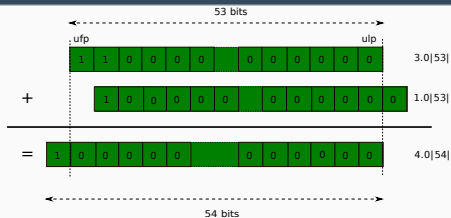
Saturn



Neptune

Transfer Functions Concrete Addition





Forward addition

$$\vec{\oplus}(x_p, y_q) = z_r$$

$$r = \text{ufp}(x + y) - \text{ufp}(\varepsilon_{x_p} + \varepsilon_{y_q} + \varepsilon_+)$$

example $\vec{\oplus}(3.0|53|, 1.0|53|) = 4.0|54|$

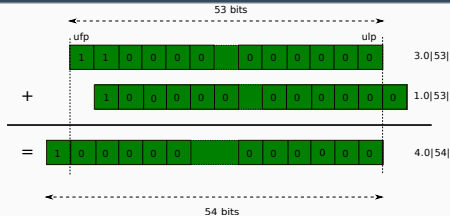
Backward addition

$$\overleftarrow{\oplus}(z_r, y_q) = (z - y)_p$$

$$p = \text{ufp}(z - y) - \text{ufp}(\varepsilon_{z_r} - \varepsilon_{y_q} - \varepsilon_+)$$

example $\overleftarrow{\oplus}(4.0|23|, 1.0|53|) = 3.0|25|$

Transfer Functions Concrete Addition



Forward addition

$$\vec{\oplus}(x_p, y_q) = z_r$$

$$r = \text{ufp}(x + y) - \text{ufp}(\varepsilon_{x_p} + \varepsilon_{y_q} + \varepsilon_+)$$

example $\vec{\oplus}(3.0|53|, 1.0|53|) = 4.0|54|$

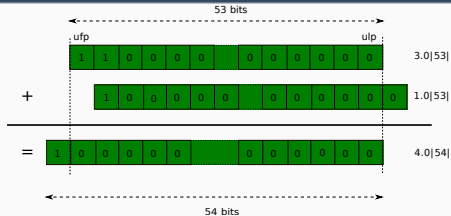
Backward addition

$$\overleftarrow{\oplus}(z_r, y_q) = (z - y)_p$$

$$p = \text{ufp}(z - y) - \text{ufp}(\varepsilon_{z_r} - \varepsilon_{y_q} - \varepsilon_+)$$

example $\overleftarrow{\oplus}(4.0|23|, 1.0|53|) = 3.0|25|$

- **Truncation error:** $\varepsilon_+ \leq 2^{\text{ufp}(z) - \text{prec}(+)}$ ($\text{prec}(+)$ precision of the operator +)



Forward addition

$$\vec{\oplus}(x_p, y_q) = z_r$$

$$r = \text{ufp}(x + y) - \text{ufp}(\varepsilon_{x_p} + \varepsilon_{y_q} + \varepsilon_+)$$

example $\vec{\oplus}(3.0|53|, 1.0|53|) = 4.0|54|$

Backward addition

$$\overleftarrow{\oplus}(z_r, y_q) = (z - y)_p$$

$$p = \text{ufp}(z - y) - \text{ufp}(\varepsilon_{z_r} - \varepsilon_{y_q} - \varepsilon_+)$$

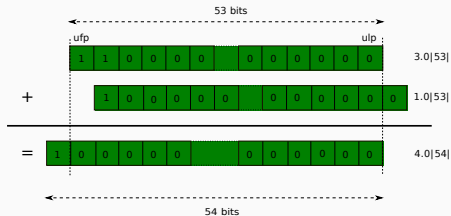
example $\overleftarrow{\oplus}(4.0|23|, 1.0|53|) = 3.0|25|$

- **Truncation error:** $\varepsilon_+ \leq 2^{\text{ufp}(z) - \text{prec}(+)}$ ($\text{prec}(+)$ precision of the operator +)

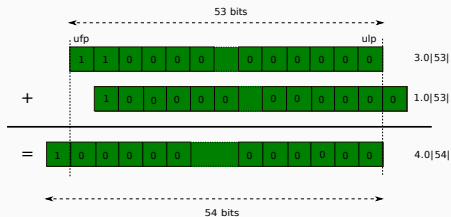
☞ Same for the multiplication, subtraction and division

☞ Technique generalizable to set of values

Constraints SMT Forward Addition Case

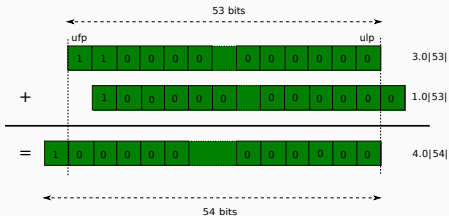


$$\varepsilon_z \leq 2^{\text{ufp}(x) - \text{nsb}_F(x) + 1} + 2^{\text{ufp}(y) - \text{nsb}_F(y) + 1} + 2^{\text{ufp}(z) - \text{prec}(+)}$$



$$\varepsilon_z \leq 2^{\text{ufp}(x) - \text{nsb}_F(x) + 1} + 2^{\text{ufp}(y) - \text{nsb}_F(y) + 1} + 2^{\text{ufp}(z) - \text{prec}(+)}$$

$$\text{ufp}(\varepsilon_z) \leq \max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+)) + \text{carry}(z, x, y)$$



$$\varepsilon_z \leq 2^{\text{ufp}(x) - \text{nsb}_F(x) + 1} + 2^{\text{ufp}(y) - \text{nsb}_F(y) + 1} + 2^{\text{ufp}(z) - \text{prec}(+)}$$

$$\text{ufp}(\varepsilon_z) \leq \max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+)) + \text{carry}(z, x, y)$$

$$\text{nsb}_F(z) = \text{ufp}(x+y) - \max(\text{ufp}(x) - \text{nsb}_F(x), \text{ufp}(y) - \text{nsb}_F(y), \text{ufp}(z) - \text{prec}(+)) - \text{carry}$$

example $\vec{\oplus}(3.0|53|, 1.0|53|) = 4.0|54|$

$$\text{nsb}_B(x) = \text{ufp}(z - y) - \text{ufp}(z) + \text{nsb}(z)$$

example $\overleftarrow{\oplus}(4.0|23|, 1.0|53|) = 3.0|25|$

$$0 \leq \text{nsb}_B(\ell) \leq \text{nsb}(\ell) \leq \text{nsb}_F(\ell), \quad \forall \ell \in \text{Lab}$$

- **Forward multiplication** $\vec{\otimes}$ between x and y in z

$$\vec{\otimes}(x, y) = z \text{ where } \text{nsb}_F(z) = \text{ufp}(x \times y) - \text{ufp}(y \cdot \varepsilon(x) + x \cdot \varepsilon(y) + \varepsilon(x) \cdot \varepsilon(y) + \varepsilon_x)$$

example $\vec{\otimes}(4.0|53|, 1.0|53|) = 4.0|53| \implies \text{nsb}_F(z) = 53$

²³D. Ben Khalifa, M. Martel and A. Adjé. "POP: A Tuning Assistant for Mixed-Precision Floating-Point Computations". FTSCS'19

- **Forward multiplication** $\vec{\otimes}$ between x and y in z

$$\vec{\otimes}(x, y) = z \text{ where } \text{nsb}_F(z) = \text{ufp}(x \times y) - \text{ufp}(y \cdot \varepsilon(x) + x \cdot \varepsilon(y) + \varepsilon(x) \cdot \varepsilon(y) + \varepsilon_x)$$

example $\vec{\otimes}(4.0|53|, 1.0|53|) = 4.0|53| \implies \text{nsb}_F(z) = 53$

- **Backward multiplication** $\overleftarrow{\otimes}$ between x and y in z

$$\overleftarrow{\otimes}(z, y) = (z \div y) \text{ where } \text{nsb}(x) = \text{ufp}(z \div y) - \text{ufp}\left(\frac{y \cdot \varepsilon(z) - z \cdot \varepsilon(y)}{y \cdot (y + \varepsilon(y))} - \varepsilon_x\right)$$

example $\overleftarrow{\otimes}(4.0|23|, 1.0|53|) = 4.0|25| \implies \text{nsb}_B(z) = 53$

²³D. Ben Khalifa, M. Martel and A. Adjé. "POP: A Tuning Assistant for Mixed-Precision Floating-Point Computations". FTSCS'19

- **Forward multiplication** $\vec{\otimes}$ between x and y in z

$$\vec{\otimes}(x, y) = z \text{ where } \text{nsb}_F(z) = \text{ufp}(x \times y) - \text{ufp}(y \cdot \varepsilon(x) + x \cdot \varepsilon(y) + \varepsilon(x) \cdot \varepsilon(x) + \varepsilon_x)$$

example $\vec{\otimes}(4.0|53|, 1.0|53|) = 4.0|53| \implies \text{nsb}_F(z) = 53$

- **Backward multiplication** $\overleftarrow{\otimes}$ between x and y in z

$$\overleftarrow{\otimes}(z, y) = (z \div y) \text{ where } \text{nsb}(x) = \text{ufp}(z \div y) - \text{ufp}\left(\frac{y \cdot \varepsilon(z) - z \cdot \varepsilon(y)}{y \cdot (y + \varepsilon(y))} - \varepsilon_x\right)$$

example $\overleftarrow{\otimes}(4.0|23|, 1.0|53|) = 4.0|25| \implies \text{nsb}_B(z) = 53$

Generalization

- Same for the subtraction and division
- Technique generalizable to set of values²³

²³D. Ben Khalifa, M. Martel and A. Adjé. "POP: A Tuning Assistant for Mixed-Precision Floating-Point Computations". FTSCS'19

POP (both versions)

- Constraint generation by static analysis
- Optimized formats given by solvers
- Mixed-precision: FP8, FP16, FP32, FP64, FPxx
- Programmer accuracy requirement
- Supports arrays, conditions, loops, no function yet

Precimonious

- Dynamic analysis by delta-debugging search
- Type configurations rely on inputs tested **only**
- FP32 and FP64
- Given programmer error threshold (10^{-4} , 10^{-6} , 10^{-8} , ...)
- C program input

²⁴C. Rubio González, C. Nguyen, H. D. Nguyen, J. Demmel, W. Kahan, K. Sen, D. H. Bailey, C. Iancu, D. Hough. "Precimonious: tuning assistant for floating-point precision". SC'13.