# EIFFEL: Inferring Input Ranges of Significant Floating-point Errors via Polynomial Extrapolation

Zuoyan Zhang
*Information Engineering University*
Zhengzhou, China
zhangzuoyan523@163.com

Bei Zhou
*Information Engineering University*
Zhengzhou, China
beibei_0812@126.com

Jiangwei Hao
*Information Engineering University*
Zhengzhou, China
haojiangweitimo@foxmail.com

Hongru Yang
*Information Engineering University*
Zhengzhou, China
hpcyhr@163.com

Mengqi Cui
*Information Engineering University*
Zhengzhou, China
saki1340@163.com

Yuchang Zhou
*Information Engineering University*
Zhengzhou, China
zyc_1013@163.com

Guanghui Song
*Information Engineering University*
Zhengzhou, China
sheensong@163.com

Fei Li
*Information Engineering University*
Zhengzhou, China
feili2022@hotmail.com

Jinchen Xu
*Information Engineering University*
Zhengzhou, China
atao728208@126.com

Jie Zhao
*Information Engineering University*
Zhengzhou, China
yaozhujiajie@gmail.com

*Abstract*—Existing search heuristics used to find input values that result in significant floating-point (FP) errors or small ranges that cover them are accompanied by severe constraints, complicating their implementation and restricting their general applicability. This paper introduces an error analysis tool called EIFFEL to infer error-inducing input ranges instead of searching them. Given an FP expression with its domain $\mathcal{D}$, EIFFEL first constructs an error data set by sampling values across a smaller domain $\mathcal{R}$ and assembles these data into clusters. If more than two clusters are formed, EIFFEL derives polynomial curves that best fit the bound coordinates of the error-inducing ranges in $\mathcal{R}$, extrapolating them to infer all target ranges of $\mathcal{D}$ and reporting the maximal error. Otherwise, EIFFEL simply returns the largest error across $\mathcal{R}$. Experimental results show that EIFFEL exhibits a broader applicability than ATOMU and S³FP by successfully detecting the errors of all 70 considered benchmarks while the two baselines only report errors for part of them. By taking as input the inferred ranges of EIFFEL, Herbie obtains an average accuracy improvement of 3.35 bits and up to 53.3 bits.

## I. INTRODUCTION

FP errors are an infamous problem in software development, often caused by the inaccurate representation of real numbers using limited precision and amplified by the arithmetic operations between these numbers. Detecting them is important since they can lead to catastrophic software failures [1], [2]. FP error detection methods can be static, dynamic or a hybrid of both. Static analysis [3]–[5] is used to provide sound but over-approximated bounds, while dynamic approaches first localize the input values of significant errors [6], [7] or smaller ranges that cover these inputs [8], [9], and next dynamically measure and report the maximal error.

This paper investigates dynamic FP error detection, which can be illustrated by Fig.1 that searches an input range $r$ or the input value $p$ within a space $\mathcal{D} = \{(x,y) : c_1 \wedge c_2 \wedge c_3\}$, where $c_1$, $c_2$ and $c_3$ are constraints of the two variables $x$ and $y$. The points within $r$ are input values that may trigger significant errors, the maximum of which, $p$, is reported as the final result. $\mathcal{D}$ may be complex and large; $r$ and $p$ could both be many. Hence, while reducing the number of input values (from those of $\mathcal{D}$ to those of $r$ or even $p$) to be measured, such methods require a well-defined heuristic to guide the search, and several heuristics [6]–[10] have been proposed in the past.



Fig. 1: Guided search.

While these search heuristics are backed by solid theoretical foundations, their implementations are usually constrained by practical factors, *making their general applicability restricted.* For example, Herbie [6] randomly samples a limited set of input values to localize error-inducing input values, but it only targets numerically unstable programs due to the large number of FP numbers in $\mathcal{D}$. S³FP [9] adopts a binary guided random testing (BGRT) to find such inputs, but this search heuristic is heavyweight and sometimes loses its effectiveness when $\mathcal{D}$ is large. ATOMU [7] leverages condition numbers to localize inputs of high errors, which requires a considered FP
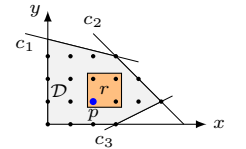
expression to be twice continuously differentiable and makes the computation of condition numbers more difficult [11], [12]. The implementation of ATOMU is thus only able to deal with single-variate FP expressions in practice.

Similar implementation dilemmas also exist in other search heuristics [8], [10], [13]. To address this issue, we present an error analysis approach called EIFFEL, the core idea of which is to *infer error-inducing ranges instead of searching them in* $\mathcal{D}$. Given an FP expression, EIFFEL builds an error data set by sampling input values across a domain $\mathcal{R}$ smaller than $\mathcal{D}$ and measuring the *ulp* errors [14], making it possible to perform data clustering and capture the error distribution features. The data points with larger errors are then partitioned into clusters using the DBSCAN algorithm [15], depending on the number, $num$, of which EIFFEL detects errors as follows:

- If $num \geq 3$, EIFFEL derives two polynomial curves that best fit the lower and upper bounds of an error-inducing range. The remaining ranges in $\mathcal{D}$ are inferred by extrapolating the obtained curves, with the maximum error among all ranges reported as the final result;
- If $num \leq 2$, EIFFEL returns the maximal error across $\mathcal{R}$.

In either case, EIFFEL localizes the input values of significant errors into (one) smaller range(s). In contrast to relying on solid but complex theoretical foundations like prior work, our approach is empirically driven, which simplifies and generalizes our implementation, thereby allowing EIFFEL to detect errors in a wider set of practical scenarios.

We use 66 functions of FPBench [16] and another four benchmarks extracted from real-life numerical programs [17], [18] to conduct experiments, comparing the results with that reported by S³FP [9] and ATOMU [7]. The outcomes indicate that EIFFEL exhibits a wider applicability by detecting the errors of all considered benchmarks, while the two baseline tools only report errors for a partial set of the 70 benchmarks. By comparing the quality and quantity of EIFFEL's inferred input ranges with Regina [19] and PSAT [13], our experiments demonstrate the effectiveness of our range inference strategy. We also feed the inferred input ranges to Herbie [6], which achieves an average accuracy improvement of 3.35 bits and up to 53.3 bits by rewriting FP expressions.

In summary, the contributions of this work are as follows.

- We address error detection as a range inference problem rather than modeling it as a guided search, overcoming the limitations of several search heuristics.
- We combine the research of data analysis and software engineering, demonstrating a practical application of polynomial extrapolation in the field of FP error analysis.
- We design and implement EIFFEL based on the above insights and extensively evaluate EIFFEL using 70 benchmarks, achieving better results than the state of the art.

## II. BACKGROUND KNOWLEDGE

### A. *Floating-point Errors and Their Measurements*

The FP representation with limited precision is a choice of computers to denote infinite real numbers in practice. This paper follows IEEE 754 [20] to define such a representation, which is computed using a formula $\text{sign} \times \text{mantissa} \times 2^{\text{exponent}}$, where the sign is a binary value always at the most significant bit. The number represented is positive when the sign is zero or negative otherwise. For a 64-bit (or 32-bit) FP representation, the mantissa or the significand offers a precision of 53 (or 24) bits by occupying the 52 (or 23) least significant bits, with the additional one implicitly inferred by the exponent residing in the 11 (or 8) bits between the sign and mantissa. Special values like $\pm\infty$ and NaN, which are interpreted as not a number and used to denote the result of certain operations like dividing a number by zero, are also expressible by setting all bits of the exponent as ones.

However, a specific FP representation can only express a limited set of numbers, but real numbers are infinite. This mismatch makes that there always exist many computation results that cannot be accurately represented, and such a real number has to be rounded to its closest exactly-represented value, leading to a difference, i.e., the *rounding error*, between its true value. Rounding errors not only exist while representing a real number but are also accumulated and propagated by FP arithmetic operations, which can lead to catastrophic results in practice [1], [2]. Hence, detecting such errors in programs is of vital importance for software development.

To detect FP errors, an approach should be able to measure them. Suppose that the ground-truth value that can be obtained by computers or the oracle is denoted as $o(f)$, and the computed result is represented as $f$. The *absolute error* is measured as

$$error_{abs} = |o(f) - f|. \tag{1}$$

However, absolute errors cannot reflect their significance over $o(f)$. One can measure the *relative error* using

$$error_{rlt} = \left| \frac{o(f) - f}{o(f)} \right|. \tag{2}$$

Relative errors are broadly used by existing error detection tools [6], [7], [9]. Similar to relative errors, there exists another measuring method that makes use of *ulp* to compute an error, i.e., the unit value of the last digit of an FP number. Such an *ulp error* and can be computed as

$$error_{ulp} = \left| \frac{o(f) - f}{\text{ULP}(o(f))} \right| = \left| \frac{o(f) - f}{\frac{\left| d_0.d_1\cdots d_n - \frac{o(f)}{2^{\text{exponent}}} \right|}{2^{n-1}}} \right|, \tag{3}$$

where ULP computes the *ulp* of $o(f)$. There are different *ulp* error definitions [14]. We follow the one defined in [21], with the formula expressed by the denominator of final expression of Eq (3). We assume that $o(f)$ is represented using an $n$-bit mantissa $d_0.d_1\cdots d_n$, with the leading bit $d_0$ implicitly encoded. By computing the absolute error as a multiple of the oracle's *ulp*, Eq (3) reports much larger numbers (see the $y$ axes of Fig. 2) than absolute and relative errors.

One can render different plots of error distribution for an FP expression when using different error measure methods. As an example, Fig. 2 depicts the error distribution of an example function $f(x) = \frac{1}{\sqrt{x+1}+\sqrt{x}}$ when using the above three errors.
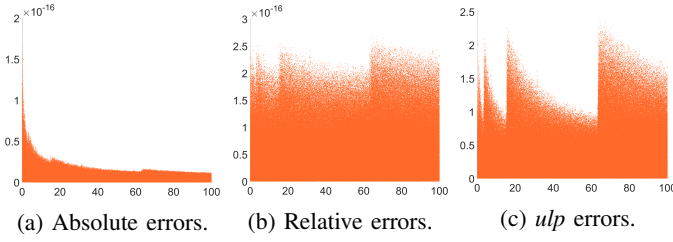
(a) Absolute errors.  (b) Relative errors.  (c) *ulp* errors.

Fig. 2: Error distributions of $f(x) = \frac{1}{\sqrt{x+1}+\sqrt{x}}$.

## B. DBSCAN Clustering

Clustering is a concept of data mining and analysis. It refers to the process of assigning a set of data to different clusters according to their features. Many clustering approaches have been developed before but we concentrate on the DBSCAN (Density Based Spatial Clustering of Applications with Noise) algorithm [15] because it has many features that best fit our purpose: (1) it has no restrictions on how the processed data is structured; (2) the data set can include noise data points; (3) the number of clusters needs not to be specified in advance; (4) the formed clusters can be of arbitrary shapes. Other algorithms like the well-known $k$-means algorithm [22], [23] miss one or several of the above requirements.
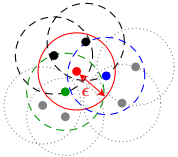


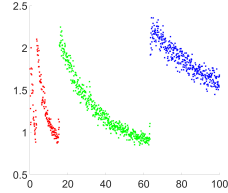Fig. 3: An example used to explain the DBSCAN algorithm, with $MinPts = 4$.



Fig. 4: The clustering result of Fig.2c.

Fig. 3 illustrates the DBSCAN algorithm. Given a set of data points, it initializes a cluster using a randomly selected (red) point that has at least $MinPts$ points (the black, blue and green ones), falling into its (red solid) circle defined by a radius $\epsilon$. The cluster is then magnified by repeatedly inspecting the $n$ points covered by each circle of its internal point: the cluster grows by taking in all of the $n$ points if $n < MinPts$ (e.g., those covered by the blue or green dashed circle) or stops growing otherwise (e.g., the black dashed circles). This iterative step is then applied to the added points (i.e., the gray ones) until the cluster does not grow any more. Those (orange or purple) points not covered by the cluster are then either considered by other clusters or not grouped and thus treated as noises. Hence, one needs to select good values for $MinPts$ and $\epsilon$ when using the DBSCAN algorithm. Fig. 4 shows the result of the application of the DBSCAN algorithm on Fig.2c, with $MinPts$ and $\epsilon$ set to 31 and 1.62, respectively.

## C. Curve Fitting and Extrapolation

Curve fitting is also a technique of data mining and analysis. Given a set of (black) data points in Fig. 5, this technique constructs a (red solid) curve that has a best fit for them by building a mathematical function that approximately fits the data points.

We assume a curve can be expressed as a polynomial function. Once established, this curve can estimate the $y$ value of an input other than those black points used for curve fitting, e.g., the green star point in the original domain $\mathcal{R}$; is can also be extended, as shown by the red dotted curve, to estimate the $y$ value of the (cyan rectangle) input that is outside $\mathcal{R}$ but in $\mathcal{D}$, predicating how the red dotted curve fluctuates across the extended range. This process is referred to as extrapolation.
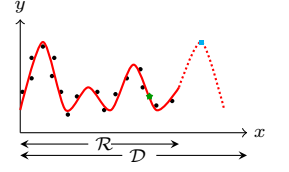


Fig. 5: curving fitting and extrapolation. $x$ axis is the input range; $y$ represents the values.

## D. Design Rationale of EIFFEL

If we consider the $x$ and $y$ axes of Fig.5 as the $\mathcal{D}$ and the error dimension in Fig.2, one can approximate its boundary line via curve fitting and easily determine the $x$ values of the curve peaks, i.e., the significant errors. However, this requires an error plot to be built by sampling a limited number of input values across $\mathcal{D}$, which would not appropriately reflect the curve fluctuation and thus miss some peaks.

To address this issue, we construct an error plot by sampling a dense set of values from a much smaller $\mathcal{R}$, which reflects the error fluctuation more exactly than sampling across $\mathcal{D}$. Fitting the boundary line of an error plot thus has more opportunities to capture the significant errors than a simple dynamic sampling approach [6]. Unfortunately, the number of data points along the boundary line of an error plot may still be many, which would make the implementation tedious like BGRT [9]. Hence, we cluster the data points near the boundary line of an error plot using the DBSCAN algorithm.

The DBSCAN algorithm requires us to automatically determine its $MinPts$ and $\epsilon$ parameters, the latter of which is related to the Euclidean distances between error points. An error point $a$ in one of the plots shown in Fig.2 is represented as its coordinate $(x_a, error_a)$, and its Euclidean distance to another point $b$, denoted as $(x_b, error_b)$, can be computed as $\sqrt{(x_a - x_b)^2 + (error_a - error_b)^2}$. $error_a$ and $error_b$ could be numbers very close to zero if we use Eq (1) or Eq (2) to measure errors, which would make the difference between the error dimension always be zero and thus does not contribute to the computation of Euclidean distances. On the contrary, they are often much larger numbers when using Eq (3) as the measuring method, which can effectively compute Euclidean distances between error points. Considering this fact, we use *ulp* errors to implement our approach.

Once the curve across $\mathcal{R}$ is fitted, it can be extrapolated to infer the $x$ values of remaining peaks in $\mathcal{D}$. In particular, we infer a small range covering such $x$ values to mitigate the inexactness of the curve fitting. The errors of each inferred range are dynamically measured, with the maximal of which reported as the final result. As each step can be implemented

using existing techniques, our work can easily be extended to deal with multi-variate FP expressions, thus exhibiting a wider general applicability than ATOMU [7].

## III. ERROR ANALYSIS USING EIFFEL

Now we explain the details of EIFFEL. We first suppose that the input FP expression is single-variate and will generalize the approach to multi-variate cases in § III-E. We restrict its error analysis within the basic blocks of control constructs like loops and conditional statements, which can be unrolled when necessary and thus are still amenable to EIFFEL. A numerical program instruction can include any FP arithmetic operators and transcendental functions. A function is allowed to be segmented and can take other functions as its inputs.

Given an FP expression, EIFFEL sets $\mathcal{D}$ either with the user-provided domain or assumes $\mathcal{D}$ as the entire feasible range of the FP representation, which is a large domain to search. This usually takes place when the programmer has little background knowledge about numerical analysis. With $\mathcal{D}$ is appropriately set, EIFFEL works as follows.

### A. Data Set Construction

The error data set is the collection of data points in a plot like Fig.2c. Traditionally, the input values used to measure the errors are often selected dynamically across $\mathcal{D}$ [6], [8], which may miss some significant errors when $\mathcal{D}$ is large. EIFFEL collects the *ulp* error distribution across a smaller $\mathcal{R}$, which involves two issues: determining $\mathcal{R}$ and deciding the number of input values to compute *ulp* errors.

FP numbers are non-uniformly distributed across the real number range line. For example, a smaller range $[-1, 1]$ covers 49.95% or 49.61% of all exactly-represented FP numbers when using the IEEE 64- or 32-bit FP representation. Hence, we should let $\mathcal{R}$ be as close to zero as possible such that it can include more exactly-represented numbers, which can better reflect the error distribution of a given FP expression. In EIFFEL, we set $\mathcal{R}$ to $[0, 100]$ by default when $\mathcal{D}$ is the positive half of a real number range line or its sub-domain; otherwise, we let $\mathcal{R}$ be $[-100, 0]$ when $\mathcal{D}$ falls into the negative half. These ranges include 50.27% and 52.17% of all exactly-represented, positive/negative numbers when using IEEE double and single precision. If $\mathcal{D}$ is unknown, we first execute EIFFEL for the positive half and next the negative half. The features observed from $\mathcal{R}$ can thus be used to infer smaller ranges in $\mathcal{D}$. Without loss of generality, we suppose $\mathcal{D}$ be a domain in the positive half of a real number range line and $\mathcal{R} = [0, 100]$ in the following context.

EIFFEL samples $s$ input values from $\mathcal{R}$ and measures their *ulp* errors. Many tools can be used to compute the oracle, among which we choose MPFR [24] due to its acknowledged power. $s$ has an impact on the accuracy of the error distribution, which becomes more accurate with the increase of $s$ and vice versa. We set $s$ to 500,000 by following S³FP [9]. This can obtain a sufficiently accurate error distribution without significant overhead, dynamically sampling with which consumes only 0.17 seconds.

Considering real numbers are distributed uniformly and $s$ is instantiated by a large value, we let the $s$ samples uniformly distribute across $\mathcal{R}$. Indeed, this selection does not perfectly match the distribution of FP numbers in $\mathcal{R}$. We will provide a complementary strategy to alleviate this weakness in § III-D.

### B. Boundary Extraction and Data Clustering

While a larger $s$ can better exhibit the error distribution across $\mathcal{R}$, it also can make the DBSCAN algorithm tedious and EIFFEL impractical. To illustrate the impact of $s$ on the overhead of data clustering, we specify the $MinPts$ and $\epsilon$ parameters of the DBSCAN algorithm using the values in Fig.4. We record the clustering overhead of the DB-SCAN algorithm by changing the value of $s$, obtaining



Fig. 6: The overhead of data clustering and determining the $MinPts$ and $\epsilon$ parameters under different $s$ or $c$ values.

the result shown by the blue curve in Fig.6. The DBSCAN algorithm did not terminate within one hour when $s$ gets larger than 100000, which we treat as a DNF execution. One can see that the clustering overhead grows with the increase of $s$. To demonstrate that the clustering overhead is only influenced by $s$ but not $MinPts$ or $\epsilon$, we also tune the values of these two parameters but observe almost the same result.

To address these conflicting demands, we still let $s$ be a large number as set in § III-A. As the $s$ sampled inputs are uniformly distributed in $\mathcal{R}$, we gather every $g = 500$ samples in one group and only preserve the one that has the largest *ulp* error in this group. Instead of using $s$ data points, we use the $\frac{s}{g}$ largest *ulp* errors of each above group to process the error data set, which obtains the result shown in Fig.7a. Note that how the *ulp* errors fluctuate across $\mathcal{R}$ is determined by the upper boundary of the plot. One can see that what Fig.7a does is just removing the errors much lower than the boundary of Fig.2c. As a result, grouping the $s$ samples this way not only reduces the number of data points but also preserves the exactness of using $s$ samples.
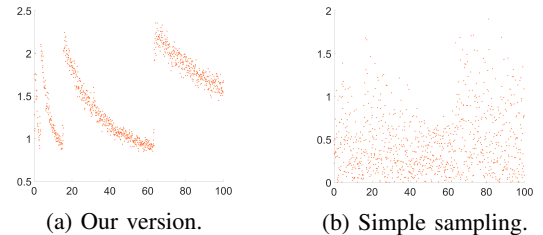


(a) Our version.  (b) Simple sampling.

Fig. 7: The error distribution of $\frac{1}{\sqrt{x+1}+\sqrt{x}}$ using $c$ inputs.

Setting $s$ and $g$ as above makes $c \leq \frac{s}{g}$ equal to 1000. Besides, grouping every $g$ values and computing the largest also do not require much overhead, which can always be
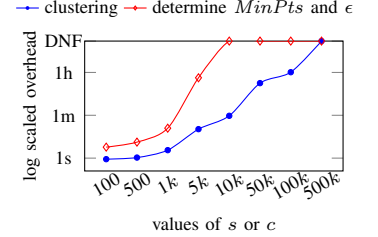
finished in seconds. $c$ can be less than $\frac{s}{g}$ because we do not preserve the largest error of one group if it is less than 0.5 *ulp*, which can be omitted [25]. Note that first densely sampling and next filtering out smaller errors via grouping is essential, and the error distribution cannot be plot directly by simply sampling $c$ inputs from the beginning, since it would obtain an inaccurate error distribution result. As an example, Fig. 7b shows the result of plotting the error distribution of $\frac{1}{\sqrt{x+1}+\sqrt{x}}$ using $c = 1000$ inputs. One can observe that it misses the error fluctuation across $\mathcal{R}$. Curve fitting and extrapolation based on this plot would introduce much inaccuracy.

The DBSCAN algorithm can now be applied to an error plot like Fig. 7a. However, we currently know neither value of $MinPts$ and $\epsilon$. Given a data set composed of $c$ points, $MinPts$ can be specified by an integer $i$ that satisfies $2 \leq i \leq c$. We do not consider $i = 1$ since this value does not need clustering. Hence, the problem is essentially to determine the values of $\epsilon$, denoted as $\epsilon_i$, under each possible value of $i$, which is iterated from 2 to $c$. Once an $\epsilon_i$ is determined, it also corresponds to an $MinPts_i$. That is to say, we can obtain $c - 1$ pairs of $\epsilon_i$ and $MinPts_i$ by iterating $i$.

We denote each data point as $p_j$ ($1 \leq j \leq c$), for which we find a circle centered at $p_j$ and with a radius $e_j$ such that there are $i$ points, including $p_j$ itself, in this circle. For each $p_k$ other than $p_j$, it is possible to compute its Euclidean distance to $p_j$, represented as $d_k$. To minimize $e_j$, we sort these Euclidean distances in a descending order and return the points whose distance ranking in the top $i - 1$, among which the one with the longest Euclidean distance to $p_j$ is used to set $e_j$. As each $p_j$ should be evaluated, we obtain $c$ different $e_j$'s for a given $i$. We use their average mean to set an $\epsilon_i$, computed as

$$\epsilon_i = \frac{\sum_{j=1}^{c} r_j}{c} = \frac{\sum_{j=1}^{c} \left( \max_{k \in \text{top } i-1} d_k \right)}{c} =$$

$$\frac{\sum_{j=1}^{c} \left( \max_{k \in \text{top } i-1} \sqrt{(x_j - x_k)^2 + (error_j - error_k)^2} \right)}{c}, \quad (4)$$

Still consider the data points in Fig. 3, which is reproduced in Fig. 8 for illustrative purpose. We assume $i = 4$ and still use the red solid circle whose radius is the $\epsilon_i$ to be computed. To let the circles of each $p_j$ cover at least four points, the green and orange points can find two radii $e_g$ and $e_o$ that form the green and orange dashed circles, respectively. We consider one point is covered by a circle when it is fully included but not on the boundary line. The green and orange circles are smaller than the red solid one. However, the blue point needs a larger radius $e_b$ to cover at least four points, which is bigger than the red solid one. Hence, computing the average mean of these radii is an approach fair to all points, because the data points are not labeled and should be treated equally. The average mean is finally used to set as the radius of the red solid circle.
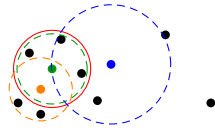


Fig. 8: An example used to compute $\epsilon_i$.

This $\epsilon_i$ is now used as the radius of each circle centered at $p_j$, and the number of data points covered by this circle is computed, which is denoted as $n_j$. Similar to $e_j$, there should be $c$ different $n_j$'s for a given $i$, and we also compute $MinPts_i$ as the average mean of these $n_j$'s, i.e.,

$$MinPts_i = \left\lfloor \frac{\sum_{j=1}^{c} n_j}{c} \right\rfloor, \quad (5)$$

where the floor operator is mandatory as $MinPts_i$ is an integer.

Finally, we select the best ones from the $c - 1$ pairs of values. To achieve this, we use each of these $c - 1$ pairs of values to perform data clustering, and record the number of obtained clusters, denoted as $num_i$. By doing so, we can maintain a

TABLE I: The lookup table.

| $i$ | $\epsilon_i$ | $MinPts_i$ | $num_i$ |
|---|---|---|---|
| 2 | $\epsilon_2$ | $MinPts_2$ | $num_2$ |
| $\cdots$ | $\cdots$ | $\cdots$ | $\cdots$ |
| $m$ | $\epsilon_m$ | $MinPts_m$ | $num_m$ |
| $m+1$ | $\epsilon_{m+1}$ | $MinPts_{m+1}$ | $num_{m+1}$ |
| $m+2$ | $\epsilon_{m+2}$ | $MinPts_{m+2}$ | $num_{m+2}$ |
| $\cdots$ | $\cdots$ | $\cdots$ | $\cdots$ |
| $c$ | $\epsilon_c$ | $MinPts_c$ | $num_c$ |

lookup table, as shown by Table I. Each parameter will be instantiated by a specific value. In the earlier iterations, $num_i$ could be large because its two parameters are too small, which partitions the $c$ data points into more clusters. $num_i$ decreases with the growing of $i$ but gradually becomes stable. We set a condition that a value $m$ of $i$ exists such that $num_m = num_{m+1} = num_{m+2}$ holds. We just let $MinPts = num_i$ and $\epsilon = \epsilon_m$, since the result is already stable.

The above procedure obtains appropriate values for $MinPts$ and $\epsilon$. Indeed, the efficiency of this procedure is also influenced by the total number of data points, because it evaluates all of them to compute the Euclidean distances and iterates as many times as the number of the data points. To demonstrate this impact, we also collect the overhead of the above procedure under different values of $c$, with the result shown by the red line of Fig. 6. The overhead is less than one minute (1m) when $c$ is smaller than 1000, which is affordable in practice.

### C. Curve Derivation and Polynomial Extrapolation

With its parameters determined, the DBSCAN algorithm partitions Fig. 7a into $num = 3$ clusters, as shown in Fig. 9. The (star) peaks can now be extracted to perform curve fitting. Once found, the (thick solid gray) curve can be extrapolated to infer the each peak whose $x$ values fall in $\mathcal{D}$ but outside $\mathcal{R}$, i.e., the black points. The one with the maximal $error$ value among all peaks is returned as the final result. Note that data clustering is essential. The error plots shown throughout this paper are only used to help readers better understand how EIFFEL works but are not visible to our tool. Without data clustering, EIFFEL would not be able to extract these peaks, thus failing to perform curve fitting.

Detecting the maximal error this way faces two challenges. First, each peak point is represented as a 2D coordinate, deriving a mathematical function for the gray curve is thus difficult. EIFFEL may fail to find such a mathematical function, or would require high execution overhead and introduce significant inexactness even if such a function can be found. Second, it is very likely to happen that the extrapolated points still
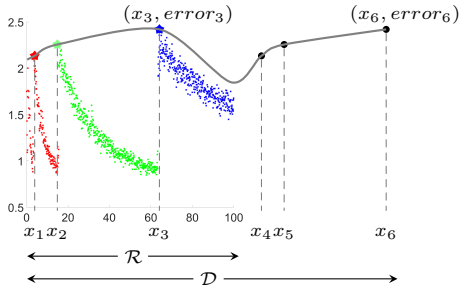
Fig. 9: Fitting and extrapolating peak points.

follow the same distribution as that of the stars, which makes the error values of the extrapolated points equal to those of the stars. For example, suppose that $(x_3, error_3)$ has the largest error among the three stars. The largest of those extrapolated points should be $(x_6, error_6)$ and $error_3 = error_6$ holds. Extrapolating the built curve would thus make no sense.

To resolve the first issue, we do not perform curve fitting on the peak points; instead, we record their $x_i$ values, i.e., $x_1$, $x_2$ and $x_3$, derive their distribution along $\mathcal{R}$ and infer the $x$ values ($x_4$, $x_5$, $x_6$ and more if any) of the extrapolated points. This simplifies curve fitting by converting the 2D coordinates into a 1D form, but extrapolation based on which may still be overfitting, i.e., the second issue. To address this, we assume the $x_i$ values of these peak points form a geometric sequence or a geometric progression, written as

$$x_i = x_1 \times q^{i-1} \qquad \text{s.t.} \qquad i \geq 1, \qquad (6)$$

where $q$ is the common ratio. Note that Eq (6) is the polynomial function of $x_i$ over $i$, and extrapolating it can be used to infer those unknown $x_i$'s without performing curve fitting on coordinates like $(x_i, err_i)$. While simplifying the algorithmic design of EIFFEL, this process also forms the reason why our approach is called polynomial extrapolation. We make this assumption based on the results observed from extensive experiments. Fig. 10 shows two FPbench benchmarks whose high error-inducing $x$ values exhibit such a property. Most practical benchmarks also exhibit a similar property but we did not show them here due to the limited space.



(a) NMSEexample36; $q \approx 1.59$.  (b) NMSEproblem331; $q \approx 2.03$.
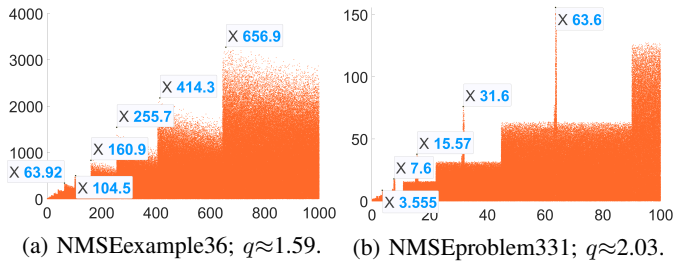
Fig. 10: The geometric progression property of two examples.

Considering this assumption is made based on experimental observations, we compute a lower bound $x_i - r$ and an upper bound $x_i + r$, which form a range $[x_i - r, x_i + r]$ that covers each $x_i$. EIFFEL instantiates $r$ using 1.0 but it can be assigned

using other values. Such a range is small with respect to the default $\mathcal{R}$. In particular, the lower bound is instantiated by $x_1$ itself when $x_1 < 2r$. Such a handling is used to avoid the case where the starting term of a geometric progression gets close or equal to zero, which either makes the computation of the geometric progression's following terms tedious or violates the definition of a geometric progression. Finally, we derive two curves that best fit the lower and upper bounds of each range.

As an example, Fig. 11 shows the curve fitting and extrapolation results of Fig. 9. EIFFEL first obtains the $x_i$ values of the peaks shown in Fig. 9 and next derives bounds for each $x_i$, based on which the value of $q$ in Eq (6) that best fits each lower or upper bound is approximated, leading to the two curves shown in the left plot. We suppose that $\mathcal{D}$ is $[0, 100000]$ and let EIFFEL extrapolate the built curves, with addition five pairs of lower and upper bounds extrapolated.
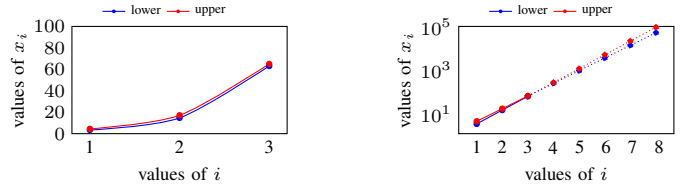


Fig. 11: Fitting (left) and extrapolating (right) the bounds of $x_i$. The $y$ axis of the right plot is log scaled.

### D. Error Detection

As all ranges are inferred, dynamically inspecting the largest error in each of such smaller ranges is possible, among which the maximum can be returned as the final result. Nonetheless, curve fitting cannot be executed when $num = 1$, which implies there exist only one $x_i$ that is not sufficient to derive a curve. Also, the inferred ranges should be considered as not reliable when $num = 2$, since deriving a curve using only two points has a higher risk of inaccuracy.



(a) $f(x) = \frac{-x^3}{6}$.  (b) $f(x) = \sqrt{\frac{e^{2x}-1}{e^x-1}}$.  (c) $f(x) = \frac{1}{x+1}$.
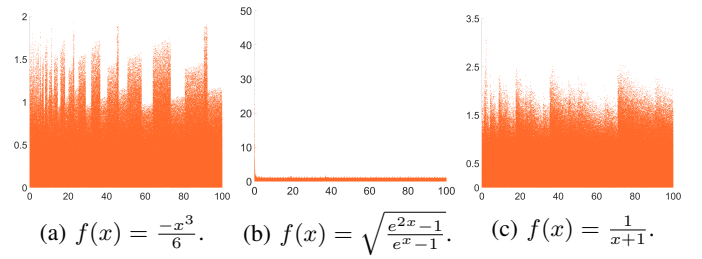
Fig. 12: The error distribution plot of different examples.

When these two cases are encountered, we just return the dynamically detected, maximal error across $\mathcal{R}$ as the result. This is reasonable because an FP expression that makes the DBSCAN algorithm form one or two clusters usually observes its maximal error in $\mathcal{R}$. Consider Fig. 12a as an example. Due to its insignificant fluctuations, the DBSCAN algorithm assigns the larger error points into one cluster. We can consider that such insignificant fluctuations also exist in $\mathcal{D}$, because

this FP expression is numerically stable. Another example is Fig. 12b, which is numerically unstable. If the significant errors are considered as noises, EIFFEL obtains only one cluster; otherwise, they form two clusters. In both cases, the maximal error of $\mathcal{R}$ is very likely to be the result of any domains. Hence, EIFFEL can determine the error detection methods according to the values of $num$.

Another feature we make use of is the monotonicity of an FP expression: we can detect the errors in a small range $[lb_{\mathcal{D}}, up_{r_1}]$ or $[lb_{r_N}, up_{\mathcal{D}}]$ if an FP expression is monotonous, where $lb_{\mathcal{D}}$ and $up_{\mathcal{D}}$ are the lower and upper bounds of $\mathcal{D}$, $up_{r_1}$ is the upper bound of the leftmost range, and $lp_{r_N}$ is the lower bound of the rightmost range. Monotonicity is determined by checking whether the maximal errors of each range are ordered in a(n) ascending/descending order. Combining $num$ and the monotonicity produces four quadrants shown in Fig. 13a, and an FP expression always falls into one of them.



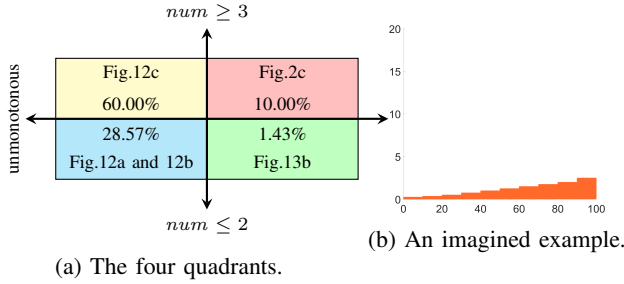(a) The four quadrants.

(b) An imagined example.

Fig. 13: Determining error detection methods according to the values of $num$ and the monotonicity of an expression.

We show examples in each quadrant for illustrative purpose and mark the percentage statistics by counting the 70 expressions used in our experiments. The colors are used to denote how the maximal error is determined. Yellow expresses that a sufficient number of clusters are obtained, and the extrapolated result across $\mathcal{D}$ is returned. Red inspects the errors in $[lb_{\mathcal{D}}, up_{r_1}]$ or $[lb_{r_N}, up_{\mathcal{D}}]$. Cyan indicates that the maximal error across $\mathcal{R}$ is reported.

EIFFEL may report false positives for some examples, e.g., Fig. 12c. EIFFEL obtains more than three clusters, and its extrapolation shows that there also exist other larger errors as shown in Fig. 12c. One can conclude from its expression that there should have no such ranges in the positive half real number range line, though it has significant errors at $x = -1.0$. However, one can notice that the errors are all smaller than 2 *ulp*, which is considered as insignificant according to the GNU C library manual [25]. EIFFEL reports false negatives when a monotonous FP expression falls into the green quadrant, which EIFFEL partitions into one or two clusters. A manually crafted example (which does not exist in practice) is shown in Fig. 13b, whose errors increase steady but slowly.

### E. Generalization for Multi-variate Scenarios

Handling multi-variate expression is similar to the single-variate case, but two adaptations should be introduced. First, since curve fitting on multi-dimensional data is difficult, we do not extract the largest errors of each group during the grouping optimization but project the multi-dimensional plot onto the variable space, which produces the variable coordinates of the errors. For example, Fig. 14a is the 3D error distribution plot of the $Hypot(x, y)$ benchmark extracted from FPBench [16]. EIFFEL obtains the projection result in Fig. 14b. The DBSCAN algorithm is then applied to these points without changes, producing three clusters.



(a) Error distribution.
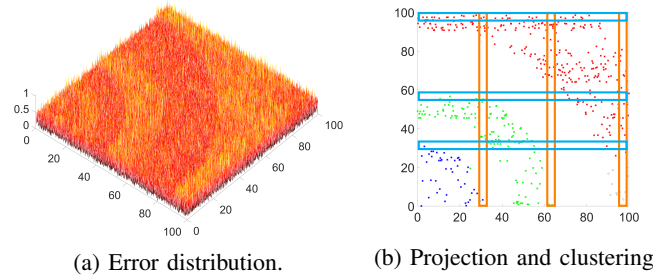
(b) Projection and clustering.

Fig. 14: Error distribution of $Hypot(x, y)$ and its projection.

Second, as the coordinates of each point in the projection are still multi-dimensional, we perform curve fitting along one dimension each time. This always produces (hyper-)rectangular ranges but is still more effective than existing approaches. For instance, we first derive curves along the $x$ dimension of Fig. 14b, which produces the orange rectangular ranges. The error points with the maximal $x$ values in each cluster are used for curve fitting. Next, we derive curves along the $y$ dimension, obtaining the cyan rectangular ranges. Also, the error points with the maximal $y$ values in each cluster are used for curve fitting. Both versions are delivered to the error detection step, and the one with the larger error is finally reported.

## IV. EXPERIMENTAL RESULTS

EIFFEL uses 360 lines of Python v3.8.10 code to determine $num$ and perform curve fitting/extrapolation, and about 1500 lines of C/C++ code to implement the remaining components. Its code is available at https://github.com/zuoyanzhang/EIFFEL. The experimental environment is Ubuntu 20.04.4 LTS based on the Linux 5.14.0-1051-oem kernel, running on an Intel Xeon E5-6230 v4 CPU. EIFFEL generates a C program that invokes a call to MPFR [24] to measure errors and is compiled by GCC 9.4.0 with options "-lm -lmpfr" enabled.

The benchmarks are summarized in Table II, out of which 66 are extracted from FPBench [16]. We only consider the expressions in the basic block of a control construct, and the if conditionals can be unrolled when necessary as handled in the similar way of Regina [19]. Benchmarks whose expressions are equivalent to one of Table II are not included for the sake of simplicity. Besides, those expressions containing the $pow(x, y)$ function are not considered, because its values easily exceed the maximal representable numbers of a computer when using a large $\mathcal{D}$. However, it does not imply that EIFFEL cannot handle this function, whose maximal error can be easily detected under a small $\mathcal{D}$. The remaining four (multi-variate) benchmarks are extracted from real-life numerical programs,

TABLE II: Summary of the benchmarks. $\mathcal{D}$ is set using large and reasonable ranges, and used for each variable dimension in multi-variate cases, with invalid input values already excluded. We also report the two parameters of the DBSCAN algorithm.

| no. | name | $\mathcal{D}$ | MinPts | $\epsilon$ | no. | name | $\mathcal{D}$ | MinPts | $\epsilon$ | no. | name | $\mathcal{D}$ | MinPts | $\epsilon$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | predatorPrey | $[-10^5, 10^5]$ | 15 | 0.82 | 24 | NMSproblem345 | $[-10^5, 10^5]$ | 998 | 3.99e05 | 47 | complex_sine_cosine | $[10^{-2}, 708]$ | 994 | 113.45 |
| 2 | sqrt_add | $[0, 10^5]$ | 26 | 1.36 | 25 | NMSsection311 | $[-10^5, 708]$ | 62 | 75.70 | 48 | jetEngine | $[10^{-2}, 10^5]$ | 47 | 12.37 |
| 3 | verhulst | $[-10^5, 10^5]$ | 24 | 1.25 | 26 | carbonGas | $[-10^5, 10^5]$ | 10 | 0.54 | 49 | rump's example | $[10^{-2}, 10^5]$ | 930 | 81.48 |
| 4 | test05_nonlin1_r4 | $[-10^5, 10^5]$ | 999 | 1.22e03 | 27 | sine | $[-10^5, 10^5]$ | 998 | 4.30e04 | 50 | doppler1 | $[10^{-2}, 10^5]$ | 81 | 13.35 |
| 5 | test05_nonlin1_test2 | $[-10^5, 10^5]$ | 8 | 0.61 | 28 | sqroot | $[-10^5, 10^5]$ | 998 | 5.33e04 | 51 | sum | $[-10^5, 10^5]$ | 31 | 7.90 |
| 6 | exp1x | $[-10^5, 708]$ | 20 | 1.07 | 29 | sineOrder3 | $[-10^5, 10^5]$ | 998 | 5.33e04 | 52 | rigidBody1 | $[-10^5, 10^5]$ | 982 | 103.00 |
| 7 | exp1x_log | $[-10^5, 708]$ | 20 | 1.07 | 30 | bsplines3 | $[-10^5, 10^5]$ | 937 | 74.94 | 53 | rigidBody2 | $[-10^5, 10^5]$ | 14 | 12.40 |
| 8 | logexp | $[-10^5, 708]$ | 299 | 24.60 | 31 | nonlin2 | $[10^{-2}, 10^5]$ | 19 | 15.86 | 54 | turbine1 | $[-10^5, 10^5]$ | 14 | 11.05 |
| 9 | intro-example-minxed | $[-10^5, 10^5]$ | 10 | 0.73 | 32 | hypot | $[10^{-2}, 10^5]$ | 12 | 9.36 | 55 | turbine2 | $[-10^5, 10^5]$ | 11 | 10.70 |
| 10 | NMSexample31 | $[-10^5, 4.50e15]$ | 18 | 2.14 | 33 | x_by_xy | $[10^{-2}, 10^5]$ | 18 | 8.91 | 56 | turbine3 | $[-10^5, 10^5]$ | 12 | 11.97 |
| 11 | NMSexample34 | $[-10^5, 10^5]$ | 998 | 2.31e10 | 34 | NMSproblem335 | $[10^{-2}, 10^5]$ | 15 | 6.68 | 57 | test01_sum3 | $[0, 10^5]$ | 31 | 7.91 |
| 12 | NMSexample35 | $[-10^5, 10^5]$ | 57 | 3.17 | 35 | NMSproblem332 | $[10^{-2}, 10^5]$ | 52 | 12.21 | 58 | sphere | $[0, 10^5]$ | 13 | 4.35 |
| 13 | NMSexample36 | $[0, 4.54e15]$ | 210 | 37.77 | 36 | floudas3 | $[10^{-2}, 10^5]$ | 109 | 5.31 | 59 | azimuth | $[0, 10^5]$ | 23 | 8.79 |
| 14 | NMSproblem331 | $[0, 9.27e15]$ | 203 | 15.54 | 37 | himmibeau | $[10^{-2}, 10^5]$ | 27 | 5.50 | 60 | delta4 | $[0, 500]$ | 27 | 9.41 |
| 15 | NMSproblem333 | $[0, 6.87e10]$ | 137 | 221.70 | 38 | carthesianToPolar_theta | $[10^{-2}, 10^5]$ | 47 | 7.54 | 61 | delta | $[0, 500]$ | 21 | 8.44 |
| 16 | NMSproblem334 | $[-2.25e15, 0]$ | 92 | 31.79 | 39 | NMSexample33 | $[10^{-2}, 10^5]$ | 26 | 8.35 | 62 | kepler0 | $[0, 200]$ | 24 | 7.74 |
| 17 | NMSproblem336 | $[0, 1.32e16]$ | 157 | 17.88 | 40 | i4 | $[10^{-2}, 10^5]$ | 19 | 11.00 | 63 | kepler2 | $[0, 200]$ | 21 | 8.44 |
| 18 | NMSproblem337 | $[-10^5, 708]$ | 492 | 7.27e03 | 41 | i6 | $[10^{-2}, 10^5]$ | 981 | 104.40 | 64 | Shoelace formula | $[0, 200]$ | 14 | 5.47 |
| 19 | NMSexample37 | $[-10^5, 708]$ | 20 | 2.36 | 42 | test03_nonlin2 | $[10^{-2}, 10^5]$ | 20 | 6.87 | 65 | matrixDeterminant | $[0, 200]$ | 18 | 7.30 |
| 20 | NMSexample38 | $[0, 7.87e15]$ | 117 | 12.93 | 43 | polarToCarthesian,x | $[10^{-2}, 10^5]$ | 102 | 5.30 | 66 | matrixDeterminant2 | $[0, 200]$ | 12 | 6.03 |
| 21 | NMSexample39 | $[-10^5, 10^5]$ | 998 | 6.97e04 | 44 | polarToCarthesian,y | $[10^{-2}, 10^5]$ | 10 | 5.87 | 67 | pov-ray | $[0, 10^5]$ | 15 | 6.93 |
| 22 | NMSproblem341 | $[-10^5, 10^5]$ | 998 | 2.41e10 | 45 | NMSproblem346 | $[10^{-2}, 10^5]$ | 20 | 9.24 | 68 | polyIDX0 | $[0, 10^5]$ | 15 | 6.23 |
| 23 | NMSproblem344 | $[-10^5, 304]$ | 13 | 0.70 | 46 | complex_square_root | $[10^{-2}, 10^5]$ | 28 | 11.30 | 69 | polyIDX1 | $[0, 10^5]$ | 19 | 6.75 |
| | | | | | | | | | | 70 | polyIDX2 | $[0, 10^5]$ | 22 | 6.10 |

one (*bench* 67) from POV-Ray [17] that uses ray-tracing to render a 3D image and the other three (*bench* 68-70) from a polynomial support vector classifier [18].

The experiments are designed and conducted to answer the following research questions (RQs):

- **RQ1**: How generally applicable and effective is EIFFEL when compared with existing error detection methods?
- **RQ2**: How does EIFFEL differ from the state-of-the-art search heuristics?
- **RQ3**: How are the quality and quantity of the input ranges inferred by EIFFEL different from similar approaches?
- **RQ4**: How scalable is EIFFEL when compared with prior work?

To address **RQ1**, we detect the maximal errors of the considered benchmarks using EIFFEL and compare the results with those reported by S$^3$FP [9] and ATOMU [7]. The differences between EIFFEL and their search heuristics, i.e., **RQ2**, are comprehensively discussed after the comparison. As for **RQ3**, we compare the quality and quantity of EIFFEL's inferred input ranges to those of Regina [19] and PSAT [13], both of which are tools that find input ranges instead of localizing specific input values. To tackle **RQ4**, we record the execution time and compare with some state-of-the-art tools.

### A. Comparison of Detected Errors and General Applicability

To demonstrate the general applicability of our method, we compare EIFFEL with two state-of-the-art tools, S$^3$FP [9] and ATOMU [7], both detecting and reporting maximal relative errors. EIFFEL leverages *ulp* errors for implementation, but it can report the more widely used relative errors, which are listed in Table III.

*1) General Applicability on Error Detection:* S$^3$FP requires to specify a TIMEOUT parameter for controlling the time

budget that constrains its BGRT, which we set to $2 \times 10^6$, limiting the BGRT to 60 minutes as used in its publication [9]. For ATOMU, we redefine $\mathcal{D}$ of *bench* 10, 13-17, and 20 as [0, ATOMU's $x$] or [ATOMU's $x$,0] depending on the sign of ATOMU's $x$, which represents the input value inferred by ATOMU because the originally defined $\mathcal{D}$ does not include this input value. We increase $\mathcal{D}$ to cover ATOMU's $x$ to validate that EIFFEL can also find a range that covers this input value.

EIFFEL successfully detects the maximal errors of all considered benchmarks. In contrast, S$^3$FP returns empty results for 27 of them while ATOMU is only able to report errors for the first 30 single-variate examples. Increasing TIMEOUT allows S$^3$FP to obtain better results but it still cannot report error data for those examples it returns empty results.

We did not show the error data of S$^3$FP and ATOMU since users are more interested in knowing whether such tools can successfully report significant errors instead of the specific numbers. In the best case, EIFFEL detects a larger maximal error (2.62E+05) than S$^3$FP (2.25E-14) for *bench* 15, and or a more significant error (262146.5) than ATOMU (262143) for *bench* 15; in the worst case, EIFFEL shows a smaller error (3.95E-09) than S$^3$FP (3.73E-06) for *bench* 66, and performs worse (1.19E-11) than ATOMU (1.39E+00) for *bench* 27. Nonetheless, the distance between the input values found by EIFFEL (-3.07859999999999445208) and ATOMU (-3.0786423044815128) is very small. EIFFEL observes larger errors than S$^3$FP and ATOMU for most other cases.

*2) Difference between BGRT:* The effectiveness of the BGRT heuristic used by EIFFEL depends heavily on the time budget of the method. On the contrary, our inference process consumes much fewer time to report the result, as shown in Table III, which ranges from 8.02 seconds (*bench* 29) to 1196.44 seconds (*bench* 59). As EIFFEL iteratively performs

TABLE III: Maximal relative errors reported by EIFFEL. We mark a benchmark as either ✓ if S³FP or ATOMU successfully returns an error for it, or × otherwise. The last two columns (time and $N$) report the overhead of EIFFEL in seconds to detect this benchmark and the number of input ranges.

| no. | S³FP | ATOMU | EIFFEL | time | N | no. | S³FP | ATOMU | EIFFEL | time | N |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | ✓ | × | 3.94E-16 | 80.00 | 25 | 36 | × | N/A | 1.58E-11 | 175.44 | 35 |
| 2 | × | × | 3.05E-16 | 56.34 | 13 | 37 | ✓ | N/A | 1.90E-15 | 147.68 | 38 |
| 3 | ✓ | ✓ | 1.02E-01 | 52.30 | 17 | 38 | × | N/A | 2.94E-16 | 191.41 | 28 |
| 4 | ✓ | ✓ | 4.49E-13 | 16.81 | 2 | 39 | × | N/A | 1.67E-05 | 176.03 | 23 |
| 5 | ✓ | ✓ | 1.67E-16 | 61.61 | 25 | 40 | × | N/A | 1.84E-16 | 237.62 | 42 |
| 6 | ✓ | ✓ | 1.10E-01 | 22.01 | 2 | 41 | × | N/A | 1.60E-02 | 23.40 | 1 |
| 7 | × | ✓ | 1.00E+00 | 29.18 | 2 | 42 | ✓ | N/A | 2.39E-16 | 188.07 | 35 |
| 8 | × | ✓ | 1.00E+00 | 27.64 | 2 | 43 | ✓ | N/A | 7.06E-11 | 20.08 | 1 |
| 9 | ✓ | ✓ | 1.66E-16 | 47.95 | 24 | 44 | ✓ | N/A | 1.09E-06 | 657.05 | 75 |
| 10 | × | ✓ | 9.98E-01 | 17.77 | 1 | 45 | × | N/A | 4.82E-10 | 472.02 | 30 |
| 11 | × | ✓ | 1.00E+00 | 24.97 | 2 | 46 | × | N/A | 2.27E-16 | 225.17 | 39 |
| 12 | × | ✓ | 2.46E-09 | 156.71 | 26 | 47 | × | N/A | 8.17E-15 | 27.00 | 1 |
| 13 | × | ✓ | 2.89E+00 | 101.82 | 1 | 48 | ✓ | N/A | 1.07E-15 | 17.46 | 1 |
| 14 | ✓ | ✓ | 1.12E+00 | 68.59 | 1 | 49 | ✓ | N/A | 5.07E-11 | 18.26 | 1 |
| 15 | ✓ | ✓ | 2.62E+05 | 60.20 | 1 | 50 | ✓ | N/A | 5.89E-16 | 292.87 | 54 |
| 16 | × | ✓ | 3.15E+00 | 27.93 | 1 | 51 | × | N/A | 4.07E-16 | 149.80 | 27 |
| 17 | × | ✓ | 9.30E-01 | 59.20 | 1 | 52 | ✓ | N/A | 1.91E-13 | 17.84 | 2 |
| 18 | × | ✓ | 3.59E+16 | 21.15 | 2 | 53 | ✓ | N/A | 4.67E-12 | 331.74 | 56 |
| 19 | × | ✓ | 1.00E+00 | 103.90 | 33 | 54 | ✓ | N/A | 2.00E-11 | 274.11 | 42 |
| 20 | ✓ | ✓ | 1.64E+00 | 53.88 | 1 | 55 | ✓ | N/A | 4.38E-16 | 471.28 | 75 |
| 21 | ✓ | ✓ | 2.30E+00 | 22.66 | 2 | 56 | ✓ | N/A | 6.76E-12 | 436.90 | 64 |
| 22 | ✓ | ✓ | 1.00E+00 | 21.46 | 2 | 57 | × | N/A | 4.64E-16 | 153.11 | 27 |
| 23 | × | ✓ | 2.93E-01 | 24.99 | 2 | 58 | ✓ | N/A | 3.26E-10 | 266.74 | 31 |
| 24 | ✓ | ✓ | 1.00E+00 | 28.02 | 2 | 59 | × | N/A | 2.55E-07 | 1196.44 | 69 |
| 25 | × | ✓ | 9.93E-02 | 17.23 | 2 | 60 | ✓ | N/A | 5.70E-08 | 198.41 | 24 |
| 26 | ✓ | ✓ | 1.92E-10 | 155.57 | 65 | 61 | ✓ | N/A | 8.72E-09 | 213.73 | 22 |
| 27 | ✓ | ✓ | 1.19E-11 | 8.43 | 2 | 62 | ✓ | N/A | 3.97E-08 | 173.21 | 21 |
| 28 | ✓ | ✓ | 7.94E-12 | 8.18 | 2 | 63 | ✓ | N/A | 2.12E-09 | 178.04 | 20 |
| 29 | ✓ | ✓ | 6.73E-12 | 8.02 | 2 | 64 | ✓ | N/A | 2.59E-08 | 179.97 | 23 |
| 30 | ✓ | × | 2.13E-16 | 15.76 | 2 | 65 | ✓ | N/A | 1.10E-08 | 244.44 | 22 |
| 31 | ✓ | N/A | 3.53E-13 | 349.00 | 59 | 66 | ✓ | N/A | 3.95E-09 | 277.10 | 24 |
| 32 | × | N/A | 1.84E-16 | 254.92 | 44 | 67 | ✓ | N/A | 4.80E-10 | 283.52 | 24 |
| 33 | ✓ | N/A | 1.66E-16 | 281.46 | 53 | 68 | ✓ | N/A | 2.76E-08 | 235.03 | 39 |
| 34 | × | N/A | 7.70E-06 | 275.66 | 40 | 69 | ✓ | N/A | 8.33E-09 | 290.20 | 45 |
| 35 | × | N/A | 6.26E-05 | 241.12 | 30 | 70 | ✓ | N/A | 1.76E-09 | 241.05 | 38 |

its approach along one dimension in the multi-variate cases, its overhead is much higher than the single-variate scenarios. Besides, the number of input ranges found by EIFFEL also impacts its execution overhead, since we currently detect errors in each range sequentially. One can see that $N$ in Table III is at a scale of one to several dozens for those benchmarks whose execution overhead is of hundreds of seconds or more. On average, EIFFEL consumes 160.88 seconds for each benchmark, achieving a 22.38× speedup over that of S³FP.

*3) Difference between ACES:* ACES refers to the atomic-condition-based evolutionary search heuristic used by ATOMU. The advantage of our approach over ACES is straightforward: EIFFEL is able to handle multi-variate expressions but ATOMU cannot. We believe this is due to the use of condition numbers, since their computation has more strict requirements [7], [12] and is often considered more different than evaluating an FP expression itself [11]. ATOMU fails to report errors for *bench* 1, 2, and 30 with no information output.

### B. Quality and Quantity of Inferred Input Ranges

The two baselines used in §IV-A represent the state of the art of error detection tools. However, both S³FP and ATOMU only localize the source of significant errors to specific input values.

As such, we compare the quality and quantity of the input ranges inferred by EIFFEL with two FP rewriting engines– Regina [19] and PSAT [13]. These two approaches either infer or search input ranges like EIFFEL does before rewriting an FP expression.

We let Regina take as inputs the 70 benchmarks considered, for all of which it unfortunately fails to infer input ranges. It seems that the reason is due to the large size of $\mathcal{D}$ we use in our experiments. However, this tool still does not work when $\mathcal{D}$ is reduced to [0, 200]. It only works when setting a domain using very small sizes defined in its publication, which does not make sense for our approach because a densely sampling across such small domains directly outputs ranges or specific values of significant errors.

Similar to ATOMU, PSAT only finds the input ranges for the 30 single-variate benchmarks. It did not terminate within one hour for all of the multi-variate cases. EIFFEL determines that 22 out of the single-variate benchmarks either are monotonous or observe their maximal errors in $\mathcal{R}$. In the former case, EIFFEL always produces a single input range, whose size could be defined as small as possible. In the later case, we do not infer input ranges because we just dynamically measure errors across $\mathcal{R}$. For the remaining eight benchmarks, EIFFEL also infers multiple input ranges, the quality and quantity of which, together with those of PSAT, are reported in Table IV. Generally speaking, EIFFEL obtains more input ranges than PSAT, with average sizes comparable to those of PSAT.

TABLE IV: Comparison of the number $N$ of input ranges inferred by EIFFEL and PSAT, and their average range size.

| no. | N | | average size | | no. | N | | average size | |
|---|---|---|---|---|---|---|---|---|---|
| | EIFFEL | PSAT | EIFFEL | PSAT | | EIFFEL | PSAT | EIFFEL | PSAT |
| 1 | 15 | 2 | 0.956 | 8.382 | 9 | 14 | 2 | 0.966 | 4.651 |
| 2 | 7 | 2 | 0.861 | 0.729 | 12 | 15 | 5 | 0.604 | 0.540 |
| 3 | 12 | 2 | 10.438 | 6.332 | 19 | 27 | 4 | 0.556 | 0.362 |
| 5 | 14 | 2 | 1.177 | 0.582 | 26 | 21 | 3 | 0.839 | 1.314 |

### C. Compatibility with Program Rewriting Engines

The inherit inexactness of curve fitting is one reason why EIFFEL infers input ranges instead of specific values. Another reason is because such input ranges can be taken in by modern program rewriting engines, with each input range treated as a sub-domain or regime where an original FP expression should be rewritten such that the overall accuracy across the entire domain is improved. Regina is such a rewriting engine, but we do not consider it here since it uses static analysis to compute an over-approximated error bound rather than dynamically measuring them as EIFFEL does. PSAT is also a system of this kind that rewrites an FP expression using its higher precision version in regimes. By multiplying the average range size and its range number $N$ in Table IV, one can conclude that EIFFEL always produces a larger regime length than PSAT. Rewriting an FP expression using higher precision in a larger regime thus always obtains better accuracy.

We also feed the inferred input ranges to Herbie v1.6 [6], which rewrites an FP expression in regimes using its equivalent
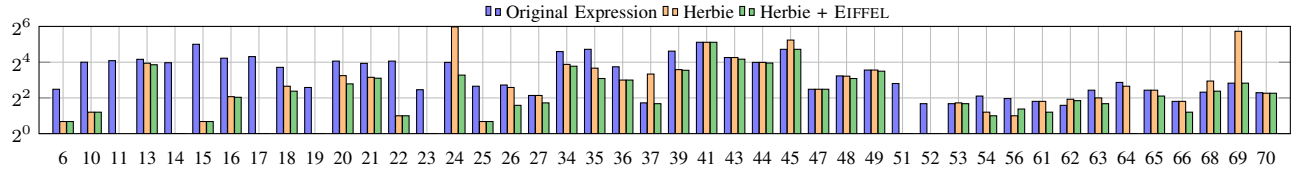
Fig. 15: Comparison of the maximal errors (lower is better). $x$ axis is the benchmark numbers; $y$ axis is log-2 scaled. Herbie, with or without EIFFEL, exhibits 0-bit errors for *bench* 11, 14, 17, 19, 23, 51, and 52. The data are thus not shown.

expression with lower errors. Fig. 15 shows the maximal errors in number of bits of different versions. Herbie samples 8000 values to localize the input values of significant errors, which often fails to produce effective regimes. By rewriting FP expressions, Herbie reduces the maximal errors of most benchmarks but increases them for *bench* 24, 37, 45, 68, and 69. By taking as the input ranges inferred by EIFFEL, Herbie always obtains effective regime strategies, further reducing the maximal errors of these benchmarks. In summary, the Herbie variant with our inferred input ranges obtains an average accuracy improvement of 3.35 bits and up to 53.3 bits (*bench* 24) over the original Herbie version. Moreover, our work helps Herbie achieve a mean accuracy improvement of 5.77 bits and up to 30.5 bits (*bench* 15) over the original expressions.

### D. Comparison of Overhead

The overhead of EIFFEL has been reported in Table III, which includes the execution times of various steps introduced in §III. To help readers better understand how these steps contribute to the overall overhead, we show the time breakdown in Fig. 16. The error detection through dynamically measuring in each inferred input range (§ III-D) consumes most of the execution overhead. A simulated result show that the overhead could be reduced to an average mean of 8.05 seconds (2.00 at least and 50.91 at most) if the error detection in each range is fully parallelized, which we leave as a future task.
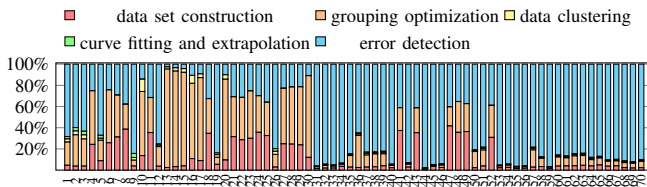


Fig. 16: Breakdown of the execution overhead.

We also compare the overhead of EIFFEL with S³FP and ATOMU, as depicted in Fig. 17. As ATOMU relies on condition numbers to localize input values of significant errors and does not need to compute oracles, its overhead is constant and produces the lower (cyan) horizontal line. On the contrary, S³FP always consumes up its time budget, yield-
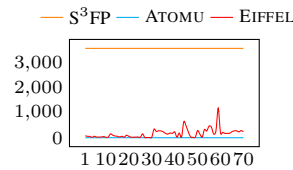


Fig. 17: Comparison of overhead. $y$ axis is the time in seconds.

ing the higher (orange) horizontal line. The (red) curve between these two lines is the overhead of EIFFEL.

## V. RELATED WORK

Many dynamic approaches for FP error detection are modeled as guided searches of error-inducing ranges [8], [10] or values [7], [9]. We evaluate EIFFEL by comparing it with two state-of-the-art tools, with the wider general applicability over their search heuristics discussed. Those not empirically evaluated [8], [10] are also challenged by the issues introduced in §I. By combining static and dynamic analysis, the work of Krämer et al. [26] also infers input ranges with reasonable errors to provide a sound error bound, while EIFFEL only infers input ranges of significant errors.

FP arithmetic is another source of numerical errors [27]. A thread of dynamic approaches seek FP arithmetic causing high errors based on the binary instrumentation tool Valgrind [28]. FpDebug [29] is the pioneer work of this category, which detects errors caused by catastrophic cancellations and locates operations by performing higher-precision FP computations side by side. By targeting on the high overhead issue caused by higher-precision computation in FpDebug, Bao and Zhang [30] propose an approach to search operations of significant errors via tagging and tracing potential inaccurate values, and the overhead of their approach is further reduced by RAIVE [31]. PFPSanitizer [32] and EFTSanitizer [33] also focus on reducing the overhead of FP error detection.

FpDebug is also enhanced to fix errors of operations sensitive to specific precision [34]. Verror [35] and Herbgrind [27] are another two methods based on Valgrind. Verror uses little instrumentation to alleviate the overhead of FpDebug, and Herbgrind determines whether an FP expression incurs high errors when given an error-triggering input. The recent error detector for computation with Posits [36] is an enhancement of Herbgrind. These methods and EIFFEL can all be used for debugging, but EIFFEL works on source code while these approaches extract dynamic information from a binary program.

By leveraging atomic conditions, ATOMU [7] is also able to localize FP operations that contribute most to the final errors. It exhibits more lightweight overhead than EIFFEL by avoiding the need to compute oracles, which, however, can be computed by tools like MPFR [24] or static approaches [37]. Condition numbers also used in [12], which require the estimated function to be differential and are thus considered as more difficult [11]. Like ATOMU, the approach of Fu et al. [12]

also seems to restricted to single-variate functions. A similar oracle-free technique [38] is also used to repair FP programs.

Based on abstract interpretation [39], interval [40] or affine arithmetic [41], and symbolic execution [42], most static analysis methods [4], [5], [43]–[46] try to tighten the worst-case error bounds, which often produce overly pessimistic results. EIFFEL is orthogonal to them in that (1) it is dynamic, (2) the maximal error found is not over-approximated, and (3) it is a combination of research techniques of data analysis and software engineering. There also exist methods [47], [48] that use symbolic execution to find concrete inputs that expose FP exceptions or large errors.

## VI. DISCUSSIONS AND FUTURE WORK

In this paper, we presented EIFFEL to infer input ranges of significant FP errors by combining the research of data analysis and software engineering. We first use *ulp* errors to construct the error data set, based on which the error data points are clustered into clusters. The number these clusters, combined with the monotonicity of an FP expression, allows EIFFEL to detect errors in $\mathcal{R}$ or the inferred ranges, both localizing the error input values of significant errors into smaller intervals. EIFFEL is finally generalized to handle multi-variate expressions. The experimental results demonstrate the general applicability of EIFFEL, the quality and quantity of its inferred input ranges, and its compatibility to rewriting engines. There are several interesting points that worth discussing.

*a) Compatibility with Other Error Measuring Methods:* Measuring errors is used to automatically compute the two parameters of the DBSCAN algorithm. The main reason why Eq (1) and Eq (2) are not used resides in their bad suitability to appropriately compute Euclidean distance between error data points. One can scale up the error axes of Fig. 2a and Fig. 2b to eliminate this weakness, which is exactly what Eq (3) does.

Suppose that such scaling is possible, e.g., multiplying each error value by $10^{16}$ for Fig. 2a and Fig. 2b. Absolute errors still face the limitation described below Eq (1). Since EIFFEL needs to approximate curves to infer unknown ranges, an error distribution plot with steeper slopes is preferred, which benefits both data clustering and curve fitting by reducing the number of data points. Fig. 2c is thus a better choice than Fig. 2b with respect to this selection criterion. As indicated by Eq (2) and Eq (3), a relative error divides the absolute error by $o(f(x))$ but the denominator becomes the ULP($o(f(x))$) for the *ulp* error. While $o(f(x))$ always changes, ULP($o(f(x))$) is only related to the exponent part and can stay unchanged for different $f(x)$ values. This difference converts Fig. 2b into Fig. 2c. Using relative errors would come with increased execution overhead and/or reduced number of clusters produced by the DBSCAN algorithm.

*b) Effectiveness of Data Clustering:* Data clustering is used to extract error features from the error data set. The DBSCAN clustering is a good fit for our work and plays an important role. Combining Eq (4) and Eq (5) into Table I makes it possible to fully automate EIFFEL. Determining the parameters by querying Table I might fail, but we fortunately

have not yet come across with such scenarios. The possible solution to this issue is to develop a better technique to determine *MinPts* and $\epsilon$. Other data analysis algorithms could also be good alternatives, provided they are with the advantages of the DBSCAN algorithm. In particular, it would be better that an alternative algorithm can avoid the need of the grouping optimization (§III-B), which can further simplify the workflow of EIFFEL, because, this optimization sometimes consumes a large portion of the execution overhead in Fig. 16.

*c) Soundness of Curve Function:* Approximating a curve allows us to concentrate on building a simple but effective mathematical function to model the distribution of error-inducing input ranges. Indeed, extrapolating such polynomial curves to infer input ranges is far from a perfect fit for error analysis, but the experimental results demonstrate its effectiveness. Importantly, integrating curve fitting, extrapolation and data clustering uncovers a new direction for this field.

The inexactness of Eq (6) impels us to consider inferring error-inducing input ranges rather error-triggering input values. This approach might still introduce false positives and also false negatives, but the latter case was not experienced. As for the theoretical underpinning of Eq (6), we believe it has a relation with the use of *ulp* errors. Functions built via a more rigorous mathematical derivation, possibly not polynomial, are expected; those that can directly locate error-triggering input values are appreciated.

*d) Limitations and Future Plans:* EIFFEL in its current form has two weaknesses. First, the handling of multi-variate expressions is still simple, making the results of EIFFEL on such benchmarks still improvable. Optimizing EIFFEL to better deal with such cases is our next plan. Second, the theoretical foundation of Eq (6) is still missing. We are now investigating on this issue. In addition, we also intend to parallelize the execution of EIFFEL to further reduce its overhead.

While EIFFEL demonstrates that applying data analysis techniques to FP error detection is practical, our work also opens some interesting research directions to follow in this field. For example, training a neural network to fit an FP error curve in Fig. 9 instead of using polynomial extrapolation is also possible. Extending EIFFEL to combine the research fields of software engineering and machine learning this way should be a promising research problem. Similarly, making use of data analysis techniques to extract and encode the domain-specific knowledge, e.g., Eq (6), should be a better choice than an empirical observation. In particular, this would probably help EIFFEL find more exact FP errors.

## REFERENCES

[1] K. Quinn, "Ever had problems rounding off figures," *This stock exchange has. The Wall Street Journal*, p. 37, 1983.

[2] R. Skeel, "Roundoff error and the patriot missile," *SIAM News*, vol. 25, no. 4, p. 11, 1992.

[3] A. Bessey, K. Block, B. Chelf, A. Chou, B. Fulton, S. Hallem, C. Henri-Gros, A. Kamsky, S. McPeak, and D. Engler, "A few billion lines of code later: Using static analysis to find bugs in the real world," *Commun. ACM*, vol. 53, no. 2, pp. 66–75, feb 2010. [Online]. Available: https://doi.org/10.1145/1646353.1646374

[4] E. Goubault and S. Putot, "Static analysis of finite precision computations," in *Proceedings of the 12th International Conference on Verification, Model Checking, and Abstract Interpretation*, ser. VMCAI'11. Berlin, Heidelberg: Springer-Verlag, 2011, pp. 232–247.

[5] S. Simić, A. Bemporad, O. Inverso, and M. Tribastone, "Tight error analysis in fixed-point arithmetic," *Form. Asp. Comput.*, vol. 34, no. 1, sep 2022. [Online]. Available: https://doi.org/10.1145/3524051

[6] P. Panchekha, A. Sanchez-Stern, J. R. Wilcox, and Z. Tatlock, "Automatically improving accuracy for floating point expressions," in *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI'15. New York, NY, USA: Association for Computing Machinery, 2015, pp. 1–11. [Online]. Available: https://doi.org/10.1145/2737924.2737959

[7] D. Zou, M. Zeng, Y. Xiong, Z. Fu, L. Zhang, and Z. Su, "Detecting floating-point errors via atomic conditions," *Proc. ACM Program. Lang.*, vol. 4, no. POPL, dec 2019. [Online]. Available: https://doi.org/10.1145/3371128

[8] X. Yi, L. Chen, X. Mao, and T. Ji, "Efficient automated repair of high floating-point errors in numerical libraries," *Proc. ACM Program. Lang.*, vol. 3, no. POPL, jan 2019. [Online]. Available: https://doi.org/10.1145/3290369

[9] W.-F. Chiang, G. Gopalakrishnan, Z. Rakamaric, and A. Solovyev, "Efficient search for inputs causing high floating-point errors," in *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPoPP '14. New York, NY, USA: Association for Computing Machinery, 2014, pp. 43–52. [Online]. Available: https://doi.org/10.1145/2555243.2555265

[10] D. Zou, R. Wang, Y. Xiong, L. Zhang, Z. Su, and H. Mei, "A genetic algorithm for detecting significant floating-point inaccuracies," in *Proceedings of the 37th International Conference on Software Engineering - Volume 1*, ser. ICSE'15. IEEE Press, 2015, pp. 529–539.

[11] N. J. Higham, *Accuracy and stability of numerical algorithms*. SIAM, 2002.

[12] Z. Fu, Z. Bai, and Z. Su, "Automated backward error analysis for numerical code," in *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*, ser. OOPSLA 2015. New York, NY, USA: Association for Computing Machinery, 2015, pp. 639–654. [Online]. Available: https://doi.org/10.1145/2814270.2814317

[13] X. Wang, H. Wang, Z. Su, E. Tang, X. Chen, W. Shen, Z. Chen, L. Wang, X. Zhang, and X. Li, "Global optimization of numerical programs via prioritized stochastic algebraic transformations," in *Proceedings of the 41st International Conference on Software Engineering*, ser. ICSE'19. IEEE Press, 2019, pp. 1131–1141. [Online]. Available: https://doi.org/10.1109/ICSE.2019.00116

[14] J.-M. Muller, "On the definition of ulp (x)," INRIA, LIP, Tech. Rep. 5504, 2005.

[15] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, ser. KDD'96. AAAI Press, 1996, pp. 226–231.

[16] N. Damouche, M. Martel, P. Panchekha, C. Qiu, A. Sanchez-Stern, and Z. Tatlock, "Toward a standard benchmark format and suite for floating-point analysis," in *Numerical Software Verification*, S. Bogomolov, M. Martel, and P. Prabhakar, Eds. Cham: Springer International Publishing, 2017, pp. 63–77.

[17] T. Plachetka, "Pov ray: persistence of vision parallel raytracer," in *Proc. of Spring Conf. on Computer Graphics, Budmerice, Slovakia*, 1998, pp. 123–129.

[18] D. Lohar, M. Prokop, and E. Darulova, "Sound probabilistic numerical error analysis," in *Integrated Formal Methods*, W. Ahrendt and S. L. Tapia Tarifa, Eds. Cham: Springer International Publishing, 2019, pp. 322–340.

[19] R. Rabe, A. Izycheva, and E. Darulova, "Regime inference for sound floating-point optimizations," *ACM Trans. Embed. Comput. Syst.*, vol. 20, no. 5s, sep 2021. [Online]. Available: https://doi.org/10.1145/3477012

[20] I. S. Association, "Ieee standard for floating-point arithmetic," *IEEE Std 754-2019 (Revision of IEEE 754-2008)*, pp. 1–84, 2019.

[21] D. Goldberg, "What every computer scientist should know about floating-point arithmetic," *ACM Comput. Surv.*, vol. 23, no. 1, pp. 5–48, mar 1991. [Online]. Available: https://doi.org/10.1145/103162.103163

[22] S. Lloyd, "Least squares quantization in pcm," *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 129–137, 1982.

[23] J. MacQueen, "Classification and analysis of multivariate observations," in *5th Berkeley Symp. Math. Statist. Probability*, 1967, pp. 281–297.

[24] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier, and P. Zimmermann, "Mpfr: A multiple-precision binary floating-point library with correct rounding," *ACM Trans. Math. Softw.*, vol. 33, no. 2, pp. 13–es, jun 2007. [Online]. Available: https://doi.org/10.1145/1236463.1236468

[25] S. Loosemore, R. M. Stallman, R. McGrath, A. Oram, and U. Drepper, "The gnu c library reference manua," 2019. [Online]. Available: https://www.gnu.org/software/libc/manual/2.35/html_node/Errors-in-Math-Functions.html

[26] J. Krämer, L. Blatter, E. Darulova, and M. Ulbrich, "Inferring interval-valued floating-point preconditions," in *Tools and Algorithms for the Construction and Analysis of Systems*, D. Fisman and G. Rosu, Eds. Cham: Springer International Publishing, 2022, pp. 303–321.

[27] A. Sanchez-Stern, P. Panchekha, S. Lerner, and Z. Tatlock, "Finding root causes of floating point error," in *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI 2018. New York, NY, USA: Association for Computing Machinery, 2018, pp. 256–269. [Online]. Available: https://doi.org/10.1145/3192366.3192411

[28] N. Nethercote and J. Seward, "Valgrind: A framework for heavyweight dynamic binary instrumentation," in *Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI '07. New York, NY, USA: Association for Computing Machinery, 2007, pp. 89–100. [Online]. Available: https://doi.org/10.1145/1250734.1250746

[29] F. Benz, A. Hildebrandt, and S. Hack, "A dynamic program analysis to find floating-point accuracy problems," in *Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI '12. New York, NY, USA: Association for Computing Machinery, 2012, pp. 453–462. [Online]. Available: https://doi.org/10.1145/2254064.2254118

[30] T. Bao and X. Zhang, "On-the-fly detection of instability problems in floating-point program execution," in *Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications*, ser. OOPSLA '13. New York, NY, USA: Association for Computing Machinery, 2013, pp. 817–832. [Online]. Available: https://doi.org/10.1145/2509136.2509526

[31] W.-C. Lee, T. Bao, Y. Zheng, X. Zhang, K. Vora, and R. Gupta, "Raive: Runtime assessment of floating-point instability by vectorization," in *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*, ser. OOPSLA 2015. New York, NY, USA: Association for Computing Machinery, 2015, pp. 623–638. [Online]. Available: https://doi.org/10.1145/2814270.2814299

[32] S. Chowdhary and S. Nagarakatte, "Parallel shadow execution to accelerate the debugging of numerical errors," in *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2021. New York, NY, USA: Association for Computing Machinery, 2021, pp. 615–626. [Online]. Available: https://doi.org/10.1145/3468264.3468585

[33] S. Chowdhary and S. Nagarakatte, "Fast shadow execution for debugging numerical errors using error free transformations," *Proc. ACM Program. Lang.*, vol. 6, no. OOPSLA2, oct 2022. [Online]. Available: https://doi.org/10.1145/3563353

[34] R. Wang, D. Zou, X. He, Y. Xiong, L. Zhang, and G. Huang, "Detecting and fixing precision-specific operations for measuring floating-point errors," in *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, ser. FSE 2016. New York, NY, USA: Association for Computing Machinery, 2016, pp. 619–630. [Online]. Available: https://doi.org/10.1145/2950290.2950355

[35] F. Févotte and B. Lathuiliere, "Verrou: a cestac evaluation without recompilation," *SCAN 2016*, p. 47, 2016.

[36] S. Chowdhary, J. P. Lim, and S. Nagarakatte, "Debugging and detecting numerical errors in computation with posits," in *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and*

*Implementation*, ser. PLDI 2020.   New York, NY, USA: Association for Computing Machinery, 2020, pp. 731–746. [Online]. Available: https://doi.org/10.1145/3385412.3386004

[37] J. Chen, Y. Bai, D. Hao, L. Zhang, L. Zhang, B. Xie, and H. Mei, "Supporting oracle construction via static analysis," in *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE '16.   New York, NY, USA: Association for Computing Machinery, 2016, pp. 178–189. [Online]. Available: https://doi.org/10.1145/2970276.2970366

[38] D. Zou, Y. Gu, Y. Shi, M. Wang, Y. Xiong, and Z. Su, "Oracle-free repair synthesis for floating-point programs," *Proc. ACM Program. Lang.*, vol. 6, no. OOPSLA2, oct 2022. [Online]. Available: https://doi.org/10.1145/3563322

[39] P. Cousot and R. Cousot, "Abstract interpretation: Past, present and future," in *Proceedings of the Joint Meeting of the Twenty-Third EACSL Annual Conference on Computer Science Logic (CSL) and the Twenty-Ninth Annual ACM/IEEE Symposium on Logic in Computer Science (LICS)*, ser. CSL-LICS '14.   New York, NY, USA: Association for Computing Machinery, 2014. [Online]. Available: https://doi.org/10.1145/2603088.2603165

[40] T. Hickey, Q. Ju, and M. H. Van Emden, "Interval arithmetic: From principles to implementation," *J. ACM*, vol. 48, no. 5, pp. 1038–1068, sep 2001. [Online]. Available: https://doi.org/10.1145/502102.502106

[41] L. H. de Figueiredo and J. Stolfi, "Affine arithmetic: Concepts and applications," *Numerical Algorithms*, vol. 37, no. 1, pp. 147–158, Dec 2004. [Online]. Available: https://doi.org/10.1023/B: NUMA.0000049462.70970.b6

[42] R. Baldoni, E. Coppa, D. C. D'elia, C. Demetrescu, and I. Finocchi, "A survey of symbolic execution techniques," *ACM Comput. Surv.*, vol. 51, no. 3, may 2018. [Online]. Available: https://doi.org/10.1145/3182657

[43] A. Izycheva and E. Darulova, "On sound relative error bounds for floating-point arithmetic," in *2017 Formal Methods in Computer Aided Design (FMCAD)*, 2017, pp. 15–22.

[44] W. Lee, R. Sharma, and A. Aiken, "Verifying bit-manipulations of floating-point," in *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI '16. New York, NY, USA: Association for Computing Machinery, 2016, pp. 70–84. [Online]. Available: https://doi.org/10.1145/2908080.2908107

[45] A. Solovyev, M. S. Baranowski, I. Briggs, C. Jacobsen, Z. Rakamarić, and G. Gopalakrishnan, "Rigorous estimation of floating-point round-off errors with symbolic taylor expansions," *ACM Trans. Program. Lang. Syst.*, vol. 41, no. 1, dec 2018. [Online]. Available: https://doi.org/10.1145/3230733

[46] Y. Zhang, L. Ren, L. Chen, Y. Xiong, S.-C. Cheung, and T. Xie, "Detecting numerical bugs in neural network architectures," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2020.   New York, NY, USA: Association for Computing Machinery, 2020, pp. 826–837. [Online]. Available: https://doi.org/10.1145/3368089.3409720

[47] E. T. Barr, T. Vo, V. Le, and Z. Su, "Automatic detection of floating-point exceptions," in *Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, ser. POPL'13. New York, NY, USA: Association for Computing Machinery, 2013, pp. 549–560. [Online]. Available: https://doi.org/10.1145/2429069.2429133

[48] H. Guo and C. Rubio-González, "Efficient generation of error-inducing floating-point inputs via symbolic execution," in *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*, ser. ICSE'20.   New York, NY, USA: Association for Computing Machinery, 2020, pp. 1261–1272. [Online]. Available: https://doi.org/10.1145/3377811.3380359