

A Software Tool for Designing Fixed-Point Implementations of Computational Data Paths

David M. Buehler and Gregory W. Donohoe
University of Idaho, Moscow, ID 83843

This research aims to investigate the feasibility of creating a fixed-point implementation of a computation based on static analysis techniques. An underlying theory for fixed-point computations has been developed, and a software tool written to implement much of that theory.

This paper presents a brief overview of the theory, which is based on determining the semantics of the individual bits of each fixed-point value in the computation.

Following the theoretical overview, some details of the software tool are provided, emphasizing both the controls available to an algorithm implementor to guide the creation of the fixed-point implementation of the computation and the heuristic choices which are embedded in the tool's operation. The software tool provides several options for what output is generated, these options are described briefly.

Finally, we discuss the quality of fixed-point implementations our tool generates for several interesting computations and describe the challenge of statically estimating truncation error. We provide a quantitative comparison between the software tool's estimates of truncation error and actual truncation error observed.

I. Introduction

IN computational environments where the cost of floating-point circuitry is prohibitive, such as deeply embedded computing and reconfigurable computing, computations requiring values from the mathematical field of reals can often be performed using fixed-point representations. However, designing fixed-point implementations of computational algorithms has a history of being a difficult, time-consuming and error-prone task which can result in sub-optimal implementations due to unnecessarily high computational error or even incorrect results.

In conjunction with the Field Programmable Processor Array (FPPA) project,¹ the software tool `SIFOpt` (Sign, Integer, Fraction-based Optimizer) has been developed which helps algorithm designers create fixed-point implementations of computational algorithms. `SIFOpt` reads an algebraic description of the computation. The user must annotate this algebraic description with information about the fixed-point encoding of the computation's arguments. A static analysis of the computation is then performed and a fixed-point implementation for the computation is determined. Values which a fixed-point implementation designer would need to determine by hand (in the absence of a design tool) are computed by `SIFOpt`, including: scaling factors for computed variables, alignment operations for additions and subtractions, rescaling operations for multiplications and determining optimal representations for the constants which are used in the computation.

A. Fixed-Point Implementations of Computational Algorithms

A fixed-point implementation of a computational algorithm has something of a dual nature. At run time, integer computations are performed on integer values. However, at design time a scaling factor is associated with each value

Received 01 November 2004; revision received 28 March 2005; accepted for publication 04 April 2005. Copyright © 2005 by the American Institute of Aeronautics and Astronautics, Inc. All rights reserved. Copies of this paper may be made for personal or internal use, on condition that the copier pay the \$10.00 per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923; include the code 1542-9423/04 \$10.00 in correspondence with the CCC.

Table 1 Effect of changing the scaling factor associated with a signed 16-bit value.

Scaling factor:	2^{-9}	2^{-8}	2^{-7}
Maximum:	01111111.11111111 (63.998)	01111111.11111111 (127.996)	01111111.11111111 (255.992)
Minimum:	10000000.00000000 (-64)	10000000.00000000 (-128)	10000000.00000000 (-256)
Delta:	00000000.00000001 (0.00195)	00000000.00000001 (0.00391)	00000000.00000001 (0.00781)

that will be computed at run time. Thus the integer run time values represent real values, which are determined by multiplying a run time integer value with the design time scaling factor associated with it.

The choice of the scaling factor for each computed value is left for the algorithm implementor to determine. In a fixed word length computational environment, the choice of a particular scaling factor establishes both the range of real values that can be represented (the product of the run time integer and the scaling factor) and the granularity of the values that are represented. As the scaling factor is increased, value range increases and granularity decreases. As the scaling factor is decreased, value range decreases and granularity increases. Table 1 illustrates this effect for changes in the scaling factor associated with a signed 16-bit value.

At run time, if the result of a computation is outside the value range established by the choice of scaling factors, an overflow condition exists. The consequences of the overflow condition range from loss of accuracy in the result to incorrect results. On the other hand, if a non-zero run time value falls below the granularity of the computed value (which is established by the choice of scaling factors) and becomes zero, we say that an underflow condition exists. Underflow conditions tend to result in loss of granularity in a result (loss of accuracy), but can also lead to incorrect results.

The primary challenge for the algorithm designer is to assign scaling factors to each computed value which are as small as possible yet not so small that they result in overflow errors at run time. In addition, the choice of scaling factors affects the operations required to align values for addition operations, and the operations required to prescale values for multiplications (in computational environments where the result of a multiplication are not as wide as the sum of the widths of the multiplicands.)

B. Related Research – Dynamic Approaches

Recent research into creating fixed-point implementations of computations have focused on dynamic approaches. With a dynamic approach, the computation is implemented in a floating-point environment and the computed values are “instrumented” to collect statistics about their run time values.

Sample data sets are then run through the instrumented versions of the computation and run time statistics are gathered. The run time statistics are then used to assign scaling factors to each computed value. This step acts as a training step which is used to “learn” the scaling factors to use for each computed value.

Researchers from Seoul National University,² The University of Toronto³ and Aachen University of Technology⁴ use this approach to perform automatic floating-point to fixed-point conversion, taking a C program which includes floating-point variables and constants and generating a C program which uses only integer variables and constants.

Researchers at the University of Washington⁵ and MIT⁶ use this approach to minimize the data path width of FPGA and ASIC implementations of fixed-point computations.

1. Weaknesses of Dynamic Approaches

We see several weaknesses of this type of approach. One is the obvious dependence of the results on the training set used to determine the fixed-point implementation. Of more concern, however, is that the scaling factors are determined based on floating-point values which are not subjected to fixed-point truncation errors. This has the potential to lead to less than optimal fixed-point implementations.

The worst example of floating-point to fixed-point mismatch we have encountered was in analyzing our “division by repeated multiplications” algorithm⁷ using floating-point methods. For a 16-bit by 16-bit division, using 16 iterations of the algorithm’s loop, the floating-point analysis indicates that the worst-case absolute error will occur

when computing $\frac{65535}{1}$. The floating-point analysis indicates that the algorithm will compute the value 41426 (absolute error of 24109) for the answer, while our fixed-point implementation computes a value of 65519 (absolute error of 16) for the answer. The difference between these values was found to be due to the truncation error of the fixed-point implementation. This is an unusual case, but it illustrates the mismatch which can occur between a floating-point analysis of an algorithm and a fixed-point implementation.

II. The Static Approach

The static approach we use is based upon a design-time notation which tracks the partitioning of the run-time integer values into Sign, Integer and Fraction (SIF) regions. Taken together, the Integer and Fraction regions of a run time value are similar to the mantissa region of a floating-point value representation, so they (the Integer and Fraction regions together) are referred to as the mantissa region.

The user is required to provide SIF partitionings for any values which are arguments to the computation. Computational arguments lie at the bottom of the parse tree for the computation’s description. Note that any values which are “recirculated” (i.e. a value computed by the computation which is also used to compute that value) must also be provided with SIF partitioning values. The user can optionally specify SIF partitionings for values that are computed (internal values in the computation.) The primary task of the SIFOPT tool is to determine SIF partitionings for the result values of every mathematical operation performed by the computation.

A. SIF Partitionings

An SIF partitioning is a notation which specifies the semantics of the individual bits of run time integer values. At run time, the bits of an integer value can act as replicated sign bits, integer-value bits, fractional-value bits or unused right-padding bits.

This notation originated with researchers working at Kansas State University and Sandia National Laboratories in the late 1970s under the name Block Floating Point notation (BFP).⁸ Members of our group were exposed to the notation at Sandia National Laboratories and suggested it as a starting point for a fixed-point design tool.

Though the notation itself remains the same, the theory underlying the original work has been completely reworked and reintroduced with the name “SIF partitionings” in.⁹

The syntax used for SIF partitionings is described by the following specification. Characters in single quotes are literal characters, characters inside square brackets are optional, and a horizontal bar separates list entries of an exclusive-or choice between a list of characters.

$$\text{'(['+'|'-']v_S'/'v_I'/'v_F')' ['^'v_n]}$$

- The optional ‘+’ or ‘-’ character before the value of v_S indicates if the run time value is known to be positive (in the case of a ‘+’) or negative (in the case of a ‘-’).
- v_S is a non-negative integer indicating how many of the run time value’s bits (beginning from the most-significant bit) are replicated sign bits.
- v_I is a non-negative integer indicating how many of the bits following the replicated sign bits represent integral information.
- v_F is a non-negative integer indicating how many of the bits following the integral bits represent fractional information.
- v_n is an optional integer value parameter which indicates an additional shift of the binary point by v_n positions. (If $v_n < 0$, the binary point is shifted to the left, otherwise it is shifted to the right.)

In the remainder of the body of this paper (not in section headings) we shall use the term “SIF” to mean an “SIF partitioning” as defined above.

Given an SIF and the word length of the value it is associated with, we can determine both: how many of the integer’s bits are unused, right-padding bits, and the scaling factor associated with the value.

B. Operations on SIF Partitionings

We view SIFs as a data type that the SIFOPT tool manipulates to produce a fixed-point implementation of the requested computation.

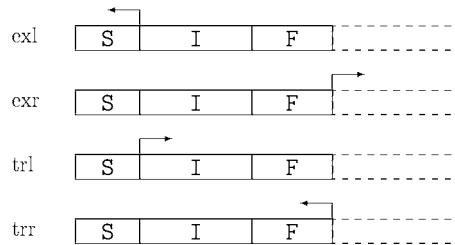


Fig. 1 Operations which move the boundary of the mantissa region.

Given any data type, it is useful to identify the operations which can be performed on values of that type. The fundamental operations which can be performed on SIFs have been identified and discussed thoroughly in,⁹ an overview is provided here.

The SIF_{OPT} tool provides the algorithm designer with the ability to insert design knowledge into the process of determining SIFs for computed values by annotating the description of the computation with functions which represent performing each of the SIF operations described in the following sections.

1. Moving the Boundaries of the Mantissa Region

There are four operations on SIFs which adjust a boundary of the mantissa region. These operations are not usually required for implementing a computation, and should only be used in “extraordinary circumstances” to represent designer knowledge of non-computational behaviors. An illustration of each of these operations can be found in figure 1.

Expand Left (exl) – moves the boundary between the sign bits and the mantissa to the left, causing some bits which were sign bits to be interpreted as mantissa bits.

Expand Right (exr) – moves the boundary between the mantissa region and the unused (padding) region to the right, causing some bits which were interpreted as unused padding bits to be interpreted as mantissa bits.

Truncate Left (trl) – moves the boundary between the mantissa region and the sign region to the right, causing some bits which were interpreted as mantissa bits to be interpreted as sign bits.

Truncate Right (trr) – moves the boundary between the mantissa region and the unused (padding) region to the left, causing some bits which were interpreted as mantissa bits to be interpreted as padding bits.

2. Moving the Binary Point

Moving the boundary between the integer and fractional regions has the effect of multiplying and dividing the real value represented by the run time integer by powers of two, with no run-time operation. Note that the binary point can lie outside the mantissa region, and even beyond the edge of the machine word.

Shift Binary Point Right (sbpr) – shift the location of the binary point to the right. The value mapped to by the integer and scaling factor is divided by 2^n .

Shift Binary Point Left (sbpl) – shift the location of the binary point to the left. The value mapped to by the integer and scaling factor is multiplied by 2^n .

3. Run time Shift with SIF Partitioning Adjustment

The SIF operations presented above are all design time boundary manipulations which do not cause any run time operation to be performed. At run time, the “shift” operations play an important role in fixed-point implementations.

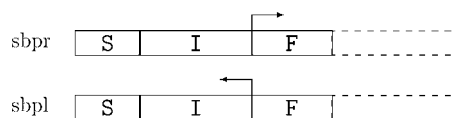


Fig. 2 Operations which move the binary point location.

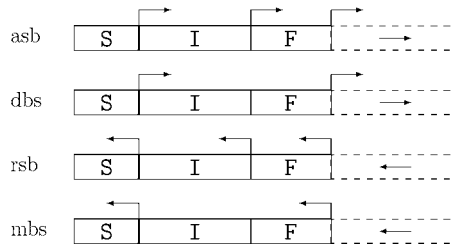


Fig. 3 Operations corresponding to run-time shift operations.

Any run time shift of a value must be accompanied by an adjustment of the value's SIF, or the real value mapped to by the run time integer and the non-adjusted SIF will not correspond to the computation being performed.

When the integer value is *right* shifted at run time, we can interpret the result as either adding sign bits (binary point shifts with integer), or dividing by a power of two (binary point retains its previous location.) When the integer value is shifted *left* at run time, we can interpret the result as either removing sign bits (binary point shifts with the integer) or as multiplying by a power of two (binary point retains its original location.)

Add Sign Bits (asb) – All boundaries are shifted the same amount as the run time shift. The real value represented by the integer stays the same (up to truncation error.)

Divide By Shifting (dbs) – The location of the binary point is held constant while the mantissa boundaries are shifted along with the run time value. The real value represented by the integer is divided by 2^n .

Remove Sign Bits (rsb) – All boundaries are shifted along with the run time value. The real value represented by the integer stays the same (up to overflow error.)

Multiply By Shifting (mbs) – The location of the binary point is held constant while the mantissa boundaries are shifted along with the run time value. The real value represented by the integer is multiplied by 2^n .

III. SIFOpt

The SIFOpt software tool creates fixed-point implementations based on a high-level description of the computation that is to be performed. This description must have annotations specifying the fixed-point formats of the computation's inputs. In addition, the designer can annotate the high-level description with information that is specific to the run time data on which the computation will be performed.

SIFOpt performs a static analysis of the described computation, extracting a computation tree and determining the following:

- Scaling factors (in the form of SIFs), (integer) value ranges and the estimated maximum absolute error (truncation error) for each computed value
- Alignment operations required for addition and subtraction operations
- Rescaling operations required for multiplications
- Optimal integer equivalent values for constants

SIFOpt can produce several types of output. The default output is to print the computation tree annotated with the computed fixed-point information. SIFOpt can also produce integer-only C code. Finally, SIFOpt can produce code which will use a C++ mixed-point class which has been created as part of this work (and which will be described later in this chapter.)

IV. Description of a Numerical Computation

An algebraic language has been implemented in which the algorithm designer specifies the numerical computation to be performed. This language supports two kinds of named values: variables and constants. Named values must be declared before (or when) they are used. Variables which are input values of the computation must be annotated with SIFs, and can optionally be annotated with value range information. Variables representing computed values in the computation can be annotated with SIFs and/or value ranges as well. This allows the algorithm designer to specify information about run time dependencies which cannot be determined by the static analysis.

```
var (0/3/13) [0:6.2832] x; // Input range is 0 to 2*pi
var s = ((((-0.005643 * x + 0.08863) * x + -0.4016) * x + 0.2862) * x +
0.8612) * x + 0.01593;
```

Fig. 4 Example SIFOpt input file for polynomial evaluation. This version puts the whole computation on one line, mimicking the form of the mathematical equation.

```
var (0/3/13) [0:6.2832] x; // Input range is 0 to 2*pi
var p1 = 0.08863 + -0.005643 * x;
var p2 = -0.4016 + x*p1;
var p3 = 0.2862 + x*p2;
var p4 = 0.8612 + x*p3;
var s = 0.01593 + x*p4;
```

Fig. 5 Example SIFOpt input file for polynomial evaluation. This version breaks the commands into simple multiply/accumulate-like steps.

As a simple example, if we wish to compute the value of the polynomial:

$$-0.005643x^5 + 0.08863x^4 + -0.4016x^3 + 0.2862x^2 + 0.8612x + 0.01593$$

we might first apply Horner's scheme to get:

$$(((((-0.005643x) + 0.08863)x + -0.4016)x + 0.2862)x + 0.8612)x + 0.01593$$

Assuming the argument is an unsigned 16-bit value with a scaling factor of 2^{-13} ranging over the real value 0 to 2π , figures 4 through 6 present three possible SIFOpt input files which will yield fixed-point implementations of the calculation.

A. Dependency Checks

The SIF, value range and word length values entered by the user are checked for violations of known dependencies between them. The following conditions cause warnings to be printed, but allow optimization to proceed.

- The user specified a range minimum that is 0 or greater, and an SIF with S-Sign \neq '+'.
- The SIF specified by the user is suboptimal for the value range specified by the user. This is the case when there is an SIF with fewer mantissa bits than the provided SIF which can represent the entire specified value range. In this case, analysis of the algorithm continues using the suboptimal range and SIF.

The following conditions cause an error message to be printed, and SIFOpt exits after parsing the rest of the input file:

- The user specified an SIF with S-Sign = '+' and a value range with a minimum value less than zero.
- The user specified an SIF and a value range which are incompatible. This is the case if the maximum value range representable by the SIF is a (strict) subrange of the provided value range.

```
var (0/3/13) [0:6.2832] x; // Input range is 0 to 2*pi

known c1 = -0.005643;
known c2 = 0.08863;
known c3 = -0.4016;
known c4 = 0.2862;
known c5 = 0.8612;
known c6 = 0.01593;

var s = (((c1 * x + c2) * x + c3) * x + c4) * x + c5) * x + c6;
```

Fig. 6 Example SIFOpt input file for polynomial evaluation. This version uses the ability to associate symbolic names with real values, simplifying the representation of the calculation.

- The user specified an *SIF* with more bits than the word length specified, or more than the default word length, if none is specified.

All warning and error messages inform the user of one possible correction for the detected violation. When possible, the suggested correction preserves the scaling factor of the user-provided *SIF* value.

B. Expressions

Expressions available to the algorithm designer are similar to algebraic expressions available in most high-level languages.

1. Binary Mathematical Operators

`SIFOpt` provides the three binary mathematical operators: ‘+’, ‘-’ and ‘*’, which provide addition, subtraction and multiplication, respectively. The ‘*’ operator has higher precedence than the ‘+’ and ‘-’ operators. Arguments of a string of ‘+’ and ‘-’ operators are grouped left-to-right. Parenthesis (‘(’, ‘)’) can be used to group binary operations for clarity or to override precedence. Division will be added in the near future.

2. Unary Mathematical Operator

The ‘-’ unary negation mathematical operator is also provided, and has higher precedence than multiplication. At this moment it is only available for literal real values.

3. Word Length Specification Operator

The word length specification operator “@*n*” is provided as a postfix operator, and has higher precedence than unary negation.

4. Built-in Functions

All of the unary *SIF* operations introduced in the earlier section “Operations on *SIF* partitionings” (`exl`, `exr`, `trl`, `trr`, `sbpl`, `sbpr`, `asb`, `rsb`, `dbb` and `mbs`) are available in the computation description language.

In addition, the functions `C(expr)` and `NC(expr)` are provided to specify that the outermost binary operation of the enclosed expression will result in a carry-out or will not result in a carry out, respectively.

V. Static Analysis

`SIFOpt` creates a fixed-point implementation of a computational algorithm by performing a static analysis of the computation, which is specified by a `SIFOpt` input file. An “optimization tree” (which mirrors the input file’s parse tree) is constructed by performing a post-order traversal of each statement in the parse tree. All values in the algorithm description must have *SIF*s associated with them before they are used in an expression. This limits the type of computations which can be input to `SIFOpt` to computations which have parse trees that can be expressed as directed acyclic graphs (DAGs).

A. Single Assignment

Single assignment¹ languages¹⁰ meet the requirement that the input algorithm be expressible as a DAG. The input language for `SIFOpt` now restricts the types of algorithms specified to single assignment algorithms.

B. Propagation of *SIF* Partitionings and Value Ranges

The most fundamental task that `SIFOpt` performs is to compute *SIF*s and value ranges for the results of individual mathematical operations. These result values are computed from the *SIF*s and value ranges of the arguments of each mathematical operation by performing a post-order traversal of the computation tree.

Value ranges are propagated via Interval Arithmetic¹¹ by default. An option is available to cause `SIFOpt` to exhaustively compute value ranges for nodes of the computation tree which have common variables in the two child trees (in these cases the range determined using Interval Arithmetic can be sub-optimal.) Exhaustive computation is

¹As the name implies, single-assignment languages allow variables to be assigned to only once.

not in general computationally feasible, but we have found it useful for the size of the problems we are currently addressing.

Using these methods of value range propagation, computed values will not overflow if the input values stay within the bounds specified by the designer. However, computed values may be subject to unnecessarily reduced precision. Reduced precision can be caused by correlations between the run time input values, or dependencies between values in the computation tree (if interval arithmetic is used to compute value ranges.) These problems can be alleviated by application of the $C()$ and $NC()$ functions or specification of value ranges for result values, but use of these language features make run time overflow errors possible.

1. Placement of the Result within the Word Length

In cases where there is flexibility in the placement of the result value within the result word, `SIFOpt` aligns operands using the following heuristics.

For addition (and subtraction) the heuristics used are:

1. If there is only one alignment which minimizes truncation error in the result, that alignment is used.
2. Otherwise, if the operands can be aligned by shifting only one of the operands, only that operand is shifted.
3. Otherwise (both operands have to be shifted to align the binary point and there are multiple ways this can be achieved) the number of sign bits in the result are maximized.

Stated simply: the first priority is minimizing the number of shift operations which must occur. If both arguments must be shifted, the priority is to maximize the number of sign bits in the result. Maximizing (as opposed to minimizing) the number of sign bits in the result is a second priority (to minimizing the number of shifts) because any unexpected overflow which might occur will not be immediately harmful if there are additional sign bits into which the result value can overflow.

For multiplication, the heuristics used are:

1. If there is only one prescaling which minimizes truncation error in the result, that prescaling is used.
2. Otherwise, if the values are already appropriately scaled, then no scaling operations are applied.
3. Otherwise if the values can be prescaled by shifting only one of the multiplicands, then that multiplicand is prescaled to be as small as possible while minimizing the number of truncated mantissa bits.
4. Otherwise both multiplicands are prescaled to be as small as possible while minimizing the number of truncated mantissa bits.

The first priority of the heuristic for multiplication is to minimize the number of shift operations required for prescaling. The second priority is to maximize the number of sign bits in the result.

For addition, the static analysis stores information about the alignment operation implemented and how much flexibility there was in the decision about how to align the addends in the addition tree node. This allows the optimizer to re-visit the node and request a change in the alignment operation(s), in an attempt to minimize the number of shift operations performed *overall* by the implementation. However, this feature is currently “stubbed out” in the `SIFOpt` code and has no effect.

For multiplication, the first case may occur in such a way that we can preserve one more bit in one multiplicand than in the other. A careful analysis of the accumulated truncation amount values for each multiplicand might indicate a preference for truncating the extra bit from one of the multiplicands, but we felt that this analysis might be more misleading than helpful. Instead, we have made the arbitrary choice to preserve an extra bit in the left hand argument of the multiplication.

C. Constants

Constant values in the computation can be optimized quite effectively. Rounding is performed on all constants.

All constant values, whether literal values (e.g. “3.14”) or named constants (knowns) can be followed in a `SIFOpt` algorithm specification by an `SIF` that specifies the format to be used. In particular, the number of sign bits to be used can be specified. If the user specifies an `SIF` for a constant value, and the requested value for the number of sign bits is smaller than the number `SIFOpt` determines using the heuristics to be described below, then `SIFOpt` ignores the `SIF` and optimizes the constant independently of the following heuristics. Therefore, when used with constant values, providing a v_S value of zero tells `SIFOpt` to compute an optimal `SIF` for the constant.

1. *Optimal SIF Partitioning for a Constant*

Given an *SIF* and a word length, we define their optimal *SIF* to have the minimum number of sign bits possible (1 if the constant is negative, 0 otherwise), and as many unused bits as there are 0's to the right of the right-most 1 in the constant's (rounded) integer representation within the given word length.

2. *Optimizing Constant Values for Addition and Subtraction*

Ideally, when we wish to compute the sum of a constant and a variable we will not be required to shift the variable before performing the addition. When this is possible, the value computed for the constant will have the same scaling factor as the variable's scaling factor, and will use only as many bits of resolution as are available in the machine word.

Using this heuristic it is possible that the desired constant will fall outside both the range of significant bits for the variable and the size of the word length, and thus become 0. In this case a warning is printed. In cases such as this, the philosophy followed with *SIFOpt* is that the designer should be alerted to the fact that they've added a value having only n fractional bits with a value whose most significant bit lies more than n places to the right of the binary point. However, a means is provided for the designer to override the way *SIFOpt* normally works. As noted earlier in the section "Constants", an *SIF* can follow any constant value, and setting the number of sign bits to 0 in that *SIF* tells *SIFOpt* to compute the optimal number of sign bits on its own. Providing an *SIF* for the constant is the only way to cause a variable to be left shifted (remove sign bits) when adding with a constant value.

There are two cases which can cause sign bits to be added to the variable before summing with a constant. First is the case in which sign bits must be added for overflow protection. Second is the case in which the constant value extends to the left of the variable and the variable's sign bit region. This case occurs when equation 1 holds. In both cases, the variable is right-shifted as little as possible. The decision to shift as little as possible made it easy to get the maximum resolution in the fixed-point representation of the constant.

$$\log_2(\text{var}) < \log_2(\text{const}) \quad (1)$$

3. *Optimizing Constant Values for Multiplication*

Multiplication by constant values is implemented by first computing an optimal *SIF* for the constant value (up to 32 bits) then using the heuristics stated previously to determine any required rescaling operations. Recall however, that an arbitrary decision was made that when the number of mantissa bits allowed in the multiplicands is odd, an extra bit would be removed from the right-hand argument of the multiplication operator.

When one of the arguments is a constant we know the values of all of the bits in its integer representation, therefore we can optimize the choice of which argument has extra mantissa bit preserved by checking the actual value of the bit which might get truncated from the constant and truncating it if it is a zero. On the other hand, if the bit is a one in the constant, *SIFOpt* chooses to preserve it on the philosophy that it is better to preserve a known 1 bit then to preserve a bit which may or may not be a 1.

This leads to the second optimization. If we detected that the constant's least significant preserved bit is a zero, we should check the next-least significant bit, and if it is also a zero then preserve an additional bit from the variable, and so on.

When no rescaling is required, the constant will be implemented so that the result has the maximum number of sign bits possible.

D. *Estimating Truncation Error*

The static analysis performs one additional function: estimation of truncation error. With every right shift, the maximum truncation error is assumed to occur (i.e. it is assumed that all truncated mantissa bits held values of 1.) The truncated bits are scaled by the scaling factor, so truncation error is computed in terms of the real values being computed. One further assumption must be made: for multiplications, the truncation estimate assumes that each multiplicand is the maximum of its value range. This results in the largest magnification of the accumulated truncation errors, leading to pessimism in the computed results.

VI. SIFOpt Output

SIFOpt can generate three types of output. The first type of output is comprised of information about the fixed-point implementation that has been computed. The second type of output is integer-only C code. The third type is C++ code which uses our “mixed-point” C++ class.

A. Fixed-Point Implementation Information

The fixed-point implementation information output from SIFOpt is composed of two sections.

The first section displays the optimization tree created by SIFOpt. Every node in the computation is displayed on one or more lines, indented in tree hierarchy. Variable declarations display the SIF and real value range that was either declared or determined for the variable. Known declarations display their real value, their optimal integral value and optimal SIF. Expressions display information about the computation being performed, including SIF, value range, absolute error, information about operator alignment and prescaling operations implemented, clipping and padding information.

A short extract of the optimization tree output from SIFOpt follows. This extract covers the statement `out1 = w2 * in1 + out0;` from a convolution implementation. This is not verbatim output from SIFOpt, some cleanup has been performed by hand. The values placed between balanced ‘\’ and ‘/’ characters are the estimated maximum truncation error amounts.

```
Assignment
  Variable - out1: (1/0/31)^-1, [-0.499985:0.499985] \1.52583e-06/
=
  Operation: AddSub - k = 0, kmin = 0, kmax = 0
  BinaryOp: '+' result is: (1/0/31)^-1, [-0.499985:0.499985] \1.52583e-06/
    Operation: Mult
    BinaryOp: '*' result is: (1/0/31)^-1, [-0.349988:0.349988] \1.52583e-06/
      Known - w2 = 0.35 as 45875 = 0xb333 in 32 bits with (+16/0/16)^-1 =
        0.349998 \1.52588e-06/
      ShiftResize - m_delta = 0, m_shift = 16 (17/0/15),
        [-0.999969:0.999969]
        Variable - in1: (1/0/15), [-0.999969:0.999969]
      Variable - out0: (2/0/30)^-2, [-0.149997:0.149997]
```

After the optimization tree, SIFOpt prints a list of the variables in the computation with: the SIFs, the range of real values that can be represented, the range of values of the underlying integer (optional) and maximum estimated absolute error values (if non-zero). The output from the convolution example is shown below.

```
in0 - (1/0/15) [-2147418112:2147418112]->[-0.999969:0.999969]
in1 - (1/0/15) [-2147418112:2147418112]->[-0.999969:0.999969]
in2 - (1/0/15) [-2147418112:2147418112]->[-0.999969:0.999969]
in3 - (1/0/15) [-2147418112:2147418112]->[-0.999969:0.999969]
out0 - (1/0/30)^-2 [-1288463974:1288463974]->[-0.149997:0.149997] \1.52583e-06/
out1 - (1/0/31)^-1 [-2147418112:2147418112]->[-0.499985:0.499985] \3.05166e-06/
out2 - (1/0/31) [-1825302119:1825302118]->[-0.849973:0.849973] \4.57796e-06/
out3 - (1/0/31) [-2147418113:2147418111]->[-0.999969:0.999969] \6.10403e-06/
```

B. Integer-only C Code

The command-line switch “-C” can be used to tell SIFOpt to generate integer-only C code. The output from the convolution example is given below.

```
int in0;
int in1;
int in2;
int in3;
int out0 = (39322 * ((in0) >> 16));
int out1 = ((45875 * ((in1) >> 16)) + ((out0) >> 1));
int out2 = (((45875 * ((in2) >> 16))) >> 1) + ((out1) >> 1));
int out3 = (((39322 * ((in3) >> 16))) >> 2) + out2);
```

C. Mixed-Point C++ Code

The command-line switch “-M” can be used to tell `SIFOpt` to generate output that targets my “mixed-point” C++ class. The output has been slightly cleaned up to remove some redundant parenthesis.

```
mixedpoint in0;
mixedpoint in1;
mixedpoint in2;
mixedpoint in3;
mixedpoint out0 = (mixedpoint::constant( 0x999a, -18, 0.15, 0) * (in0 >> 16));
mixedpoint out1 =
    (mixedpoint::constant( 0xb333, -17, 0.35, 0) * (in1 >> 16)) + (out0 >> 1);
mixedpoint out2 =
    ((mixedpoint::constant( 0xb333, -17, 0.35, 0) * (in2 >> 16)) >> 1) + (out1 >> 1);
mixedpoint out3 =
    ((mixedpoint::constant( 0x999a, -18, 0.15, 0) * (in3 >> 16)) >> 2) + out2;
```

VII. Mixed-Point C++ Class

The mixed-point class is intended to be used to compare a fixed-point implementation of a computation with the same computation performed using floating-point values. As a computation is performed using the mixed-point data type, both a fixed-point value and a floating-point value are computed for every mathematical operation performed.

A. Mixed-Point Data

Each mixed-point variable carries the following information:

- An integral value: I .
- A value for the exponent of the fixed-point scaling factor (an integer): e .
- A long double-precision truncation error value: t .
- A double-precision floating-point value: R .

The mixed-point class currently has no notion of word length. Values of algorithms implemented with limited word length will be found in the least significant bits of mixed-point values, so overflow errors caused by reduced word length will not be modeled correctly.

Accessor functions are provided for the values of I , t , R and $I \cdot 2^e$ (which is the real value mapped to by the fixed-point representation.)

Most of the mathematical operators have been overloaded to work with the mixed-point class. The exceptions are “/=", “/”, “%”, “%=”, “++” and “- -”. The shift operators have all been overloaded, and are interpreted as adding and removing sign bits. The function `mixedpoint::ldexp(const mixedpoint&, int)` is provided for manipulating the scaling factor exponent, similar to the `ldexp` function defined in the C language.¹²

In normal use equation 2 will hold. This equation is only an approximation because both t and R are subject to their own truncation errors, most likely at different resolutions.

$$R \approx (I \cdot 2^e) + t \quad (2)$$

B. Truncation Error

Right shift operations are the only source of truncation error at runtime. Any time a right-shift operation is performed, the bits which will be shifted out of the word are extracted into a floating-point value, multiplied by the scaling factor and added to any existing truncation error already held by the variable. These values give us very accurate (up to the precision of a double floating-point value) information about actual run-time truncation error encountered by the fixed-point implementation.

Constant values can be a source of design-time truncation error. Constants are declared by providing the integer (I) and scaling-factor exponent (e), along with a double-precision value (R). Truncation error is computed in the mixed-point constant constructor as the difference between R and $I \cdot 2^e$. Note that is truncation error is design time truncation error – not run time truncation error.

C. Use of the Mixed-Point Implementation

Mixed-point implementations have been used for the purpose of debugging and evaluating fixed-point algorithm implementations generated by `SIFOpt`.

VIII. Results

`SIFOpt` has been used to create fixed-point implementations of many algorithms. A few of the interesting results are presented here. Each algorithm was actually implemented in C++ using the mixed-point class described earlier, along with additional “instrumentation” to collect maximum and minimum values and run time truncation error amounts. The tests were run on a system with 32 bit integer values, and in most cases (exceptions noted) `SIFOpt` was set to allow the integer values to grow that large.

There are two measures of interest in gauging the quality of a fixed-point implementation of a computation as determined by `SIFOpt`. First is the question of whether the fixed-point implementation utilizes the full range of bits for each computed value when truncation occurs. Second is how closely the truncation error estimate is to actual run time truncation error amounts.

A. 8-Weight Convolution

Several 8-weight convolutions were implemented with a variety of weight distributions, but always with the weight values summing to 1.0. The frequency responses of the resulting filters were computed, and “rail-to-rail” inputs for a variety of the frequency responses were used to test the fixed-point implementations. The input values were encoded as 16-bit values.

Several implementations, differing in internal data path widths, were created: 32-bit internal data paths, 32-bit multiplication results reduced to 16 bits post-multiply and 16-bit internal data paths (multiplications have 16-bit arguments and 16-bit results.)

In each case, the implementation generated by `SIFOpt` used the full range of bits for each computed value – at least for signals of a frequency that should be “passed”.

Implementations having a 32-bit internal data path (hence requiring very few truncations of bits) had truncation error closely matching the estimates computed by `SIFOpt`. In the other cases, `SIFOpt`’s estimates were pessimistic by several orders of magnitude.

(Note that there are no static leaf dependencies for any of the computed values, so the exhaustive value range computation method is not applicable to this computation.)

B. Division by Repeated Multiplications

The division by repeated multiplications algorithm⁷ was implemented to check the accuracy of an 8-bit by 8-bit division with just eight numerator computations. The computation tree has static dependencies between the leaves of some computed values. Before the exhaustive value range computation was implemented the designer had to identify computed values for which `SIFOpt` incorrectly determined value ranges. These nodes were easily identified due to known characteristics of the computation, and protection of the computations using the `NC()` function resulted in fixed-point implementations which used the full range of bits for each computed value. The exhaustive computation of value ranges eliminates the need for the designer to protect any computations with `NC()` functions, and results in an improved fixed-point implementation at the ninth numerator computation and beyond.

After the third numerator computation in the algorithm, large numbers of bits begin to be truncated. Up to that point in the computation, `SIFOpt`’s absolute error estimation tracked with run time truncation amounts. Beyond that point, the error estimation and run time error amounts began to diverge. Figure 7 displays the divergence between the estimated truncation error values computed by `SIFOpt` and actual truncation error amounts observed at runtime.

C. Goertzel’s Fourier Transform Algorithm

The challenge with implementing Goertzel’s Fourier Transform algorithm was that the computation has feedback from the output to the inputs. Because `SIFOpt` requires us to provide an `SIF` for each input, we had to determine an `SIF` to use as a starting point. This was achieved by implementing the computation in a floating-point environment and running the computation on sample data sets. The resulting maximum magnitude value was then used to fix an `SIF` for the algorithm inputs and `SIFOpt` was used to generate a fixed-point implementation with this information.

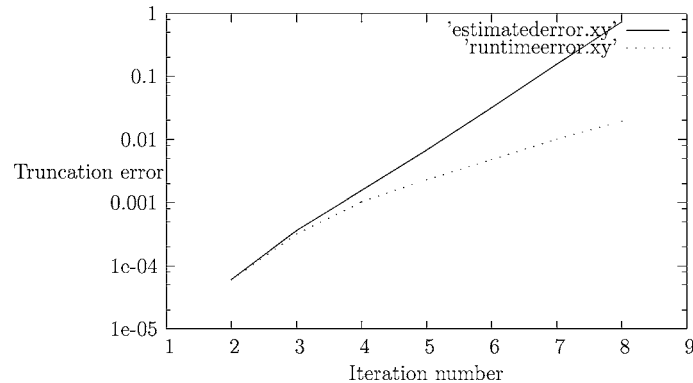


Fig. 7 Graph of the truncation error estimated by SIFOpt versus runtime observed error for each iteration of the division by repeated multiplications computation.

Special care had to be taken to ensure that the scaling factor of the resulting output matched the scaling factor of the feedback values. In addition, because the initial calculation of a value range for the outputs was performed with floating-point values, we took special care to check that the fixed-point implementation was not subject to significant truncation error which would push the result values out of that computed range.

Due to the fact that we pre-computed the maximum feedback value using datasets that matched our testing and the very short computation performed, the fixed-point implementation was virtually guaranteed to have little truncation error and to use the full range of bits. The interesting part of this example was having to deal with a loop in the computation tree.

IX. Conclusions

We have created a software tool to evaluate our static analysis based methodology for creating fixed-point implementations of computations. This tool has been used effectively to create fixed-point implementations of computations for our project. We have demonstrated applications including FIR filtering, division by repeated multiplications, and discrete Fourier transform algorithms. In each case, we were able to use our tool to create a fixed-point implementation in which runtime values utilize the full range of integer bit positions.

Our methodology for estimating truncation error has proved to be fairly inaccurate when large numbers of truncations occur. The estimates can be much larger than are possible during run time. We believe that the estimates can be used to compare different implementations of a computation, but the estimates cannot be relied upon to give a true indication of the magnitude of truncation error caused by the fixed-point implementation.

We intend to work on methods for noting run time dependencies between input values which will result in loss of precision, and to add methods for dealing with decision-making algorithmic structures in the near future.

Acknowledgements

This research has been supported by Dr. Pen-Shu Yeh at the NASA Goddard Space Flight Center. We are grateful for her support and feedback.

References

- ¹Donohoe, G. W., Hass, K. J., Bruder, S., and Yeh, P.-S., "A Reconfigurable Data Path Processor for Space Applications," *Proceedings of the Military and Aerospace Applications of Programmable Logic Devices 2000*, September 2000.
- ²Kum, K., Kang, J., and Sung, W., "AUTOSCALER for C: An Optimizing Floating-Point to Integer C Program Converter for Fixed-Point Digital Signal Processors," *IEEE Transactions on Circuits and Systems II*, Vol. 47, No. 9, September 2000, pp. 840–848.
- ³Aamodt, T. and Chow, P., "Embedded ISA Support for Enhanced Floating-Point to Fixed-Point ANSI C Compilation," *Proceedings of the International Conference on Compilers, Architecture and Synthesis for Embedded Systems*, 2000, pp. 128–137.

⁴Keding, H., Hürtgen, F., Willems, M., and Coors, M., "Transformation of Floating-Point into Fixed-Point Algorithms by Interpolation Applying a Statistical Approach," *Proceedings of the 9th International Conference On Signal Processing Applications and Technology 1998*, 1998.

⁵Chang, M. L. and Hauck, S., "Precis: A Design-Time Precision Analysis Tool," *Proceedings of the 2002 IEEE Symposium on Field-Programmable Custom Computing Machines*, 2002.

⁶Stephenson, M., Babb, J., and Amarasinghe, S., "Bitwidth Analysis with Application to Silicon Compilation," *Proceedings of the ACM SIGPLAN '2000 Conference on Programming Language Design and Implementation (PLDI)*, June 2000.

⁷Parhami, B., *Computer Arithmetic: Algorithms and Hardware Designs*, Oxford, 2000.

⁸Simpson, J. E., "A Block Floating-Point Notation for Signal Processes," Tech. Rep. SAND79-1823, Sandia National Laboratories, 1979.

⁹Buehler, D. M., *A Methodology for Designing and Analyzing Fixed-Point Implementations of Computational Data Paths*, Ph.D. thesis, University of Idaho, 2004.

¹⁰Najjar, W. A., Böhm, W., Draper, B. A., Hammes, J., Rinker, R., Beveridge, J. R., Chawathe, M., and Ross, C., "High-Level Language Abstraction for Reconfigurable Computing," *IEEE Computer*, Vol. 38, No. 8, August 2003.

¹¹Kearfott, R. B., "Interval Computations: Introduction, Uses and Resources," *Euromath Bulletin*, Vol. 2, No. 1, 1996, pp. 95–112.

¹²Harbison, S. P. and Steele, G. L. Jr., *C: A Reference Manual*, Tartan, Inc., 4th ed., 1995.

Carol Traynor
Associate Editor