

An optimization-based approach to LUT program transformations

Chris Wilcox^{1†*}, Michelle Mills Strout^{1†}, James M. Bieman^{1†}

Computer Science Department, Colorado State University, Fort Collins, Colorado, USA

SUMMARY

Scientific programmers can speed up function evaluation by precomputing and storing function results in lookup table (LUTs), thereby replacing costly evaluation code with an inexpensive memory access. A code transform that replaces computation with LUT code can improve performance, however, accuracy is reduced because of error inherent in reconstructing values from LUT data. LUT transforms are commonly used to approximate expensive elementary functions. The current practice is for software developers to (1) manually identify expressions that can benefit from a LUT transform, (2) modify the code by hand to implement the LUT transform, and (3) run experiments to determine if the resulting error is within application requirements. This approach reduces productivity, obfuscates code, and limits programmer control over accuracy and performance. We propose source code analysis and program transformation to substantially automate the application of LUT transforms. Our approach uses a novel optimization algorithm that selects Pareto Optimal sets of expressions that benefit most from LUT transformation, based on error and performance estimates. We demonstrate our methodology with the Mesa tool, which achieves speedups of 1.4-6.9× on scientific codes while managing introduced error. Our tool makes the programmer more productive and improves the chances of finding an effective solution. Copyright © 2013 John Wiley & Sons, Ltd.

Received ...

KEY WORDS: lookup table; performance optimization; error analysis; code generation; scientific computing; memoization; fuzzy reuse; Mesa tool

1. INTRODUCTION

Math-intensive scientific codes are often performance limited by elementary functions such as *sin*, *exp*, and *log* that consume many CPU cycles. For example, calling *cos* in the math library is 35-45× slower than floating-point addition on current architectures. A lookup table (LUT) improves the performance of function evaluation by precomputing and storing function results, thereby allowing replacement of subsequent evaluations with less expensive memory lookups [1, 2]. Practical concerns limit the size and accuracy of LUT results, so LUT methods inevitably introduce error to gain performance.

Hardware designers have long used LUTs to improve elementary function performance [3, 4]. Dedicated memory is expensive in hardware, but linear interpolation and polynomial reconstruction can provide high accuracy with small LUTs [5]. Scientific programmers use similar LUT techniques to expedite function evaluation. For example, the Fastest Fourier Transforms in the West (FFTW) libraries incorporate *cos* and *sin* tables to achieve “significant reductions in computation time” [6]. As another example, the Rapid Radiative Transfer Model (RRTM) software uses LUT transforms for exponential and tau functions, yielding a 1.75× improvement on code that consumes 25% of the execution time of a global climate model [7].

[†]E-mail: wilcox, strout, bieman@cs.colostate.edu

*Correspondence to: Chris Wilcox, Colorado State University, 1873 Campus Delivery, Fort Collins, CO 80523, USA.

Table I. Results of LUT optimization with Mesa.

(Intel Core 2 Duo, E8300, family 6, model 23, 2.83GHz, single core)

Application Name	Original Time	Optimized Time	Performance Speedup	Maximum Error	Memory Usage
Saxs Scattering (discrete)	196.2s	29.0s	6.8X	$4.06 \times 10^{-3}\%$	4MB
Saxs Scattering (continuous)	10.1s	2.5s	4.0X	$1.48 \times 10^{-4}\%$	4MB
Stillinger-Weber (simulation)	14.6s	10.4s	1.4X	$2.91 \times 10^{-2}\%$	1MB
Neural Network (logistics)	8.0s	3.6s	2.2X	$8.70 \times 10^{-2}\%$	4MB
Neural Network (hypertan)	10.9s	3.9s	2.8x	$6.30 \times 10^{-1}\%$	4MB
PRMS (slope aspect)	242ns	56ns	4.3X	$8.21 \times 10^{-6}\%$	4MB
PRMS (solar radiation)	13.7s	6.1s	2.2X	$2.97 \times 10^{-4}\%$	4MB

Lacking a methodology and tools, scientific programmers typically transform source code by manually inserting LUT code. The *ad hoc* nature of such transforms makes it difficult to predict and control application accuracy and performance. Developing performance code by hand is also a substantial effort that impacts programmer productivity [8] and can obfuscate application code [9]. We propose a methodology and algorithms that make LUT development more efficient and effective. We demonstrate our approach with Mesa [10], a tool that uses program analysis and transformation to substantially automate the application of LUT transforms.

Our interest in LUT methods started with application code that we wrote for the Small Angle X-ray Scattering (SAXS) project [11]. We reduced the execution time of our original code by manually incorporating a LUT transform for the dominant calculation. Through lengthy experimentation we achieved an application speedup of 6-7 \times while maintaining reasonable accuracy. However, the manual process was inefficient and provided limited control over accuracy and performance.

Table I shows the results achieved using the latest version of Mesa to performance-tune the SAXS code and five additional scientific applications. The performance speedup, calculated as the ratio of the original time divided by the optimized time, varies from 1.4-6.9 \times . The maximum error shown is the difference between the output of the original and optimized applications. The evaluation of these applications is described in Section 6.

In prior papers, we described versions 1.0 [12] and 1.1 [13] of Mesa. Mesa 1.0 lacked support for domain profiling and linear interpolation and required the user to manually identify candidate expressions in the source code. With Mesa 1.1 a user could apply LUT transformation to expressions in C and C++ source code by identifying the relevant statements with pragmas, thereby operating directly on application source code. The dissertation by Wilcox [14] includes a comprehensive description of our research on automated LUT transformation and the development of Mesa 2.0. This paper extends our SCAM 2012 conference paper [15] by adding information on analytic versus numerical error analysis algorithms, several optimizations associated with expression enumeration, a more thorough treatment of threats to validity and related work, and more information on trends in cache availability and the parallel efficiency of code generated by Mesa.

This paper introduces our approach to optimizing the selection of expressions for LUT transformation to maximize the benefits and minimize costs. The current version of Mesa is a source-to-source translation tool that automatically finds and evaluates candidate expressions through source code analysis and program transformation. Mesa finds the most effective set of LUT transforms by building and solving a multi-objective LUT optimization problem. Our goal is to identify sets of expressions that provide the most performance benefit with the least impact on accuracy. Error analysis and a performance model provide the criteria for choosing between potential LUT transforms.

The contributions of this research are as follows:

- A comprehensive methodology and tool support to apply LUT transforms to source code,
- error estimation and a performance model to characterize LUT transforms,
- a novel LUT optimization technique that maximizes performance and minimizes error, and
- case studies that demonstrate the effectiveness of our methodology and tool.

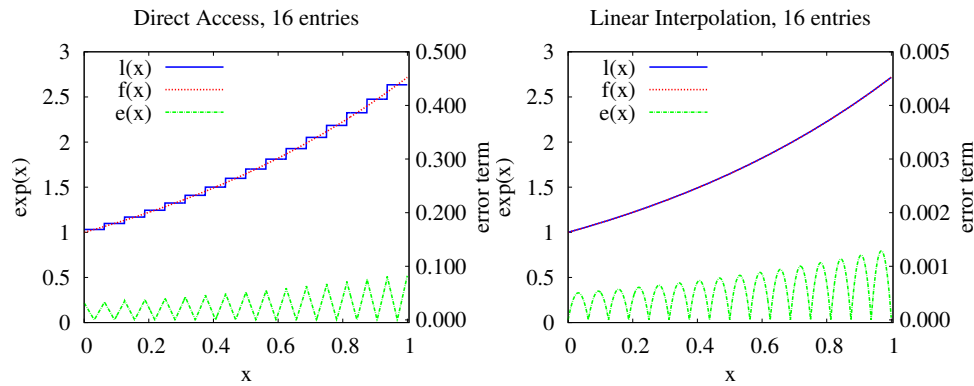


Figure 1. Table Data for Exponential Function

2. BACKGROUND

We present LUT terminology using example LUT data for $\exp(x)$, as shown in Figure 1. We restrict the input values to $0 \leq x \leq 1$ so the LUT *domain* is limited to $[0, 1]$. The exponential function is $f(x)$ and its approximation via the LUT is $l(x)$. Tables are created by partitioning the domain into uniform intervals, and assigning a LUT *entry* for each interval. The number of intervals is the LUT *size*, which is 16 for our example. The interval width represents the LUT *granularity*. Equation 1 shows the relationship between domain, granularity, and size, from which we compute a granularity of 0.0625 for both tables.

$$\text{Granularity} = \text{Domain}/\text{Size} \quad (1)$$

The performance and accuracy of a LUT transform depend on how LUT data is sampled. Common sampling methods are *direct access* and *linear interpolation*. Direct access simply finds the interval that contains the input value and returns its LUT entry. Linear interpolation selects the two closest LUT entries and combines them according to their respective distances from the exact input. Figure 1 shows direct access and linear interpolation sampling for a table with 16 entries.

LUT *error* is the absolute difference between a function and its LUT approximation. In Figure 1 the error is the distance between $f(x)$ and $l(x)$, plotted separately as $e(x)$. Each step in $l(x)$ represents a single LUT entry. For direct access the maximum absolute error is 0.0836; for linear interpolation the maximum absolute error is 0.0013. The equations for computing the maximum error are given in Section 4. The graphs illustrate that linear interpolation yields much smaller error terms (we reduced the error scale as shown on the right side of the graph by two orders of magnitude for linear interpolation to make the smaller error visible).

Regardless of the sampling method, the purpose of a LUT transform is a reduction in execution time that we refer to as LUT *benefit*. We compute the benefit as the difference between the cost of expression evaluation and the cost of LUT access. Linear interpolation has less benefit because of the extra LUT access and computation, and is therefore 1.5-1.8 \times slower than direct access on current architectures.

3. DEFINING THE LUT OPTIMIZATION PROBLEM

Following current practice a programmer must explicitly identify candidate expressions and manually apply LUT transforms. The primary contribution of this paper is the automatic selection of expressions for which LUT transforms are most effective, meaning they provide the highest performance benefit with the least error. Our algorithm (1) enumerates expressions within the subroutines specified by the programmer, (2) estimates the error and performance impact of applying a LUT transform to each of these expressions, and (3) constructs and solves a numerical optimization problem that selects the most beneficial expressions from the enumeration.

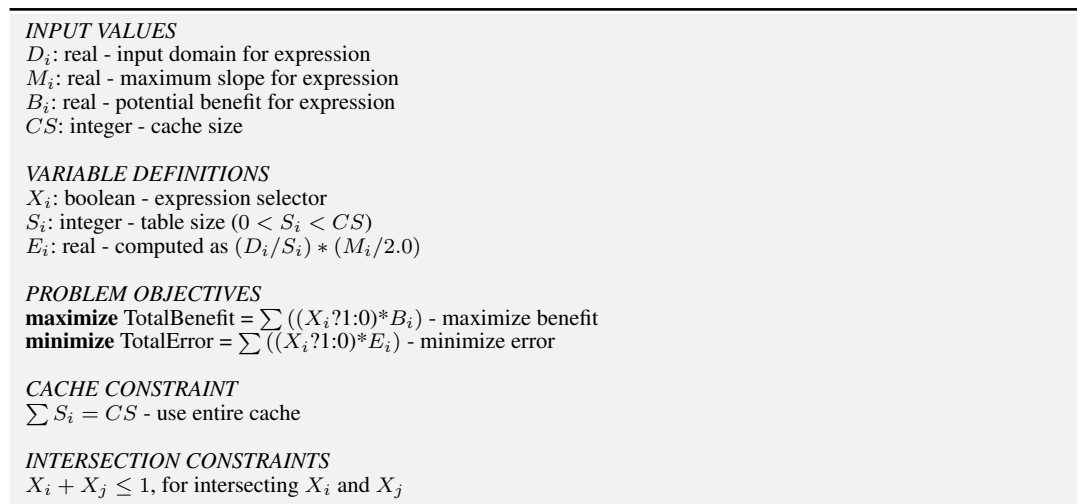


Figure 2. Mathematical definition of optimization problem.

To find the most effective set of LUT transforms, Mesa builds and solves a mixed-integer, non-linear, multi-objective optimization problem. The problem is mixed-integer because error terms are represented as real numbers and memory usage and performance benefits are integer values. The relationship between error and memory usage is non-linear. The problem is multi-objective because it attempts to minimize error while maximizing performance, thus it optimizes multiple independent objective functions. Optimization theory provides several methods for finding *Pareto Optimal* solutions to multi-objective problems, which are defined as those solutions that have the most favorable tradeoff between conflicting objectives.

In optimization terms, a solution is a set of expressions (X_i) for which we wish to minimize the sum of errors (E_i) and maximize the sum of benefits (B_i). The error is computed based on an equation that takes into account the LUT size (S_i), domain (D_i), and slope (M_i) of the function within the domain. The optimization constrains LUT data to fit within the cache size specified by the programmer. Figure 2 shows the mathematical definition of the LUT optimization problem.

Before running Mesa, the programmer identifies the optimization scope by inserting pragmas above subroutine definitions in C source code (or method definitions in C++). Mesa considers LUT transforms for expressions in the bodies of these subroutines. Mesa can optimize any number of methods in a single pass, however, the cost of analysis increases quickly as methods are added. Within the specified methods, Mesa enumerates only expressions that contain elementary function calls, and the optimization is based on a model of the error and benefit for each of these expressions.

The complexity of the LUT optimization problem is exponential, with up to 2^N solutions for N expressions. To reduce the computational complexity, we have developed an algorithm for culling the solution space when N is large. The LUT optimization problem has competing objectives, thus it produces multiple optimal solutions. Our algorithm discards suboptimal solutions, leaving the programmer to select from the much smaller set of Pareto Optimal solutions, which are ranked by accuracy and performance.

4. SOLVING THE LUT OPTIMIZATION PROBLEM

As previously stated, the LUT optimization problem is mixed-integer and non-linear. Optimization frameworks exist that handle these attributes, however, our problem is unusual because the solver must simultaneously select expressions and allocate cache for them. As a result we have not yet been able to use existing solvers including Couenne, Bonmin, and MINLP [16]. Our solution to the LUT optimization problem is divided into two parts. The first part allocates cache memory for the LUT

data associated with each solution. We call this *local optimization*. The second part selects the set of LUT transforms that are Pareto Optimal. We call this *global optimization*. We now present several algorithms that we use to construct and solve the LUT optimization problem. First we describe estimation of the error introduced by a LUT transform. Next we describe the performance model that estimates the potential benefit of a LUT transform. Finally we address our methodology for solving the local and global optimization problems.

4.1. Error Analysis

Error analysis provides an estimate of the maximum error for single LUT transforms, which we use as a proxy for application error. A better solution would be to characterize the propagated error through the entire application. Techniques such as interval analysis can bound introduced error through a sequence of operations [17], but this is beyond the scope of our work as we discuss in 7. Our methodology instead provides support for empirically measuring application error for representative inputs.

Equation (2) shows the error equation for direct access. The maximum error depends on the LUT size, input domain, and the function slope over the domain. The domain is known so we compute only the maximum slope, deferring the LUT size computation until local optimization. The error is inversely proportional to LUT size, so a $2\times$ increase in size decreases error by $2\times$.

$$\text{MaxError} = (\text{Domain}/\text{Size}) * (\text{MaxSlope}/2) \quad (2)$$

Equation (3) shows the error computation for linear interpolation [18]. The maximum delta is the maximum change in the function slope over the domain because the error for linear interpolation is related to the curvature of the function over the interval. Because the error term contains the square of the granularity, a $2\times$ increase in size decreases the error by $4\times$.

$$\text{MaxError} = (\text{Domain}/\text{Size})^2 * (\text{MaxDelta}/8) \quad (3)$$

The error equations require the maximum slope (direct access) or delta slope (linear interpolation) of the function that we are replacing with a table lookup. For elementary functions we can use analytic methods to find the maximum slope over the domain interval. This requires the integral (and sometimes the derivative) of the function, which is not trivial to compute for arbitrary expressions. For this reason Mesa uses numerical methods, which are slower but more robust. Numerical methods must traverse the entire domain to find the maximum slope or delta slope.

The most straightforward numerical method is *exhaustive* traversal of the domain. This method is highly accurate but requires extensive sampling. To improve performance we experimented with *stochastic* sampling, which greatly reduces the number of samples, but we have not found a reliable method for deciding how many samples are needed to get an accurate estimate of the slope. We found the most reliable method to be *boundary* sampling, which evaluates the slopes at the boundaries of each LUT interval. For direct access, the method evaluates the function at the left and right boundary of each LUT interval, giving us the two values needed compute the slope. For linear interpolation, the method adds an evaluation at the interval center, which gives us the three values needed to compute the delta slope. The boundary method has proven to be very accurate for tables of all sizes, and it has the advantage that the number of samples required is proportional to the LUT size instead of the domain size. A comparison of error methods is shown in [14].

4.2. Performance Modeling

Performance modeling estimates the benefit associated with each LUT transform based on the execution time of arithmetic operators, elementary functions, and memory access. Mesa incorporates a benchmark that measures the average performance of these operations for direct access and linear interpolation.

Equation 4 shows the performance model, which uses a count of the arithmetic operators and elementary function calls per expression. These counts are multiplied by the cost of each operation,

and the LUT access time is subtracted. For linear interpolation, the result is divided by the relative performance compared to direct access. This result is an estimate of the benefit of replacing the expression with a LUT access.

$$\text{Benefit} = ((\text{Cost}(\text{Op}) * \text{Count}(\text{Op})) - \text{Cost}(\text{Access})) * \text{Frequency} \quad (4)$$

For example, consider the optimization of an expression with a cosine call. On our test system the sine call takes 45ns. Subtracting 7.4ns for the LUT access gives a savings of 37.6ns per execution. We multiply this by the call frequency to get the total expression benefit. For example, a call frequency of 10^8 for the same expression would yield a benefit $37.6\text{ns} \times 10^8$ or 3.76s. When the expression is optimized with linear interpolation, the benefit would be reduced by the relative factor, which is $\sim 1.8 \times$ on our test system, leading to a benefit of 2.09s.

We do not expect the performance model to predict exact timing. Instead, it establishes the relative performance to allow comparison of solutions. Even so, we are often within 5-10% when estimating performance speedup for individual elementary functions. When optimizing code with multiple elementary function calls, the model tends to overestimate benefit, mainly because it does not model compiler optimizations. However, our case studies show that our estimates are usually within $\sim 1-2 \times$ of the actual performance benefit.

4.3. Local Optimization

The purpose of local optimization is to find the optimal allocation of cache resources for the set of expressions within each solution. Equation (5) shows the closed-form formula that computes LUT sizes for each LUT transform to minimize solution error [14]. A variant of the formula directly computes the maximum error for a single LUT access, as shown in Equation (6).

$$S_i = CS / \left(\sum_{j=1..n} \sqrt{(M_j D_j / M_i D_i)} \right) \quad (5)$$

$$E_i = M_i D_i \left(\sum_{j=1..n} \sqrt{(M_j D_j / M_i D_i)} \right) / CS \quad (6)$$

4.4. Global Optimization

We find Pareto Optimal solutions by sorting the solutions by estimated error and performance, and selecting solutions that lie on the convex hull. To find the convex hull we use a modified Graham Scan algorithm [19]. Figure 3 shows the Pareto chart for the example code presented in Section 5. Solution error is on the x-axis and solution benefit is on y-axis. To make the chart we plotted both the Pareto Optimal (circles) and suboptimal (triangle) solutions. The Pareto solutions are joined by a line called the Pareto curve. The optimal solutions C0 through C4 lie above and the left of suboptimal solutions, because they provide more performance with the same or less error.

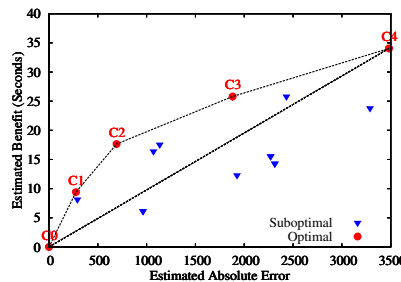


Figure 3. Pareto chart for example program.

As demonstrated in Figure 3, the Pareto curve gives insight into the effectiveness of the LUT transforms that are combined into Pareto solutions. LUT transforms that contribute lots of benefit but little error appear on the steep left side of the curve. LUT transforms that introduce lots of error but little performance appear on the flatter right side of the curve. The programmer can examine the Pareto curve and decide how much benefit is possible for the amount of error they can tolerate.

5. MESA TOOL

We have implemented error and performance models and optimization algorithms from Section 4 in a standalone tool called Mesa. Mesa incorporates ROSE compiler infrastructure [20, 21] to support static and dynamic source analysis and program transformation on C/C++ application code. Figure 4 shows the automated optimization methodology implemented by Mesa.

Figure 5 shows the example program that we use to explain our methodology and tool. The performance of this code is limited by the computation of elementary functions, and the domain of input values for the program enables us to satisfy accuracy requirements while fitting in cache memory. We now describe the six stages of the methodology displayed in Figure 4 in terms of the example code.

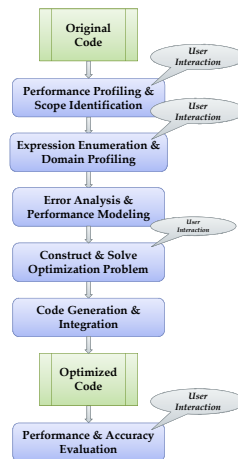


Figure 4. Methodology for automated optimization.

```

#pragma LUTOPTIMIZE
double ScatterSample(Sample sample, ...
{
S35     double dProduct;
S36     double dSum0 = 1.0;
S37     double dSum1 = 1.0;
S38
S39     // Iterate geometry
S40     for (int j = 0; j < vGeometry.size(); j++)
S41     {
S42         dProduct = (sample.x * vGeometry[j].x) ...
S43         dSum0 += exp(dProduct) + sin(dProduct);
S44         dSum1 += exp(dProduct) + cos(dProduct);
S45     }
S46
S47     // Return answer
S48     return dSum0 * dSum0 + dSum1 * dSum1;
}
  
```

Figure 5. Source listing for example program.

Table II. Enumerated expressions for example code.

Expression Identifier	Expression Description	Statement Identifiers	Input Variables
X0	exp(dProduct)	S43	dProduct
X1	sin(dProduct)	S43	dProduct
X2	exp(dProduct) * sin(dProduct)	S43	dProduct
X3	exp(dProduct)	S44	dProduct
X4	cos(dProduct)	S44	dProduct
X5	exp(dProduct) * cos(dProduct)	S44	dProduct

Original: X0 \cap X2, X1 \cap X2, X3 \cap X5, X4 \cap X5
 Introduced: X0 \cap X6, X3 \cap X6
 Inherited: X2 \cap X6, X5 \cap X6

Figure 6. Intersection constraints for example program.

5.1. Stage 1: Performance Profiling and Scope Identification

In the first stage the user *profiles performance* by manually running a tool such as **gprof** to find the most costly subroutines. For the example code, the subroutine ScatterSample consumes >90% of the execution time. The user inserts the pragma shown in Figure 5 to identify the *optimization scope* as the body of this subroutine. The pragma causes Mesa to evaluate expressions in statements S35 through S48 for potential LUT transforms.

5.2. Stage 2: Expression Enumeration and Domain Profiling

In the next stage Mesa *enumerates expressions* that are candidates for LUT transformation. Our enumeration extracts expressions from statements within the optimization scope specified by the user, rejecting those that do not match our criteria. For example, expressions must contain one or more elementary function calls, since these functions are the focus of our methodology. Enumeration extracts both individual elementary functions and more complex expressions that combine elementary functions with arithmetic operators. Complex expressions are considered because additional performance gains can occur when an expression with multiple elementary functions is handled by a single LUT. The handling of complex expressions can also reduce memory usage by with reducing the amount of LUT data.

Table II lists the enumerated expressions for statements S43 and S44 in Figure 5, labeled as X0 to X5. For each expression we show the expression identifier, syntax, statement identifier, and input variable. Extracted expressions can overlap other expressions in the same statement. For example, expressions X0 and X1 are subsumed by X2 and expressions X3 and X4 are subsumed by X5, so they cannot be simultaneously optimized. Mesa maintains *intersection constraints* that prevent overlapping expressions from being combined into a solution. Figure 6 shows the original intersection constraints for the example code.

During the enumeration stage Mesa performs a number of additional transformations that reduce memory usage. First, Mesa combines similar expressions via *expression coalescing*. Coalescing saves memory by sharing LUT code and data between two or more identical expressions. For the example code, coalescing combines the identical exponential calls in X0 and X3 to create a new expression X6. The coalesced expression introduces new constraints, since X0 and X3 are subsumed by X6, and it inherits old constraints as shown in Figure 6

Second, Mesa supports *parameter merging*, which avoids the realization of unnecessary multi-dimensional LUTs by merging input variables. To illustrate this we consider optimization of the statement $x = \exp(y/z)$. Mesa extracts expressions for the exponential by itself and the entire expression with the division. Optimizing the exponential is done with a single-dimensional LUT indexed by (y/z) , whereas the entire expression requires a multi-dimensional LUT indexed by y and z . The latter requires significantly more memory for the same level of accuracy.

Table III. Input data for optimization problem.

Expression Identifier	Expression Description	Statement Identifiers	D_i	M_i	B_i
X0	exp(dProduct)	S43	2.44	3.31	6.15s
X1	sin(dProduct)	S43	2.44	1.00	8.15s
X2	exp(dProduct) + sin(dProduct)	S43	2.44	3.67	16.40s
X3	exp(dProduct)	S44	2.44	3.31	6.15s
X4	cos(dProduct)	S44	2.44	0.95	9.40s
X5	exp(dProduct) + cos(dProduct)	S44	2.44	2.38	17.65s
X6	exp(dProduct)	S43, S44	2.44	3.31	12.30s

Finally, Mesa implements *domain conditioning* to take advantage of the cyclical nature of some elementary functions, thus saving memory by reducing the input domain of the associated LUT transform. For example, the sine and cosine functions are completely represented by the interval from 0 to 2π radians. Input values outside of this domain can be mapped back into the interval by a modulo operation or repeated addition and subtraction. We call this technique domain conditioning, also referred to in the literature as *parameter folding* or *range reduction* [22].

After expression enumeration Mesa adds *domain profiling* calls to instrument the source code to capture runtime data such as the domain of input variables and the execution frequency of statements. The user must compile and run the instrumented code to perform the dynamic analysis. At program completion the instrumented code stores profiling results for subsequent stages.

5.3. Stage 3: Error Analysis and Performance Modeling

Mesa applies the algorithms for error analysis and the performance model from Section 4 to compute the maximum slope and performance benefit. Table III shows the candidate expressions and their estimated error and benefit parameters, which is the input data for the optimization problem. D_i is the input domain, M_i is the maximum slope, and B_i is the potential benefit in seconds.

5.4. Stage 4: Solving the Optimization Problem

Mesa builds and solves the LUT optimization problem. Figure 7 is a partial listing from Mesa that shows the optimization result, ending with the presentation of Pareto Optimal solutions to the programmer. Mesa displays the number of possible solutions as 2^7 or 128 for the 7 expressions extracted from the example code, and the number of actual solutions as 29 after intersection culling. The $4\times$ decrease in solutions is not unusual when expression coalescing is enabled because of introduced intersection constraints. From these, Mesa finds five Pareto optimal solutions numbered from C0 to C4. Mesa shows these to the programmer along with the corresponding error and benefit estimates. The listing gives absolute error, benefit in nanoseconds, and lists the expressions that comprise each solution.

Figure 7 shows that the solutions range from C0, which has zero benefit and error, to C4, which has the maximum benefit. The latter combines expressions X2 ($\exp(dProduct) * \sin(dProduct)$) and X5 ($\exp(dProduct) * \cos(dProduct)$), each of which contains two elementary function calls. The coalesced exponential call X6 ($\exp(dProduct)$) does not appear in the solution so expression coalescing has not helped.

The listing further shows user selection of solution C4, which causes Mesa to replace expressions X2 and X5 with LUT accesses. The X2 expression receives a cache allocation of 2270KB. This is larger than the 1826KB allocation of X5 because X2 has a larger maximum slope. The global optimization has identified the most effective LUT transforms based on the error and performance model, for which the local optimization has allocated the ideal amount of cache.

```

128 solutions (possible)
29 solutions (actual)
5 solutions (pareto optimal)
Solution  Error      Benefit      Optimizations
C0         0.000e+00  0.000e+00
C1         2.759e+02  9.400e+09  X4
C2         6.921e+02  1.765e+10  X5
C3         1.881e+03  2.580e+10  X5,X1
C4         3.483e+03  3.405e+10  X5,X2
Select solution: 4
X5 Di=2.44 Mi=2.38 Ei=1.553e+03 Bi=1.765e+10
   Si = 1869833 (1826KB)
X2 Di=2.44 Mi=3.67 Ei=1.930e+03 Bi=1.640e+10
   Si = 2324471 (2270KB)
Mesa 2.0: Generating optimized code
Realizing X5 in statement S44
Realizing X2 in statement S43
Mesa 2.0: Optimization completed.

```

Figure 7. Mesa optimization result for example code.

5.5. Stage 5: Code Generation and Integration

After the programmer selects a solution, Mesa realizes the code for the specified set of LUT transforms, including the LUT data, constructor, destructor, initialization code, table access, and original function. Mesa then integrates the generated code into the application and replaces original expressions with LUT access calls. An optimized version of the application code is written to the file specified on the command line. The programmer rebuilds the application using the normal build process, and the resulting executable should behave in an identical manner to the original program, except for differences in performance and accuracy. Figure 8 shows a partial listing of the code generated by Mesa. To save space, we show only the code for the X5 expression ($\exp(dProduct) * \cos(dProduct)$) and the modifications to the original subroutine.

5.6. Stage 6: Accuracy and Performance Evaluation

In the last stage, the programmer evaluates the benefit and accuracy of the optimized program version against the original version. We compute performance speedup as the original execution time divided by the optimized execution time. Accuracy is evaluated as the absolute or relative deviation from original program output. The optimized version of the program shows a $8.1 \times$ speedup over the original version. We compute a relative error of $1.13 \times 10^{-4}\%$, based on the accumulation of return values from the ScatterSample subroutine.

The evaluation stage completes when the programmer has the accuracy and performance information needed to evaluate whether the optimization is worthwhile. The programmer can accept the optimization and use the code generated by Mesa, run Mesa again and select another solution, or revert to the original code. A Mesa parameter specifies the selection of a solution so that the iterative process just described can be scripted.

5.7. Mesa Workflow

Figure 9 shows the Mesa workflow. The programmer inserts pragma statements in the source code to identify the optimization scope, and runs Mesa to instrument the source code for domain profiling (Step 1). Next the programmer compiles and runs the instrumented code (Step 2) to capture domain profiles (Step 3). The programmer runs Mesa again to request program optimization (Step 4). Based on user selection of a Pareto Optimal solution, Mesa generates optimized code (Step 5) that the programmer compiles, runs, and compares with the original program to evaluate performance and accuracy differences (Step 6).

```

// Start of code generated by Mesa, version 2.0
// LUT constants
const double X5_lower = -1.2461340000e+00;
const double X5_upper = 1.1962890000e+00;
const double X5_gran = 5.2249019600e-06;
class CLut {
public:
    // LUT Constructor
    CLut() {
        double dIn, dOut;
        for (double dIn=X5_lower; dIn<=X5_upper; dIn+=X5_gran) {
            dOut = X5_orig(dInput+(X5_gran/2.0));
            X5_data.push_back(dOut);
        }
    }
    // LUT Destructor
    ~CLut() {
        X5_data.clear();
    }
    // LUT function for expression
    float X5_lut(float X5_param) {
        X5_param -= X5_lower;
        int uIndex = (int) (X5_param*(1.0/X5_gran));
        return(X5_data[uIndex]);
    }
    // Original expression
    double X5_orig(double dProduct) {
        return (exp(dProduct)+cos(dProduct));
    }
private:
    // LUT data structures
    std::vector<float> X2_data;
};
// Object instantiation
CLut clut;
// End of code generated by Mesa, version 2.0

// Expressions replaced by Mesa
S43 dSum0 += clut.X2_lut(dProduct);
S44 dSum1 += clut.X5_lut(dProduct);

```

Figure 8. Optimized code generated by Mesa.

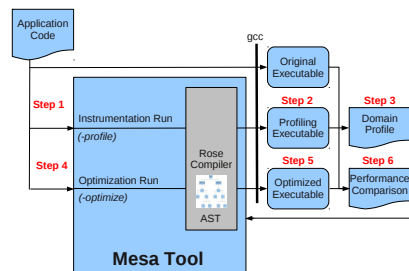


Figure 9. Diagram of Mesa workflow.

The code instrumentation and generation are implemented by calling the Rose libraries to read the original code and parse it into an abstract syntax tree (AST). Mesa analyzes the portion of the AST identified by the pragma to find expressions that may benefit from LUT transforms, then instruments or optimizes the code by modifying the AST and calling Rose to unparse it back into C/C++ application code.

Table IV. Application results from Mesa optimization.

(Intel Core 2 Duo, E8300, family 6, model 23, 2.83GHz, single core)

Application Name	Lines of Code	Number of Expressions	Possible Solutions	Actual Solutions	Pareto Solutions	Processing Time
PRMS Slope Aspect	35	9	512	384	9	13.7s
PRMS Slope Aspect	35	11	2048	425	9	15.5s
PRMS Solar Radiation	7	6	64	64	8	14.1s
SAXS Discrete	60	3	8	4	3	11.2s
SAXS Discrete	60	3	8	4	3	16.5s
SAXS Continuous	30	5	32	20	4	10.8s
Stillinger-Weber	44	6	64	36	3	9.3s
Neural Network (logistics)	5	2	4	3	2	4.9s
Neural Network (hypertan)	5	1	2	2	2	2.8s

5.8. Tool Limitations

The following limitations apply to the tool, not the methodology. Mesa parses only C and C++ code and generates only C++. Mesa also only handles a single source module. Additionally there are syntactic elements that are not handled, including type casts and structure and pointer access. Mesa handles only the following functions: sin, asin, sinh, cos, acos, cosh, tan, atan, tanh, exp, log, and sqrt. In addition, Mesa optimizes assignment expressions and initializers, but will not find computations in other constructs.

6. EVALUATION

We evaluate our methodology in terms of ease of use, accuracy, and performance by using Mesa to optimize six scientific applications. Two case studies evaluate the *slope aspect* and *solar radiation* computations from the Precipitation-Runoff Modeling System (PRMS), developed by the United States Geologic Survey [23]. The third and fourth case studies come from two applications written for the previously cited SAXS project [11]. The fifth application is Stillinger-Weber, a molecular dynamics program developed and used for research at Cornell University [24]. The sixth application is neural network code [25] developed by a faculty member in our department.

6.1. Evaluation Methodology

To evaluate our methodology we apply the process that a programmer would follow to use Mesa. We run gprof to find bottleneck subroutines that we identify with a pragma. We run Mesa to instrument the application for domain profiling, and to optimize the application. Part of our evaluation is to examine the number of possible and actual solutions for each application, and we measure tool processing time. We finish by compiling and running the original and optimized code and comparing the benefit and error predicted by Mesa to the actual application performance and accuracy.

6.2. Evaluation Results

Table IV summarizes the results of our case studies. The table shows program and tool statistics: lines of code analyzed, number of expressions, number of possible, actual, and Pareto Optimal solutions, followed by the tool processing time. Only the lines of code in the function optimized by Mesa are listed. For example, the Neural Network code has a simple transfer function with only five lines of code. Refer to Table I for performance speedup and error relative to the original output. For example, Mesa extracts 9 expressions from the PRMS slope aspect code, from which it analyzes 384 solutions, and finds 9 to be Pareto Optimal. The processing time is 13.7 seconds.

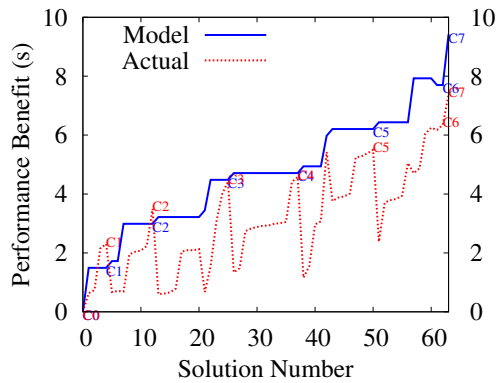


Figure 10. Performance model comparison

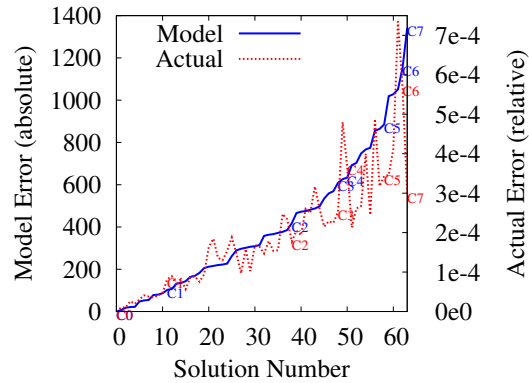


Figure 11. Error model comparison

6.3. Model Evaluation

We now evaluate our error and performance model when applied to the solar radiation application. Figure 10 shows that the performance model predicts the trend of the actual performance, but overestimates the benefit. Some solutions are clearly below the estimated performance, which may be due to instruction-level parallelism and other compiler optimizations. Figure 11 compares the estimated maximum error from the error analysis against the actual maximum error of the optimized application. The model again correctly predicts the trend but has some local variation. We attribute this to the difficulty of modeling application error. Despite being able to quantify many aspects of the error introduced by a LUT transform, a general method to compute the effect on application accuracy remains an open problem. The propagation of error through arbitrary sections of application code poses a complex numerical analysis problem that is unique to each application. This implies that some level of experimentation will still be necessary to evaluate the accuracy of a LUT transform.

6.4. Parallel Execution

Given the prevalence of multi-core architectures, we can justify sequential optimizations only if they are effective in the context of parallel execution. A critical factor for successful LUT optimization on multi-core systems is cache availability, because LUT data must reside in cache to be beneficial. A higher number of cores raises the concern that manufacturers may be unable to maintain the current levels of per core cache. However, recent cache size trends do not support this hypothesis. Figure 12 shows the combined size of L2 and L3 cache per core, on historical and recent Intel processors with up to 8 cores. The L2 caches of uniprocessor systems grew quickly from 256KB in 1995 to 1MB or more in 2004. Since that time there is no evidence of a reduction in combined L2 and L3 cache size, despite the growth in the number of cores. In our case studies we confirm that sufficient cache exists to support LUT transformation on existing multi-core systems.

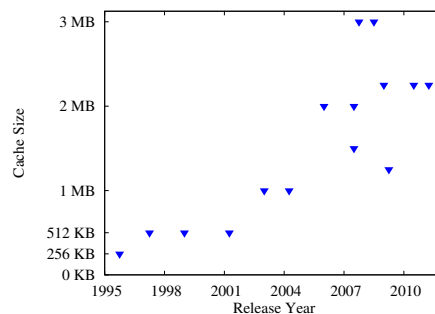
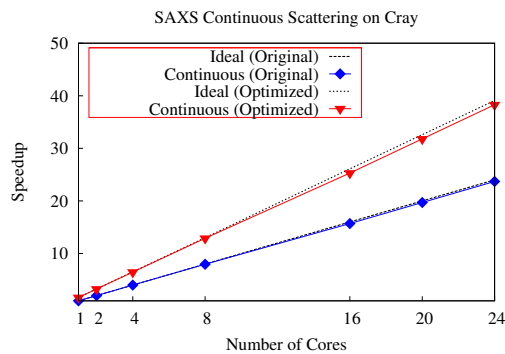


Figure 12. Combined L2 and L3 cache per core.



(Cray XT6m, AMD Opteron 6100, 2.5Ghz, 512KB L2, 6MB L3, 24 cores)

Figure 13. Parallel efficiency of SAXS continuous scattering application.

To confirm the effectiveness of LUT techniques in multi-core environments, we used Mesa to evaluate the parallel performance of LUT code generated by Mesa. Figure 13 shows parallel efficiency of 98% for continuous scattering on a 24-core Cray XT6m computer [13]. We have replicated this result on multi-core systems including 4 and 8-core Xeon, concluding that our single-core optimizations are independent from and complementary to parallelization.

6.5. Summary and Evaluation

Our case studies demonstrate that LUT optimizations are effective on the applications shown, because they show significant performance improvements while meeting the accuracy requirements of the application. Mesa also reduces the programming effort in several ways. First, our tool automates source code analysis, which frees the programmer from having to search for candidate expressions. Second, Mesa gives the programmer error and performance estimates without them having to manually instrument code. Third, our tool generates and integrates LUT code, freeing the programmer from coding. Fourth, our tool computes optimal cache allocations for LUT data using a method that would be very cumbersome by hand. Finally, the tool simplifies experimentation with different solutions.

Mesa has improved our own process for tuning applications whose performance is bound by elementary functions. Our original *ad hoc* LUT implementation for the SAXS discrete code required several weeks of development time and experimentation, even after the base algorithm was implemented and tested. Characterization of error and performance was especially time-consuming, because it required multiple runs of the entire SAXS application. In addition, we simply had no way to estimate the performance or error impact of our LUT transforms. In contrast, we can now evaluate applications with Mesa and quickly receive feedback on whether LUT methods will help the application. For example, we were able to optimize the SAXS continuous code in a matter of minutes, achieving the results shown in Section 6.

7. THREATS TO VALIDITY

The primary threat to validity for our research is external validity. Empirical research is always limited with respect to the number and scope of the applications that can be evaluated. Our empirical evaluation consists of case studies of six applications in four scientific areas, of which two were partially written by the authors. Further research is required to demonstrate applicability to other domains. However, we expect that our results will generalize to applications that have the same limitations on performance caused by elementary function calls, assuming that other environmental factors (compilers, languages, hardware) are consistent.

Other threats to validity are as follows. First, LUT methods depend on the relative performance of function evaluation versus memory access, which can change as processor architectures evolve.

Second, our performance model sometimes overestimates benefit because it does not account for compiler optimizations. Third, our error model considers only introduced error, not the error that propagates through the application. With respect to the latter threat, there is considerable precedence for methods that trade accuracy for performance. Computer precision is inherently limited, yet the existing floating-point representations have proven to be satisfactory for many scientific applications. Linderman et al. [26] argue that reducing the precision of computation has benefits, but caution that a careful analysis is necessary to maintain accuracy. Research suggests that single-precision arithmetic can sometimes be used in place of double-precision [27]. Some libraries give the programmer explicit control over accuracy and performance. For example, the Intel Math Kernel Library (MKL) supports an enhanced performance mode that throws away up to half of the significant bits for single or double-precision [28].

The existence of methods that trade accuracy for performance does not necessarily imply that all applications can do so safely. In the citations listed above, the burden of analyzing numerical stability is placed on the programmer. The numerical stability of an application is determined by many factors [29], so some applications are more sensitive to the introduction of error than others. Techniques such as interval analysis can be used to bound error on sequences of operations [17], but characterizing the error propagation for an entire application is a complex numerical analysis problem [30]. Our methodology estimates the error introduced by a set of LUT approximations, and we provide support for empirically measuring application error.

8. RELATED WORK

The reuse of computation is fundamental to a number of optimization techniques. Many programs exhibit *value locality*, in which computations are performed repeatedly with a small number of inputs [31]. Value locality creates redundancy that can be exploited to improve performance. Optimizations based on *value reuse* eliminate redundant computation by storing and reusing results from previous computations [32]. Some compilers implement value reuse by selectively caching results to avoid future computation [33]. Hardware can also exploit value reuse by caching results or sequences of instructions, or through speculative value prediction [34]. Most implementations of value reuse apply the method to a small set of precise inputs over localized areas in a program.

Value reuse can be exploited even without value locality. Consider a program that repeatedly evaluates one or more expressions with identical or similar inputs. Caching previously computed results can provide a benefit even when the computation is distributed throughout the program, if the reuse is applied globally. Performance gains depend on the reuse inherent to the program, which provides the opportunity to eliminate redundant computation. We use the term *precise reuse* when input values must exactly match [35]. Precise reuse provides the same accuracy as the original expression, if the cache has sufficient precision. Alvarez et al. introduced *fuzzy reuse*, in which input values can match imprecisely [36]. Fuzzy reuse improves performance by sharing cached results between multiple inputs, at the expense of a loss of accuracy. LUT transformation for function evaluation is based on fuzzy reuse.

Memoization is a related technique that reuses previous evaluation results to avoid future computation. Memoization employs a mapping function to map input values to cached results. The mapping function can require precise reuse or the comparison can be imprecise or fuzzy [36]. Memoization differs from LUT methods in that results are cached only as needed, instead of computing the entire table in advance [35], thus requiring extra work to determine if a cached result exists for the current computation.

Some of the history of hardware LUTs is presented in Section 1. Software LUTs have few academic references, but some books [2, 37] discuss the topic. We have found one tool that supports software LUT transforms [38]; it is a standalone compiler that analyzes mathematical expressions written in a language that is similar to MATLAB, and transforms these expressions either into an FPGA design or C/C++ code. Although the compiler described in the Zhang et al. paper does not operate directly on source code, their work provides a unique discussion of software LUT issues that is very relevant to our work.

We are not aware of other published work on domain profiling, error analysis, or performance modeling in the context of software LUTs, nor have we seen case studies that characterize the performance versus accuracy tradeoff, or that explore the cost versus benefit of the various LUT sampling methods. The issue of whether LUT data must reside in cache to make LUT transformation beneficial is especially important, but the only related work that we are aware of that discusses the impact of cache usage is Defour [39]. Defour combines polynomial reconstruction with small tables, and suggests that LUT data must fit into L1 cache to be effective. Our research has shown a substantial performance gain even when LUT data resides in L2 or L3 cache [12].

9. CONCLUSION

This paper demonstrates a novel approach to LUT optimization that is supported by error analysis and performance estimation. Our methodology and Mesa tool substantially automates the application of LUT optimization to scientific programs. Our approach is effective at speeding up code that is performance limited by elementary function calls. Case studies demonstrate speedups from 1.4-6.8 \times with reasonable accuracy. Automation improves programmer productivity by reducing the effort required to identify and implement LUT transforms, and by providing information that helps the programmer make the critical tradeoff between error and performance. The Mesa tool provides an alternative to current *ad hoc* practices that require significant programmer effort. This paper extends [15] with additional information on analytic versus numerical algorithms for error analysis, optimizations performed during expression enumeration, trends in cache availability, and the parallel efficiency of code generated by Mesa.

ACKNOWLEDGMENT

We acknowledge support for our research from Award Number 1R01GM096192 from the National Institute Of General Medical Sciences. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institute Of General Medical Sciences or the National Institutes of Health. Our research is also supported by grant number DESC0003956 from the Department of Energy, with additional seed funding from the Vice President of Research and the Office of the Dean of the College of Natural Sciences at Colorado State University (CSU) and from a Department of Energy Early Career grant DE-SC3956. This research utilized the CSU ISTeC Cray HPC System supported by NSF Grant CNS-0923386.

REFERENCES

1. Lee DU, Abdul Gaffar A, Mencer O, Luk W. Optimizing Hardware Function Evaluation. *IEEE Trans. Comput.* Dec 2005; **54**:1520–1531.
2. Pharr M, Fernando R. *GPU Gems 2: Programming Techniques for High-Performance Graphics and General-Purpose Computation*. Addison-Wesley Professional: Boston, MA, USA, 2005.
3. Gal S. Computing Elementary Functions: A New Approach for Achieving High Accuracy and Good Performance. *Proceedings of the Symposium on Accurate Scientific Computations*, Springer-Verlag: London, UK, 1986; 1–16.
4. Tang PTP. Table-lookup Algorithms for Elementary Functions and their Error Analysis. *Proceedings of the 10th IEEE Symposium on Computer Arithmetic*, 1991.
5. Piñeiro JA, Bruguera JD, Muller JM. Faithful Powering Computation Using Table Look-Up and a Fused Accumulation Tree. *ARITH '01: Proceedings of the 15th IEEE Symposium on Computer Arithmetic*, IEEE Computer Society: Washington, DC, USA, 2001.
6. Jones DL. Efficient FFT Algorithm and Programming Tricks. *Connexions* 2007; <http://cnx.org/content/m12021/1.6/>.
7. Rapid Radiative Transfer Model. 2010. http://rtweb.aer.com/rrtm_frame.html.
8. Faulk S, Porter A, Gustafson J, Tichy W, Johnson P, Votta L. Measuring HPC productivity. *International Journal of High Performance Computing Applications* 2004; **2004**:459–473.
9. Loh E, Van De Vanter ML, Votta LG. Can Software Engineering Solve the HPCS Problem? *Proceedings of the Second International Workshop on Software Engineering for High Performance Computing System Applications*, ACM: New York, NY, USA, 2005; 27–31.
10. MESA Project. 2012. <http://www.cs.colostate.edu/hpc/MESA>.

11. SAXS Project. 2010. <http://www.cs.colostate.edu/hpc/SAXS>.
12. Wilcox C, Strout M, Bieman J. Mesa: Automatic Generation of Lookup Table Optimizations. *Proceedings of the 4th International Workshop on Multicore Software Engineering, IWMSE '11*, ACM: New York, NY, USA, 2011.
13. Wilcox C, Strout M, Bieman J. Tool support for software lookup table optimization. *Scientific Programming* Dec 2011; **19**(4):213–229.
14. Wilcox C. A methodology for automated lookup table optimization of scientific applications. PhD Thesis 2012.
15. Wilcox C, Strout MM, Bieman J. Optimizing expression selection for lookup table program transformation. *Proceedings of the 12th IEEE International Working Conference on Source Code Analysis and Manipulation (SCAM)*, 2012.
16. NEOS Solvers. 2012. <http://www.neos-server.org/neos/solvers/index.html>.
17. Moore RE, Bierbaum F. *Methods and Applications of Interval Analysis*. Society for Industrial and Applied Math (SIAM): Philadelphia, PA, USA, 1979.
18. Epperson J. *An Introduction to Numerical Methods and Analysis*. John Wiley & Sons: New York, NY, USA, 2007.
19. HCormen T, Leiserson CE, LRivest R, Stein C. *Introduction to Algorithms (Second Edition)*. The MIT Press: Cambridge, MA, USA, 2001.
20. ROSE Project. 2011. <http://www.rosecompiler.org/>.
21. Liao C, Quinlan DJ, Panas T, de Supinski BR. A ROSE-Based OpenMP 3.0 Research Compiler Supporting Multiple Runtime Libraries. *IWOMP*, 2010; 15–28.
22. Sobti K, Deng L, Chakrabarti C, Pitsianis N, Sun X, Kim J, Mangalagiri P, Irick K, Kandemir M, Narayanan V. Efficient Function Evaluations with Lookup Tables for Structured Matrix Operations. *2007 IEEE Workshop on Signal Processing Systems*, 2007.
23. PRMS Project. 2010. <http://water.usgs.gov/software/PRMS>.
24. Haran M, Catherwood JA, Clancy P. Diffusion of Group V Dopants in Silicon-Germanium Alloys. *Applied Physics Letters* Apr 2006; **88**(17):173–502.
25. Neural Network Software. 2011. <http://www.cs.colostate.edu/~anderson/meOther.html>.
26. Linderman MD, Ho M, Dill DL, Meng TH, Nolan GP. Towards program optimization through automated analysis of numerical precision. *Proceedings of the 8th annual IEEE/ACM international symposium on Code generation and optimization, CGO '10*, ACM: New York, NY, USA, 2010; 230–237.
27. Buttari A, Dongarra J, Kurzak J, Luszczek P, Tomov S. Using Mixed Precision for Sparse Matrix Computations to Enhance the Performance while Achieving 64-bit Accuracy. *ACM Transactions on Mathematical Software* 2008; **34**:1–22.
28. Intel. *Intel Math Kernel Library Reference Manual*. Beaverton, OR, USA 2011.
29. Goldberg D. What every Computer Scientist should know about Floating-point Arithmetic. *ACM Computing Surveys* Mar 1991; **23**:5–48.
30. Chapra SC, Canale RP. *Numerical Methods for Engineers: With Programming and Software Applications*. 3rd edn., McGraw-Hill: New York, NY, USA, 1997.
31. Sastry SS, Bodik R, Smith JE. Characterizing Coarse-Grained Reuse of Computation. *3rd ACM Workshop on Feedback Directed and Dynamic Optimization*, 2000; 16–18.
32. Kumar KVS. Value Reuse Optimization: reuse of Evaluated Math Library Function Calls through Compiler Generated Cache. *SIGPLAN Notices* Aug 2003; **38**:60–66.
33. Ding Y, Li Z. A Compiler Scheme for Reusing Intermediate Computation Results. *Proceedings of the International Symposium on Code Generation and Optimization*, IEEE Computer Society: Washington, DC, USA, 2004; 279.
34. Connors DA, Hwu WmW. Compiler-Directed Dynamic Computation Reuse: Rationale and Initial Results. *Proceedings of the 32nd Annual ACM/IEEE International Symposium on Microarchitecture*, IEEE Computer Society: Washington, DC, USA, 1999; 158–169.
35. Hall M, Paul M. Improving Software Performance with Automated Memoization. *The Johns Hopkins APL Technical Digest* 1997; **18**:254–260.
36. Alvarez C, Corbal J, Valero M. Fuzzy Memoization for Floating-Point Multimedia Applications. *IEEE Transactions on Computers* 2005; **54**(7):922–927.
37. Hyde R. *The Art of Assembly Language*. No Starch Press: San Francisco, CA, USA, 2003.
38. Zhang Y, Deng L, Yedlapalli P, Muralidhara S, Zhao H, Kandemir M, Chakrabarti C, Pitsianis N, Sun X. A Special-Purpose Compiler for Look-up Table and Code Generation for Function Evaluation. *Design, Automation Test in Europe Conference Exhibition (DATE), 2010*, 2010; 1130–1135.
39. Defour D. Cache-optimised methods for the evaluation of elementary functions. *Technical Report 2002-38*, Ecole normale superieure de Lyon, Lyon, France 2002.