# Algorithmic C (AC) Math Library Reference Manual

Software Version v2.0.10

August 2018

# Table of Contents

Table of Contents

Table of Contents

# Section 1.  Introduction

The Algorithmic C Math Library (*ac_math*) contains synthesizable C++ functions commonly used in Digital Signal Processing applications. The functions use the Algorithmic C data types and are meant to serve as examples on how to write parameterized models and to facilitate migrating an algorithm from using floating-point to fixed-point arithmetic where the math functions either need to be computed dynamically or via lookup tables or piecewise linear approximations.

The input and output arguments of the math functions are parameterized so that arithmetic may be performed at the desired fixed point precision and provide a high degree of flexibility on the area/performance trade-off of hardware implementations obtained during Catapult synthesis.

The hardware implementations produced by Catapult on the math functions are bit accurate. Simulation of the RTL can thus be easily compared to the C++ simulation of the algorithm.  The following sections provide a summary of the *ac_math* library:

- Summary of Functions
- Types of Approximations
- Installing the ac_math library

## 1.1.1. Using the ac_math library

In order to utilize any of the math functions, add the following include line to the source:

```
#include <ac_math.h>
```

## 1.2.  Summary of Functions

The following sections summarize the functions and classes currently supported in the ac_math library. A discussion of the approximation methods follows.

## 1.2.1. Basic Math Functions

| Function Type | Function Call | Approximation Method | Supported Data Types | | |
|---|---|---|---|---|---|
| | | | **ac_fixed** | **ac_float** | **ac_complex** |
| **Absolute Value** | *ac_abs()* | N/A | Yes | Yes | No |
| **Division** | *ac_div()* | N/A | Yes | Yes | Yes |
| **Normalization** | *ac_normalize()* | N/A | Yes | No | Yes |
| **Reciprocal** | *ac_reciprocal_pwl()* | PWL | Yes | Yes | Yes |
| **Logarithm Base e** | *ac_log_pwl()* | PWL | Yes | No | No |
| | *ac_log_cordic()* | CORDIC | Yes | No | No |
| **Logarithm Base 2** | *ac_log2_pwl()* | PWL | Yes | No | No |

| Function Type | Function Call | Approximation Method | Supported Data Types | | |
|---|---|---|---|---|---|
| | | | ac_fixed | ac_float | ac_complex |
| | *ac_log2_cordic()* | CORDIC | Yes | No | No |
| **Exponent Base e** | *ac_exp_pwl()* | PWL | Yes | No | No |
| | *ac_exp_cordic()* | CORDIC | Yes | No | No |
| **Exponent Base 2** | *ac_pow2_pwl()* | PWL | Yes | No | No |
| | *ac_exp2_cordic()* | CORDIC | Yes | No | No |
| **Square Root** | *ac_sqrt_pwl()* | PWL | Yes | Yes | Yes |
| | *ac_sqrt()* | N/A | Yes | No | No |
| **Inverse Square Root** | *ac_inverse_sqrt_pwl()* | PWL | Yes | Yes | Yes |
| **Sine/Cosine** | *ac_sincos()* | LUT | Yes | No | N/A |
| | *ac_cos_cordic()* | CORDIC | Yes | No | N/A |
| | *ac_sin_cordic()* | CORDIC | Yes | No | N/A |
| | *ac_sincos_cordic()* | CORDIC | Yes | No | N/A |
| **Tangent** | *ac_tan_pwl()* | PWL | Yes | No | N/A |
| **Inverse Trig** | *ac_atan_pwl()* | PWL | Yes | No | N/A |
| | *ac_arccos_cordic()* | CORDIC | Yes | No | N/A |
| | *ac_arcsin_cordic()* | CORDIC | Yes | No | N/A |
| | *ac_arctab_cordic()* | CORDIC | Yes | No | N/A |
| **Shift Left/Right** | *ac_shift_left* | N/A | Yes | No | Yes |
| | *ac_shift_right* | N/A | Yes | No | Yes |

# 1.2.2. AC Matrix Class

The class ac_matrix implements a 2-D container class with a template parameter to specify the data type of the internal storage.

The class has member functions to implement some common operations including

- Assignment: operator=()

- Read-Only and Read-Write Element Access: *this(<row>,<col>)

- Comparison: operator!=(), operator==()

- Piecewise Addition: operator+(), operator+=()

- Piecewise Subtraction: operator-(), operator-=()

- Piecewise Multiplication: pwisemult()

- Matrix Multiplication (nested loops): operator*()

- Matrix Transpose: transpose()

- Sum All Elements: sum()

- Scale All Elements: scale(value)

- Formatted Stream Output: ostream &operator<<()

When using the computational functions with AC Datatypes, the form that returns a value is designed in such a way as to determine the full precision required in the output type in order to preserve accuracy during the operation. So using operator+ between two 10 bit ac_fixed matrices will return an 11 bit ac_fixed matrix. If you wish to prevent the bit growth and accept the truncation, you can use the compound operators +=,-=, etc. so that the target object receives the truncated values.

In addition to the built-in member functions, the ac_math library also includes stand-alone functions for more complicated linear algebra operations as described in the next section.

# 1.2.3. Linear Algebra Functions

The ac_math library includes several linear algebra functions that operate on either ac_matrix or plain C-style arrays. These functions, when used with AC Datatypes, are designed to give the user greater control over the bit precision of internal variables and the return value.

- Matrix Multiplication

- Matrix Determinant

- Cholesky Decomposition

- Cholesky Inverse

- QR Decomposition

# 1.3.  Types of Approximations

The following sections discuss the calculation methods used by the functions in the ac_math library and the trade-offs of using one method over another in your design.

- Piecewise Linear

- Lookup Table (LUT)

- CORDIC

- Miscellaneous Math Functions

# 1.3.1. Piecewise Linear

Some of the functions available in the ac_math library are implemented as piecewise linear (PWL) approximations.

A PWL approximation of a function essentially attempts to replicate a function as a set of line segments that are the closest fit to the actual function curve.

**Illustration 1: Reciprocal PWL approximation**

An example of a PWL approximation would be the reciprocal function, for the interval of [0.5, 1). A graphical depiction is given above. The interval is divided into 6 segments. As is evident from the graph and the close nature of the fit, where the red line (PWL output) is virtually indistinguishable from the green one (Accurate output), the approximate output is very close to the accurate output.

## Constructing the PWL Lookup Table

In order to construct the PWL approximation, we need information on the accurate function output for the end points of each segment. Once that information is obtained, we calculated the slope and intercept value for that segment based on the input and output values corresponding to the segment end points. The segments are then shifted slightly in the direction opposite to the direction of concavity of the function in the interval of a particular segment, in order to further optimize the fit obtained. This is done by appropriately changing the intercept values for that segment. The slope and intercept values are then stored in separate lookup tables.

The errors obtained are generally small and tolerable for approximate applications. For instance, the reciprocal PWL approximation just mentioned has a maximum absolute error of 0.005511 over the interval of [0.5, 1), which corresponds to at least 7 error-free fractional bits in case of a fixed-point PWL output.

## Pitfalls

However, it is important to note that if PWL function calls are cascaded, this error can build up and exceed tolerances. Such a situation is encountered, for instance, if the PWL functions are used for linear algebra functions that operate on matrices. Hence, the user must be aware of this pitfall and take steps to avoid it if necessary.

## C++ Compiler

The PWL functions use default template arguments. In order to use a C++ compiler that supports this functionality, the user must use C++11 as the standard for their compilation, or a later standard, failing which a compile-time error is thrown.

## Rounding Mode for PWL Output

The internal variable to store the PWL output for all the functions is set to have rounding turned off by default (*AC_TRN*). This is because the default bitwidth of the PWL output operates at full precision, and no

rounding is required. If, however, the user wishes to reduce the number of bits and use another rounding mode, they can pass it as a template argument. For an example on how to pass this argument for all the PWL functions that are implemented, please refer to the "Function Prototypes" and "Example Function Calls" sections for the various functions as explained later in the documentation.

## 1.3.2. Lookup Table (LUT)

Some of the functions can be efficiently implemented as a lookup table. For example, the sine and cosine functions can be defined as a lookup table where the symmetry of the functions can be exploited to minimize the size of the table. In most cases the table is described using literal values expressed as C++ double precision values. The context in which the function is then used (determined by the fixed-point output precision) will determine the amount of error in the function.

## 1.3.3. CORDIC

Some of the hyperbolic and trigonometric functions have implementations based on the CORDIC algorithm. These implementations use an iterative approach that typically converges with one digit per iteration. The iterations may involve addition, subtraction, bit shifts and table lookups. Given the iterative nature of these implementations the resulting hardware will be larger and/or slower than the PWL or LUT implementations but offer the best accuracy.

## 1.3.4. Miscellaneous Math Functions

A number of math functions that do not fall in the categories listed above are also provided. These are the absolute value, division, square root and shifting functions. These functions give an exact or very accurate output.

## 1.4. Installing the ac_math library

The library consists of three directories, shown here:

```
.
|-- include
|   |-- ac_math
|   |   |-- ac_inverse_sqrt_pwl.h
|   |   |-- ac_log_pwl.h
|   |   |-- ac_math_ns.h
|   |   |-- ac_normalize.h
|   |   |-- ac_pow_pwl.h
|   |   |-- ac_reciprocal_pwl.h
|   |   |-- ac_shift.h
|   |   |-- ac_sincos_lut.h
|   |   `-- ac_sqrt_pwl.h
|   `-- ac_math.h
```

```
|-- pdfdocs
|   `-- ac_math_ref.pdf
`-- tests
    |-- Makefile
    |-- rtest_ac_inverse_sqrt_pwl.cpp
    |-- rtest_ac_log2_pwl.cpp
    |-- rtest_ac_log_pwl.cpp
    |-- rtest_ac_normalize.cpp
    |-- rtest_ac_pow_pwl.cpp
    |-- rtest_ac_reciprocal_pwl.cpp
    |-- rtest_ac_sincos_lut.cpp
    `-- rtest_ac_sqrt_pwl.cpp
```

In order to utilize this library you must have the AC Datatypes package installed and configure your software environment to provide the path to the AC Datatypes "include" directory as part of your C++ compilation arguments.

## Testing and Error Calculation

The ac_math library includes a series of unit tests that exercise each of the approximation functions across various fixed-point bit widths to ensure that the accuracy of the approximation is within a certain tolerance of the standard C++ math library equivalent (under the same input and output bit-width constraints). To exercise these tests from a Linux shell, use the following GNU make command line:

```
gmake all AC_TYPES_INC=<path to AC Datatypes include directory>
```

where the variable AC_TYPES_INC specifies the path to the install location of the AC Datatypes package. The results of the tests look something like this:

```
 TEST: ac_inverse_sqrt_pwl()   INPUT: ac_float< 5, 3, 3, AC_RND>  OUTPUT:
ac_float<64,32,10, AC_RND>    RESULT: PASSED , max error (0.083916)
 TEST: ac_inverse_sqrt_pwl()   INPUT: ac_float< 5, 1, 3, AC_RND>  OUTPUT:
ac_float<64,32,10, AC_RND>    RESULT: PASSED , max error (0.083916)
 TEST: ac_inverse_sqrt_pwl()   INPUT: ac_float< 5, 0, 3, AC_RND>  OUTPUT:
ac_float<64,32,10, AC_RND>    RESULT: PASSED , max error (0.083916)
 TEST: ac_inverse_sqrt_pwl()   INPUT: ac_float< 5,-2, 3, AC_RND>  OUTPUT:
ac_float<64,32,10, AC_RND>    RESULT: PASSED , max error (0.083916)
 TEST: ac_inverse_sqrt_pwl()   INPUT: ac_float< 5, 9, 3, AC_RND>  OUTPUT:
ac_float<60,30,11, AC_RND>    RESULT: PASSED , max error (0.083916)
```

The output shows the function under test, the input and output bit-widths and the maximum error observed over the tested range of minimum and maximum values expressible in the input type stepping by the smallest value expressible in the input type.

In addition to the PWL approximations, the LUT and CORDIC implementations are also subject to error calculations on their output, in order to evaluate whether the error is tolerable.

For comparison, the input to the function being tested is converted to a double value (or ac_complex<double>, in case of complex inputs), and this double value is passed to the equivalent C++

math library function. The output, which is a double (or ac_complex<double>, in case of complex outputs), is subject to quantization according the output type of the function being tested.

It is this quantized output that is compared to the output of the function being tested. The comparison is done either in terms of relative error or absolute error. For real outputs, the relative/absolute error between the accurate  (C++ math library) output and approximate output is calculated, based upon whether the accurate output lies above or below a pre-defined threshold; if the accurate output lies above the threshold, the relative error is calculated, if it lies below the threshold, the absolute error is calculated. For complex outputs, the error is calculated in a manner similar to the real output, with the differences being that (a) the quantization is carried out on the real and imaginary parts of the accurate output with respect to the type of the real and imaginary parts of the output of the function being tested and (b) the magnitude of the relative/absolute error is what is reported. The above approach for comparison is followed for all the ac_math functions except for the:

- *ac_log_pwl*, *ac_log2_pwl*, *ac_atan_pwl* and *ac_sincos_lut* functions where only the absolute error is calculated.

- *ac_abs*, *ac_shift_left*, *ac_shift_right* and *ac_normalize* functions, which return an exact value at the output.

# Section 2. Piecewise Linear Functions

The *ac_math* package includes the following piecewise linear functions:

- Logarithm (ac_log_pwl)
- Power (ac_pow_pwl)
- Reciprocal (ac_reciprocal_pwl)
- Square Root (ac_sqrt_pwl)
- 2.5Inverse Square Root (ac_inverse_sqrt_pwl)
- Tangent (ac_tan_pwl)
- Arctangent (ac_atan_pwl)
- Sigmoid (ac_sigmoid_pwl)

Every function normalizes the function input to a value that is within the domain of the PWL approximation. The PWL approximation produces an output for this normalized value, and then applies a sort of "denormalization" that cancels out the effect of the previous normalization on the PWL output. This gives us an approximate output for the input that was supplied to the function. The domain for each PWL approximation, number of segments used, the maximum absolute error for the output in that domain, and the minimum number of error-free fractional bits for fixed-point PWL outputs is given in the table below:

| PWL Function | Domain | Segments | Max. Abs. Error | Min. No. of Error-free Fractional Bits |
|---|---|---|---|---|
| Logarithm | [0.5, 1) | 8 | 0.001251 | 9 |
| Power | [0, 1) | 4 | 0.003718 | 8 |
| Reciprocal | [0.5, 1) | 8 | 0.003274 | 8 |
| Square Root | [0.5, 1) | 4 | 0.000582 | 10 |
| Inverse Square Root | [0.5, 1) | 8 | 0.000893 | 10 |
| Tangent | [0, π/4) | 8 | 0.002300 | 8 |
| Arctangent | [0, 1) | 4 | 0.002885 | 8 |
| Sigmoid | [0, 5) | 8 | 0.002755 | 10 |

The following subsections describes the implementation and usage of these functions in more detail.

## 2.1. Logarithm (ac_log_pwl)

The *ac_log_pwl* library provides a piecewise linear implementation for the base 2 and base e logarithmic functions, optimized to provide high performance with quick results. This function is implemented for *ac_fixed* datatypes. At a minimal cost to accuracy, it is observed that it the function is faster than other

algorithms and consumes significantly lesser area. These qualities make it useful as a basic building block in high speed IP.

# 2.1.1. The ac_log2_pwl Implementation

The header file provides piecewise linear implementation of the base two logarithm via the *ac_log2_pwl* function. The natural logarithm is computed by using the change of base property.

## Handling Zero Input

A macro-enabled AC_ASSERT is provided to alert the user if a zero input is passed to the PWL function. In addition, functionality is also provided to ensure that the ac_log_pwl function passes the minimum negative value that can be represented with the output type when a zero input is encountered.

## PWL Approximation Graph

To explain the closeness of the PWL approximation to the actual, accurate implementation of the base 2 logarithm function, the following graph compares the PWL output against the accurate function output, where the red line (PWL output) is virtually indistinguishable from the green one (Accurate output):



**Illustration 2: PWL Output vs. Accurate Output for Base2 Logarithm**

## 2.1.1. The ac_log_pwl Implementation

In order to calculate the natural logarithm of *ac_fixed* inputs, we rely upon the *ac_log2_pwl* implementation. As mentioned earlier, the change of base property is used. It can be represented by the following equation

```
log(x) = log2(x) * log(2)
```

The output of the *ac_log2_pwl* implementation is multiplied by log(2), a constant value. The product is the output of the *ac_log_pwl* function.

## 2.1.2. Function Prototypes

The following code shows the template prototypes for the different implementations:

```
template<ac_q_mode q_mode_temp = AC_TRN,
         int W1, int I1, ac_q_mode q_mode_in,  ac_o_mode o_mode_in,
         int W2, int I2, bool S2, ac_q_mode q_mode_out, ac_o_mode q_mode_out>
void ac_log2_pwl(
  const ac_fixed<W1, I1, false, q_mode_in, o_mode_in> input,
  ac_fixed<W2, I2, S2, q_mode_out, o_mode_out> &result
)
```

```
template<ac_q_mode q_mode_temp = AC_TRN,
         int W1, int I1, ac_q_mode q_mode_in,  ac_o_mode o_mode_in,
         int W2, int I2, bool S2, ac_q_mode q_mode_out, ac_o_mode q_mode_out>
void ac_log_pwl(
  const ac_fixed<W1, I1, false, q_mode_in, o_mode_in> input,
  ac_fixed<W2, I2, S2, q_mode_out, o_mode_out> &result
)
```

Note that the function only accepts unsigned inputs, which corresponds to the behavior of the actual logarithm function, which only accepts positive values for its domain. This is different than the behavior of the C math library logarithm functions, which use signed double datatypes for their inputs.

### Returning by Value

The *ac_log2_pwl* and *ac_log_pwl* functions can return their output by value as well as by reference. In order to return the output by value, the user must pass the information of the output type to the function as a template argument. For an example of how to do this, please refer to the Example Function Calls section below. The prototypes for the function call to return by value for both functions is shown below:

```
template<class T_out,
         ac_q_mode q_mode_temp = AC_TRN,
         class T_in>
T_out ac_log2_pwl(
  const T_in &input
)
```

```
template<class T_out,
         ac_q_mode q_mode_temp = AC_TRN,
         class T_in>
T_out ac_log_pwl(
    const T_in &input
)
```

## 2.1.3. Example Function Calls

An example of a function call to store the value of the base 2 and base e logarithms of a sample *ac_fixed* variable *x* in the variables y_log2 and y_log is shown below:

```
ac_fixed<20, 11, false, AC_RND, AC_SAT> x = 2.875;
typedef ac_fixed<30, 15,  true, AC_RND, AC_SAT> output_type;
output_type y_log, y_log2;

ac_log2_pwl(x, y_log2); // Approximates y_log2 = log2(x)
ac_log_pwl(x, y_log);   // Approximates y_log  = log(x)
```

### Changing Rounding/Truncation Mode for the PWL Output Variable

As mentioned in the introduction, the temporary variable that stores the PWL output has rounding turned off by default. An example in which the user changes the rounding/truncation mode is shown as follows:

```
ac_log2_pwl<AC_RND>(x, y);
ac_log_pwl<AC_RND>(x, y);
```

### Returning by Value

In order to have the function return by value, the user must pass the output type information as a template parameter. This is done as follows, for the base 2 and base e logarithmic functions:

```
y = ac_log2_pwl<output_type>(x);
y = ac_log_pwl<output_type>(x);
```

If the user wants to change the default template parameters, e. g. they want to round the output of the PWL implementation, they can do so by using the following function calls as guidelines:

```
y = ac_log2_pwl<output_type, AC_RND>(x);
y = ac_log_pwl<output_type, AC_RND>(x);
```

# 2.2.  Power (ac_pow_pwl)

The *ac_pow_pwl* function is designed to provide a quick approximation of the base 2 and natural exponentials of real numbers using a piecewise linear (PWL) implementation with 5 points/4 segments, with high accuracy. This frees us of the burden of having to calculate a more accurate output using methods such as Taylor series expansion, which requires loop unrolling, a large number of adders and hence a comparatively large area, to give 100% throughput. By contrast, the ac_pow_pwl function can give 100% throughput by merely pipelining it with an II of 1, which results in a lower area for the hardware.

## 2.2.1. The ac_pow2_pwl Implementation

The *ac_pow_pwl* library provides two functions, one which calculates the base 2 exponential, and the other which calculates the natural exponential. The natural exponential version (henceforth called the *ac_exp_pwl* implementation) depends upon the base 2 exponential version (henceforth called the *ac_pow2_pwl* implementation) for its computation, as will be explained later. Both functions only accept *real*, *ac_fixed* inputs and calculate *ac_fixed* outputs.

### PWL Approximation Graph

---

To explain the closeness of the PWL approximation to the actual, accurate implementation of the base 2 exponential function, the following graph compares the PWL output against the accurate function output:



**Illustration 3: PWL Output vs. Accurate Output for Base 2 Exponential**

# 2.2.2. The ac_exp_pwl Implementation

In order to calculate the natural exponential of ac_fixed inputs, we rely upon the *ac_pow2_pwl* implementation. The following relation is used:

```
exp(x) = 2 ^ (x * log2(e))
```

Hence, all that needs to be done is to multiply the input by *log2(e)*, store the product in a temporary variable, and then pass that to the *ac_pow2_pwl* implementation. While doing so, it is ensured that the temporary variable has enough fractional bits to represent the multiplication result. This can be required when the user uses an input *ac_fixed* type with an insufficient number of fractional bits.

## Setting the Minimum Number of Fractional Bits

In order to change the minimum number of fractional bits that are used to store the product of *the input* and *log2(e)*, the user can pass a template value with the minimum they desire. By default this value is set to 9, because that minimum value was empirically seen to provide the best accuracy for inputs that did not have enough fractional bits. For an example on how to pass this argument, please refer to the *Example Function Calls* section below.

# 2.2.3. Function Prototypes

The following code shows the template prototypes for the different implementations:

```
template<ac_q_mode pwl_Q = AC_TRN,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, ac_q_mode outQ, ac_o_mode outO>
void ac_pow2_pwl(
  const ac_fixed<W, I, S, Q, O> &input,
  ac_fixed<outW, outI, false, outQ, outO> &output
)
```

```
template<int n_f_b = 9, ac_q_mode pwl_Q = AC_TRN,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, ac_q_mode outQ, ac_o_mode outO>
void ac_exp_pwl(
  const ac_fixed<W, I, S, Q, O> &input,
  ac_fixed<outW, outI, false, outQ, outO> &output
)
```

Note that the function only accepts unsigned ac_fixed datatypes for the output. This corresponds to the behavior of the actual exponential functions, which have a range that covers all positive real values. The behavior of our PWL functions in this sense is different than the behavior of the standard C math library functions, which accept signed doubles for both the input and the output.

## Returning by Value

The *ac_pow2_pwl* and *ac_exp_pwl* functions can return their output by value as well as by reference. In order to return the output by value, the user must pass the information of the output type to the function as a template argument. For an example of how to do this, please refer to the *Example Function Calls* section below. The prototypes for the function call to return by value for both functions is shown below:

```
template<class T_out,
         ac_q_mode pwl_Q = AC_TRN,
         class T_in>
T_out ac_pow2_pwl(
    const T_in &input
)
```

```
template<class T_out,
         int n_f_b = 9,
         ac_q_mode pwl_Q = AC_TRN,
         class T_in>
T_out ac_exp_pwl(
    const T_in &input
)
```

# 2.2.4. Example Function Calls

An example of a function call to store the value of the natural exponential of a sample *ac_fixed* variable *x* in a variable *y* is shown below:

```
typedef ac_fixed<20, 11, true, AC_RND, AC_SAT> input_type;
typedef ac_fixed<30, 15, true, AC_RND, AC_SAT> output_type;

input_type x = 2.875;
output_type y_exp, y_pow2;
```

```
ac_pow2_pwl(x, y_pow2); //Approximates y_pow2 = pow(2, x.to_double())
ac_exp_pwl(x, y_exp);   //Approximates y_exp = exp(x)
```

The variable *y* hereafter contains the approximate value of the natural exponential of *x*.

## Changing the Minimum Number of Fractional Bits

As mentioned earlier, the minimum number of fractional bits for the result of the multiplication of the input to the natural exponent functions and *log2(e)* can be varied. An example of this is shown below, where the user allocates a minimum of 7 fractional bits.

```
ac_exp_pwl<7>(x, y);
```

## Changing Rounding/Truncation Mode for the PWL Output Variable

As mentioned in the introduction, the temporary variable that stores the PWL output has rounding turned off by default. An example in which the user changes the rounding/truncation mode is shown as follows:

```
ac_pow2_pwl<AC_RND>(x, y);
ac_exp_pwl<9, AC_RND>(x, y);
```

The latter function call changes the rounding/truncation mode while keeping the minimum number of fractional bits same as the default, i.e. 9.

## Returning by Value

In order to have the function return by value, the user must pass the output type information as a template parameter. This is done as follows, for the base 2 and base e exponential functions:

```
y = ac_pow2_pwl<output_type>(x);
y = ac_exp_pwl<output_type>(x);
```

If the user wants to change the default template parameters, e.g. they want to truncate the output of the PWL implementation and/or change the minimum no. of fractional bits to, say, 7, they can do so by using the following function calls as guidelines:

```
y = ac_pow2_pwl<output_type, AC_RND>(x);
y = ac_exp_pwl<output_type, 7, AC_RND>(x);
```

# 2.3.  Reciprocal (ac_reciprocal_pwl)

The *ac_reciprocal_pwl* function is designed to provide a quick approximation of the reciprocal of real and complex numbers using a Piecewise Linear (PWL) implementation with 9 points/8 segments, optimized to give a high-accuracy implementation. Many hardware division operations are calculated indirectly by first obtaining the value of the reciprocal of the denominator, and then multiplying that with the numerator. The calculation of the reciprocal can be done using PWL approximation, which is faster and requires lesser area than actual division hardware.

## 2.3.1. The ac_reciprocal_pwl Implementation

The *ac_reciprocal_pwl* library provides four overloaded functions for the calculation of the reciprocal of real and complex inputs. Each overloaded function handles a different input datatype. The four datatypes hence handled are (a) *ac_fixed*, (b) *ac_float*, (c) *ac_complex<ac_fixed>* and (d) *ac_complex<ac_float>*. It is the *ac_fixed* function that actually contains the code required for the PWL implementation. All the other functions rely upon the *ac_fixed* PWL implementation.

## Handling Zero Input

The *ac_reciprocal_pwl* library provides a macro-enabled AC_ASSERT which will produce a run-time assert when a zero input is encountered. If this assert fails to kick in, additional functionality is provided that ensures output saturation when a zero input is encountered.

## PWL Approximation Graph

To explain the closeness of the PWL approximation to the "real thing", i.e. the actual, accurate implementation of the reciprocal function, the following graph compares the PWL output against the accurate function output:



**Illustration 4: PWL Output vs. Accurate Output for Reciprocal**

# 2.3.2. Function Templates

The following are the overloaded function prototypes for the *ac_reciprocal_pwl* function for different datatypes:

```
template<ac_q_mode pwl_Q = AC_TRN,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, bool outS, ac_q_mode outQ, ac_o_mode outO>
void ac_reciprocal_pwl(
  const ac_fixed<W, I, S, Q, O> &input,
  ac_fixed<outW, outI, outS, outQ, outO> &output
);
```

```
template<ac_q_mode pwl_Q = AC_TRN,
         int W, int I, int E, ac_q_mode Q,
         int outW, int outI, int outE, ac_q_mode outQ>
void ac_reciprocal_pwl(
  const ac_float<W, I, E, Q> &input,
  ac_float<outW, outI, outE, outQ> &output
);
```

```
template<ac_q_mode pwl_Q = AC_TRN,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, bool outS, ac_q_mode outQ, ac_o_mode outO>
void ac_reciprocal_pwl(
  const ac_complex<ac_fixed<W, I, S, Q, O> > &input,
  ac_complex<ac_fixed<outW, outI, outS, outQ, outO> > &output
);
```

```
template<ac_q_mode pwl_Q = AC_TRN,
         int W, int I, int E, ac_q_mode Q,
         int outW, int outI, int outE, ac_q_mode outQ>
void ac_reciprocal_pwl(
  const ac_complex<ac_float<W, I, E, Q> > &input,
  ac_complex<ac_float<outW, outI, outE, outQ> > &output
);
```

### Returning by Value

The *ac_reciprocal_pwl* functions can return their output by value as well as by reference. In order to return the value, the user must pass the information of the output type to the function as a template argument. For an example of how to do this, please refer to the *Example Function Calls* section below. The prototype for the function call to return by value is shown below:

```
template<class T_out,
         ac_q_mode pwl_Q = AC_TRN,
         class T_in>
T_out ac_reciprocal_pwl(
  const T_in &input
);
```

## 2.3.3. Example Function Calls

An example of a function call to store the value of the reciprocal of a sample *ac_fixed* variable *x* in a variable *y* is shown below:

```
typedef ac_fixed<20, 11, true, AC_RND, AC_SAT> input_type;
typedef ac_fixed<30, 15, true, AC_RND, AC_SAT> output_type;
input_type x = -1.75;
output_type y;
ac_reciprocal_pwl(x, y);  //Approximates y = 1 / x
```

The variable y hereafter stores the approximate value of the reciprocal of *x*.

## Changing Rounding/Truncation Mode for PWL Output Variable

As mentioned in the introduction, the temporary variable that stores the PWL output has rounding turned off by default. An example of passing *AC_RND* as the rounding/truncation mode instead, for the same data-types of x and y as in the previous example, is shown below:

```
ac_reciprocal_pwl<AC_RND>(x, y);
```

## Returning by Value

As mentioned earlier, the *ac_reciprocal_pwl* functions can also return by value. In order to do so, the type information for the output must be passed explicitly to function as shown below.

```
y = ac_reciprocal_pwl<output_type>(x);
```

If the user also wishes to change the rounding mode for the temporary variable and have the function return by value, they can call the function as follows:

```
y = ac_reciprocal_pwl<output_type, AC_RND>(x);
```

# 2.4. Square Root (ac_sqrt_pwl)

The *ac_sqrt_pwl* function is a piecewise linear implementation of a square root function that has been optimized to provide a fast, low-area implementation with minimal error, making it useful as a basic building block in high speed IP blocks. This function is implemented as an overloaded function for *ac_fixed*, *ac_float and ac_complex<ac_fixed> datatypes*.

# 2.4.1. The ac_sqrt_pwl Implementation

The *ac_sqrt_pwl* library provides three overloaded functions for the calculation of the square root of real and complex inputs, each handling a different input datatype. The three datatypes hence handled are (a) *ac_fixed*, (b) *ac_float* and (c) *ac_complex<ac_fixed>*. It is the ac_fixed function that actually contains the code required for the PWL implementation. All the other functions rely upon the *ac_fixed* PWL implementation.

## PWL Approximation Graph

To explain the closeness of the PWL approximation to the "real thing", i.e. the accurate implementation of the square root function, the following graph is provided:
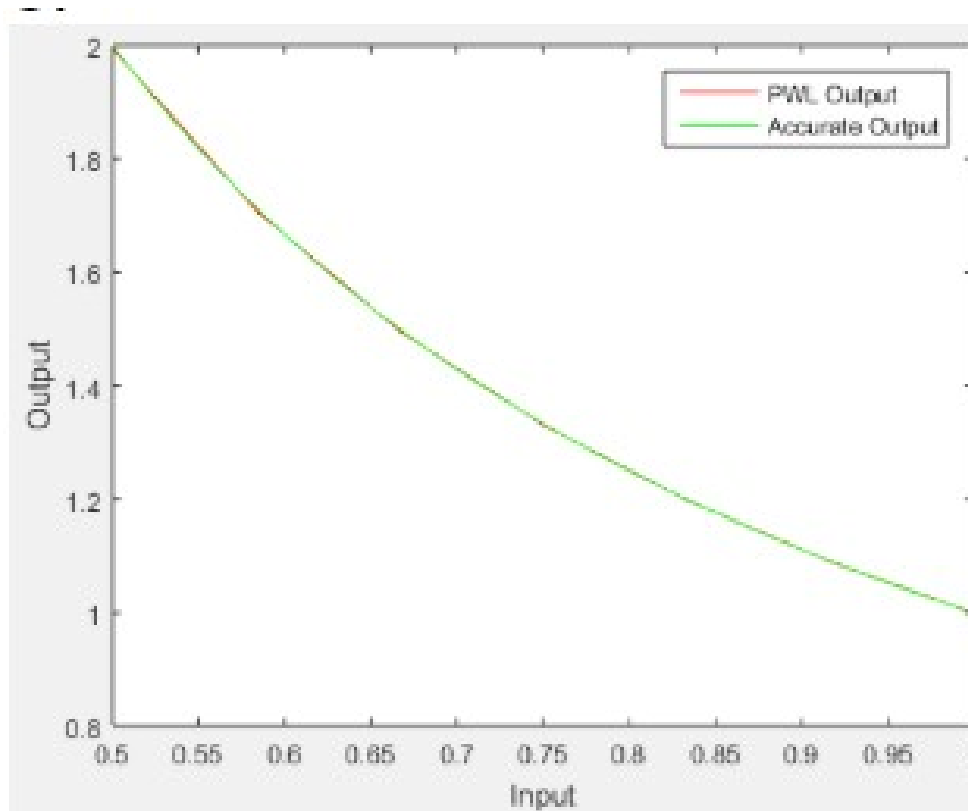
**Illustration 5: PWL Output vs. Accurate Output for Square Root**

# 2.4.2. Function Templates

The following are the overloaded function prototypes for the *ac_sqrt_pwl* function for different datatypes:

```
  template <ac_q_mode pwlQ = AC_TRN,
            int W, int I, ac_q_mode Q, ac_o_mode O,
            int outW, int outI, ac_q_mode outQ, ac_o_mode outO>
 void ac_sqrt_pwl(
   const ac_fixed <W, I, false, Q, O> input,
   ac_fixed <outW, outI, false, outQ, outO> &output
 );
```

```
  template <ac_q_mode pwlQ = AC_TRN,
            int W, int I, int E, ac_q_mode Q,
            int outW, int outI, int outE, ac_q_mode outQ>
 void ac_sqrt_pwl(
   const ac_float <W, I, E, Q> input,
   ac_float <outW, outI, outE, outQ> &output
 );
```

```
  template <ac_q_mode pwlQ = AC_TRN,
            int W, int I, ac_q_mode Q, ac_o_mode O,
            int outW, int outI, ac_q_mode outQ, ac_o_mode outO>
 void ac_sqrt_pwl(
   const ac_complex <ac_fixed <W, I, true, Q, O> > input,
   ac_complex <ac_fixed <outW, outI, true, outQ, outO> > &output
 );
```

### Returning by Value

The *ac_sqrt_pwl* functions can return their output by value as well as by reference. In order to return the value, the user must pass the information of the output type to the function as a template argument. For an example of how to do this, please refer to the "Example function calls" section below. The prototype for the function call to return by value is shown below:

```
template<class T_out, ac_q_mode pwlQ = AC_TRN, class T_in>
T_out ac_sqrt_pwl(
  const T_in &input
);
```

## 2.4.3. Example Function Calls

An example of a function call to store the value of the square root of a sample *ac_fixed* variable *x* in a variable *y*, by using both return by reference and return by value, is shown below. A line of code that illustrates the method to change the rounding mode of the intermediate PWL variable is also included.

```
typedef ac_fixed<20, 11, false, AC_RND, AC_SAT> input_type;
typedef ac_fixed<30, 15, false, AC_RND, AC_SAT> output_type;

input_type x = 1.75;
output_type y;

// Returns y = sqrt(x), and returns by reference.
ac_sqrt_pwl (x, y);

// Change the rounding mode for intermediate PWL variable to AC_RND
ac_sqrt_pwl<AC_RND> (x, y);

// The following line, returns by value instead of by reference.
y = ac_sqrt_pwl <output_type> (x);
```

## 2.5. Inverse Square Root (ac_inverse_sqrt_pwl)

The *ac_inverse_sqrt_pwl* function is a piecewise linear implementation of the inverse square root function ( 1/sqrt(x) ) that has been optimized to provide a high-performance implementation with minimal error, making it useful as a basic building block in high speed IP blocks. This function is implemented as an overloaded function for *ac_fixed, ac_float and ac_complex<ac_fixed> datatypes*. It provides a more accurate and lower-area output as compared to using the square root and reciprocal PWL approximations together to calculate the inverse square root of a number.

## 2.5.1. The ac_inverse_sqrt_pwl Implementation

The *ac_inverse_sqrt_pwl* library provides three overloaded functions for the calculation of the inverse square root of real and complex inputs. Each overloaded function handles a different input datatype. The three datatypes hence handled are (a) *ac_fixed*, (b) *ac_float* and (c) *ac_complex<ac_fixed>*. It is the *ac_fixed* function that actually contains the code required for the PWL implementation. All the other functions rely upon the *ac_fixed* PWL implementation.

## Handling Zero Input

The *ac_inverse_sqrt_pwl* library provides a macro-enabled AC_ASSERT which will produce a run-time assert when a zero input is encountered. If this assert fails to kick in, additional functionality is provided that ensures output saturation when a zero input is encountered.

## PWL Approximation Graph

To explain the closeness of the PWL approximation to the "real thing", i.e. the accurate implementation of the inverse square root function, the following graph is provided:
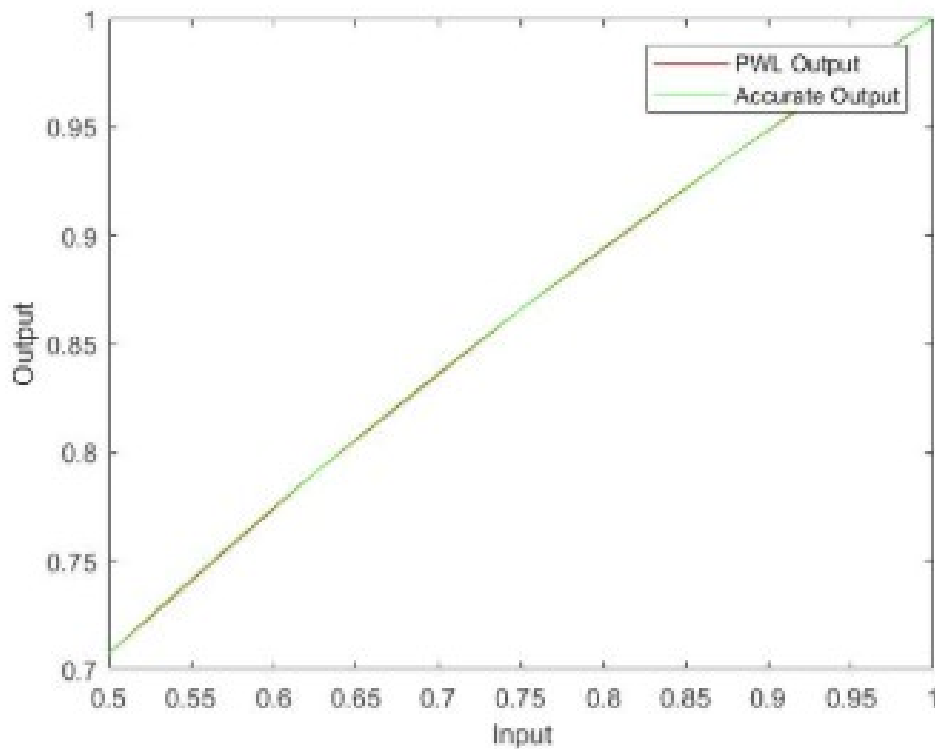


**Illustration 6: PWL Output vs. Accurate Output
for Inverse Square Root**

# 2.5.2. Function Templates

The following are the overloaded function prototypes for the *ac_inverse_sqrt_pwl* function for different datatypes:

```
template<ac_q_mode q_mode_temp = AC_TRN,
         int W1, int I1, ac_q_mode q1, ac_o_mode o1,
         int W2, int I2, ac_q_mode q2, ac_o_mode o2>
void ac_inverse_sqrt_pwl(
  const ac_fixed <W1,I1, false, q1, o1> &input,
  ac_fixed <W2, I2, false, q2, o2> &output
);
```

```
template <ac_q_mode q_mode_temp = AC_TRN,
          int W1, int I1, int E1, ac_q_mode q1,
          int W2, int I2, int E2, ac_q_mode q2>
void ac_inverse_sqrt_pwl(
  const ac_float <W1, I1, E1, q1> &input,
  ac_float <W2, I2, E2, q2> &output
);
```

```
template<ac_q_mode q_mode_temp = AC_TRN,
         int W1, int I1, ac_q_mode q1, ac_o_mode o1,
         int W2, int I2, ac_q_mode q2, ac_o_mode o2>
void ac_inverse_sqrt_pwl(
  const ac_complex <ac_fixed <W1,I1,true, q1, o1> > &input,
  ac_complex <ac_fixed <W2, I2, true, q2, o2> > &output
);
```

### Returning by Value

The *ac_inverse_sqrt_pwl* functions can return their output by value as well as by reference. In order to return the value, the user must pass the information of the output type to the function as a template argument. For an example of how to do this, please refer to the Example Function Calls section below. The prototype for the function call to return by value is shown below:

```
template<class T_out, ac_q_mode q_mode_temp = AC_TRN, class T_in>
T_out ac_inverse_sqrt_pwl(
  const T_in &input
)
```

## 2.5.3. Example Function Calls

An example of a function call to store the value of the inverse square root of a sample *ac_fixed* variable *x* in a variable *y*, by using both return by reference and return by value, is shown below.

```
typedef ac_fixed<20, 11, false, AC_RND, AC_SAT> input_type;
typedef ac_fixed<30, 15, false, AC_RND, AC_SAT> output_type;

input_type x = 1.75;
output_type y;

// Returns y = 1/sqrt(x), and returns by reference.
ac_inverse_sqrt_pwl (x, y);

// Change the rounding mode for intermediate PWL variable to AC_RND
ac_inverse_sqrt_pwl<AC_RND> (x, y);

// The following line, returns by value instead of by reference.
y = ac_inverse_sqrt_pwl <output_type> (x);
```

# 2.6.  Tangent (ac_tan_pwl)

The *ac_tan_pwl* library provides a fast, PWL implementation of the tangent function of first-quadrant angles with minimal inaccuracy. The implementation is faster than other algorithms such as the CORDIC algorithm, which require a significant number of iterations and stored table values to reach an acceptable accuracy in the output.

## Normalization

The actual tangent PWL implementation only covers the domain of [0, π/4). In order to cover a larger range of first-quadrant angle values, i. e. in order to extend the domain of the *ac_tan_pwl* function to [π/4, π/2), the input angle in radians is halved using a shift operation, and the following equation is used:

tan(x) = 2 * tan(x/2) * (1 / (1 – tan(x/2)^2))

In order to calculate the 1 / (1 – tan(x/2)^2) factor in the above equation, we use the *ac_reciprocal_pwl* library for ac_fixed inputs and outputs. Hence, the user must keep in mind that any changes made to that library while also affect the accuracy of the tangent approximation for values that equal or exceed π/4.

## Saturation

It is important to note that the output will saturate to the maximum possible value once the input has crossed the value of 1.556640625 radians, which is roughly 89.2 degrees.

## PWL Approximation Graph

To explain the closeness of the PWL approximation to the "real thing", i.e. the accurate implementation of the tangent function, the following graph is provided for the domain [0, π/4):

**Illustration 7: PWL Output vs. Accurate Output for Tangent**

# 2.6.1. Function Templates

The following is the function prototype for the *ac_tan_pwl* function for ac_fixed datatypes:

```
template<ac_q_mode pwl_Q = AC_TRN,
         int W, int I, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, ac_q_mode outQ, ac_o_mode outO>
void ac_tan_pwl(
  const ac_fixed<W, I, false, Q, O> &input,
  ac_fixed<outW, outI, false, outQ, outO> &output
);
```

## Returning by Value

The *ac_tan_pwl* function can return its output by value as well as by reference. In order to return the value, the user must pass the information of the output type to the function as a template argument. For an example of how to do this, please refer to the Example Function Calls section below. The prototype for the function call to return by value is shown below:

```
template<class T_out, ac_q_mode pwl_Q = AC_TRN, class T_in>
T_out ac_tan_pwl(const T_in &input);
```

# 2.6.2. Example Function Calls

An example of a function call to store the value of the tangent of a sample *ac_fixed* variable *x* in a variable *y*, by using both return by reference and return by value, is shown below.

```
typedef ac_fixed<20, 11, false, AC_RND, AC_SAT> input_type;
typedef ac_fixed<30, 15, false, AC_RND, AC_SAT> output_type;

input_type x = 0.25;
output_type y;
```

```
// Returns y = tan(x), and returns by reference.
ac_tan_pwl (x, y);

// Change the rounding mode for intermediate PWL variable to AC_RND
ac_tan_pwl<AC_RND> (x, y);

// The following line, returns by value instead of by reference.
y = ac_tan_pwl <output_type> (x);
```

# 2.7.  Arctangent (ac_atan_pwl)

The *ac_atan_pwl* library provides a PWL based approximation of the arctangent function for positive *ac_fixed* inputs, with the support of the reciprocal function for any values that exceed unity.

## Normalization

The actual arctangent PWL implementation only covers the domain of [0, 1). This is done not only because the domain from (1, ∞) can be covered with the reciprocal function, but because attempting to cover a larger domain with saturation for input values that cross a certain limit means that firstly, it's harder to obtain an adequate fit to the actual arctangent function, due to the nature of the function curve, necessitating more segments and a larger area to store LUT values, and secondly, a comparator must be used for saturation, resulting in a further increase in area.

The formula used for any values that exceed unity is as follows:

atan(x) =  π/2 – atan(1 / x)

Where the value of 1/x is obtained using the reciprocal PWL function.

An input of unity is handled by storing it in a variable which saturates to a value slightly less than unity, hence ensuring that the normalized input value is still within the domain of the PWL function.

## PWL Approximation Graph

To explain the closeness of the PWL approximation to the "real thing", i.e. the accurate implementation of the arctangent function, the following graph is provided for the domain [0, 1):
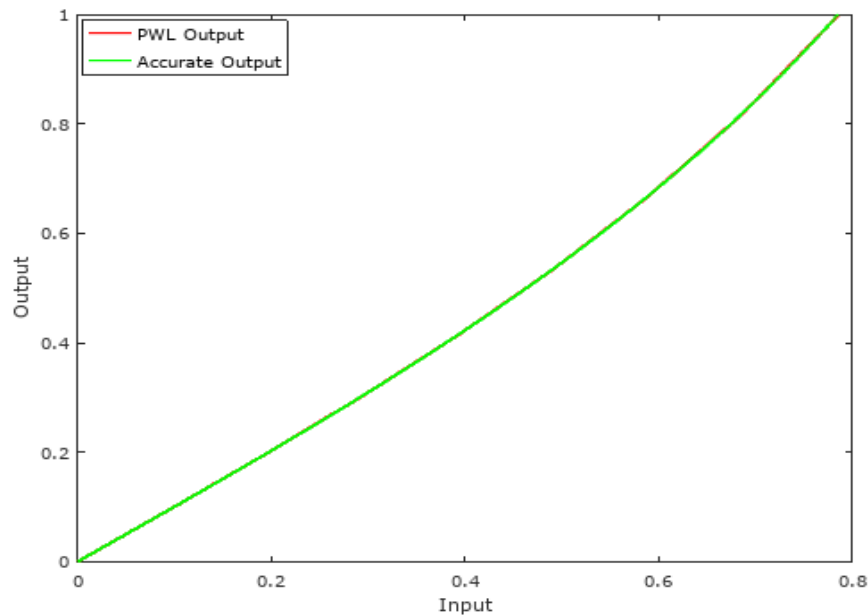
**Illustration 8: PWL Output vs. Accurate Output for Arctangent**

# 2.7.1. Function Templates

The following is the function prototype for the *ac_atan_pwl* function for *ac_fixed* datatypes:

```
template<ac_q_mode pwl_Q = AC_TRN,
         int W, int I, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, ac_q_mode outQ, ac_o_mode outO>
void ac_tan_pwl(
  const ac_fixed<W, I, false, Q, O> &input,
  ac_fixed<outW, outI, false, outQ, outO> &output
);
```

## Returning by Value

The *ac_atan_pwl* function can return its output by value as well as by reference. In order to return the value, the user must pass the information of the output type to the function as a template argument. For an example of how to do this, please refer to the Example Function Calls section below. The prototype for the function call to return by value is shown below:

```
template<class T_out,
         ac_q_mode pwl_Q = AC_TRN,
         class T_in>
T_out ac_atan_pwl(const T_in &input);
```

## 2.7.2. Example Function Calls

An example of a function call to store the value of the arctangent of a sample *ac_fixed* variable *x* in a variable *y*, by using both return by reference and return by value, is shown below.

```
typedef ac_fixed<20, 11, false, AC_RND, AC_SAT> input_type;
typedef ac_fixed<30, 15, false, AC_RND, AC_SAT> output_type;

input_type x = 0.5;
output_type y;

// Returns y = tan(x), and returns by reference.
ac_atan_pwl (x, y);

// Change the rounding mode for intermediate PWL variable to AC_RND
ac_atan_pwl<AC_RND> (x, y);

// The following line returns by value instead of by reference.
y = ac_atan_pwl <output_type> (x);
```

# 2.8.  Sigmoid (ac_sigmoid_pwl)

The *ac_sigmoid_pwl* library provides a PWL based approximation of the sigmoid function for *ac_fixed* inputs with minimal inaccuracy.

### Normalization

The actual sigmoid PWL implementation only covers the domain of [0, 5).  This pwl domain was selected because attempting to cover a larger domain makes it harder to obtain an adequate fit to the actual sigmoid function necessitating more segments and a larger area to store LUT values. As the input approaches infinity, the sigmoid function saturates and thus for inputs greater than 5, the output is saturated to the maximum possible pwl output. For negative inputs, the output will simply be 1 - the output of its positive conterpart taking advantage of the symmetry of the sigmoid function.

### PWL Approximation Graph

To explain the closeness of the PWL approximation to the "real thing", i.e. the accurate implementation of the tangent function, the following graph is provided for the domain [0, 5) :
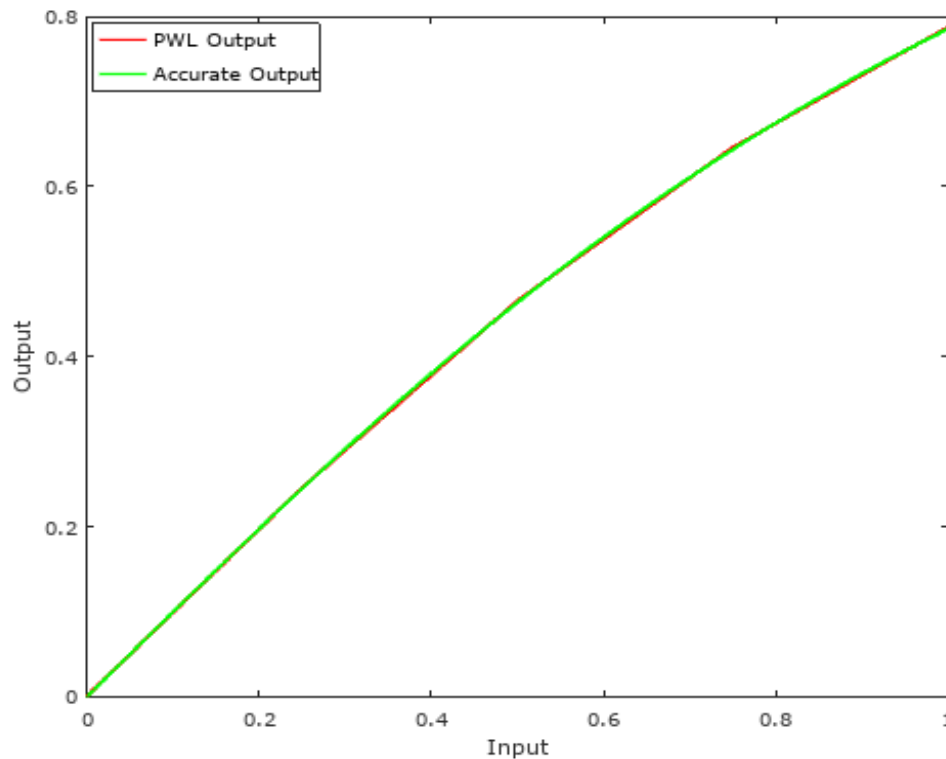
**Illustration 9: PWL Output vs. Accurate Output for Sigmoid**

# 2.8.1. Function Templates

The following is the function prototype for the *ac_sigmoid_pwl* function for *ac_fixed* datatypes:

```
template<ac_q_mode pwl_Q = AC_TRN,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, ac_q_mode outQ, ac_o_mode outO>
void ac_sigmoid_pwl(
  const ac_fixed<W, I, S, Q, O> &input,
  ac_fixed<outW, outI, false, outQ, outO> &output
);
```

## Returning by Value

The *ac_sigmoid_pwl* function can return its output by value as well as by reference. In order to return the value, the user must pass the information of the output type to the function as a template argument. For an example of how to do this, please refer to the 2.8.2Example Function Calls section below. The prototype for the function call to return by value is shown below:

```
template<class T_out,
         ac_q_mode pwl_Q = AC_TRN,
         class T_in>
T_out ac_sigmoid_pwl(const T_in &input);
```

# 2.8.2. Example Function Calls

An example of a function call to store the value of sigmoid of a sample *ac_fixed* variable *x* in a variable *y*, by using both return by reference and return by value, is shown below.

```
typedef ac_fixed<15, 7, true, AC_RND, AC_SAT> input_type;
```

```
typedef ac_fixed<30, 3, false, AC_RND, AC_SAT> output_type;

input_type x = 2.5;
output_type y;

// Returns y = sigmoid(x), and returns by reference.
ac_sigmoid_pwl (x, y);

// Change the rounding mode for intermediate PWL variable to AC_RND
ac_sigmoid_pwl<AC_RND> (x, y);

// The following line returns by value instead of by reference.
y = ac_sigmoid_pwl <output_type> (x);
```

# Section 3.  Cordic Math Functions

The CORDIC-based functions in the *ac_math* library include the following operations:

- Sine/Cosine (ac_sin_cordic/ ac_cos_cordic)

- Arcsin/Arccos (ac_arcsin_cordic/ ac_arccos_cordic)

- Arctangent (ac_atan2_cordic)

- Exponential (ac_exp_cordic/ ac_exp2_cordic)

- Logarithm (ac_log_cordic/ ac_log2_cordic)

- Power (ac_pow_cordic)

## 3.1.  Sine/Cosine (ac_sin_cordic/ ac_cos_cordic)

There are algorithms that require the computation of the *ac_sin_cordic()* and/or ac_*cos_cordic()* functions for dynamic angles (as opposed to a set of predefined angles as in FFTs which are best implemented as a table lookup). In such cases, the CORDIC algorithm is often used as a way to compute the *ac_sin_cordic()* or *ac_cos_cordic()* functions. The math library contains a sample CORDIC implementation using fixed-point data types. The available functional interfaces compute sin, cos or both scaled sin and cos with one call (more efficient hardware than having separate calls to sin and cos).

The first argument is the angle scaled by (*1/PI*). The advantage of the prescaled angle is that it makes the quadrant computation inside the CORDIC algorithm trivial and it also simplifies the complexity of the caller as in most cases the calls to sin and cos are made with angles that are multiples of PI.

AC Math v2.0.10

| **AC: ac_fixed** |
| --- |
| ```void ac_sin_cordic (``` <br> ```  ac_fixed<AW,AI,true,AQ,AO> angle_over_pi,``` <br> ```  ac_fixed<OW,OI,true,OQ,OO> &sin``` <br> ```)``` |
| ```void ac_cos_cordic (``` <br> ```  ac_fixed<AW,AI,true,AQ,AO> angle_over_pi,``` <br> ```  ac_fixed<OW,OI,true,OQ,OO> &cos``` <br> ```)``` |
| ```void ac_sincos_cordic (``` <br> ```  ac_fixed<AW,AI,true,AQ,AO> angle_over_pi,``` <br> ```  ac_fixed<OW,OI,true,OQ,OO> C,``` <br> ```  ac_fixed<OW,OI,true,OQ,OO> &sin,``` <br> ```  ac_fixed<OW,OI,true,OQ,OO> &cos``` <br> ```)``` |
| **SystemC: sc_fixed** |
| ```void ac_sin_cordic (``` <br> ```  sc_fixed<AW,AI,AQ,AO> angle_over_pi,``` <br> ```  sc_fixed<OW,OI,OQ,OO> &sin``` <br> ```)``` |
| ```void ac_cos_cordic (``` <br> ```  sc_fixed<AW,AI,AQ,AO> angle_over_pi,``` <br> ```  sc_fixed<OW,OI,OQ,OO> &cos``` <br> ```)``` |
| ```void ac_sincos_cordic (``` <br> ```  sc_fixed<AW,AI,AQ,AO> angle_over_pi,``` <br> ```  sc_fixed<OW,OI,OQ,OO> C,``` <br> ```  sc_fixed<OW,OI,OQ,OO> &sin,``` <br> ```  sc_fixed<OW,OI,OQ,OO> &cos``` <br> ```)``` |

**Table 1: Functions for Sine, Cosine and Combined Sine and Cosine with Amplitude**

Note that the interface that computes both *sine* and *cosine*, the data types for *C, sin, cos* parameters must be identical.

# 3.2. Arcsin/Arccos (ac_arcsin_cordic/ ac_arccos_cordic)

The *ac_math* library contains sample CORDIC implementation using fixed-point data types. The available functional interfaces compute the inverse trigonometric functions arcsin and arccos.

The result of these functions is scaled by (1/PI), which is often the case in practice, maintaining consistency with the sin/cos implementations.



**Illustration 10: Graph of scaled arcsin and arccos functions**

| AC: ac_fixed |
|---|
| ```
void ac_arccos_cordic(
  ac_fixed<AW,AI,true,AQ,AO> x,
  ac_fixed<OW,OI,false,OQ,OO> &arccos_x_over_pi
);

void ac_arcsin_cordic(
  ac_fixed<AW,AI,true,AQ,AO> x,
  ac_fixed<OW,OI,true,OQ,OO> &arcsin_x_over_pi
);
``` |
| **SystemC: sc_fixed\*** |
| ```
void ac_arccos_cordic(
  sc_fixed<AW,AI,AQ,AO> x,
  sc_ufixed<OW,OI,OQ,OO> &arccos_x_over_pi
``` |

```
  );

  void ac_arcsin_cordic(
    sc_fixed<AW,AI,AQ,AO> x,
    sc_fixed<OW,OI,OQ,OO> &arcsin_x_over_pi
  );
```
*Note: only supports up to 20 fractional bits, i.e., AW - AI <= 20.

# 3.3.  Arctangent (ac_atan2_cordic)

The *ac_atan2_cordic* function is a fixed-point CORDIC implementation of the functionality of the corresponding *ac_math.h* function. It takes two signed fixed-point arguments called *y* and *x* and it returns the arc tangent in the range *-PI* to *PI*. The accuracy of the computation is directly dependent on the precision of the fixed-point *atan* variable passed to the function.

| AC: ac_fixed | SystemC: sc_fixed |
|---|---|
| <pre>void ac_atan2_cordic(<br>  ac_fixed<YW,YI,true,YQ,YO> y,<br>  ac_fixed<XW,XI,true,XQ,XO> x,<br>  ac_fixed<OW,OI,true,OQ,OO> &atan<br>)</pre> | <pre>void ac_atan2_cordic(<br>  sc_fixed<YW,YI,YQ,YO> y,<br>  sc_fixed<XW,XI,XQ,XO> x,<br>  sc_fixed<OW,OI,OQ,OO> &atan<br>)</pre> |

**Table 2: Functions for Atan2**

The behavior when one or both of the arguments is zero is identical with that of the *ac_atan2_cordic* in *ac_math.h*:

- 0 when x == 0 and y == 0

- -PI/2 when x == 0 and y < 0

- +PI/2 when x == 0 and y > 0

# 3.4.  Exponential (ac_exp_cordic/ ac_exp2_cordic)

The *ac_math* library contains sample CORDIC implementations of exponentials. The *exp* function evaluates the exponential of a fixed point argument *x*, i.e., the transcendental number '*e*' raised to the power of *x*.  The *exp2* function evaluates the number 2 raised to the power of a fixed point argument *x*.

The following code shows the function declarations for the AC Numerical Datatypes (*ac_fixed*).  There are no SystemC variants.

```
void ac_exp_cordic( const ac_fixed<AW,AI,AS,AQ,AV> &x,
                    ac_fixed<ZW,ZI,false,ZQ,ZV> &z);

void ac_exp2_cordic( const ac_fixed<AW,AI,AS,AQ,AV> &x,
                     ac_fixed<ZW,ZI,false,ZQ,ZV> &z);
```

# 3.5. Logarithm (ac_log_cordic/ ac_log2_cordic)

The *ac_math* library contains sample CORDIC implementations of logarithms. The 'log' function evaluates the natural logarithm of a fixed point argument x, i.e., the logarithm with transcendental number '*e*' as its base. The log2 function evaluates the base 2 logarithm of argument *x*.

The following code shows the function declarations for the AC Numerical Datatypes (*ac_fixed*). There are no SystemC variants.

```
void ac_log_cordic( const ac_fixed<AW,AI,false,AQ,AV> &x,
                    ac_fixed<ZW,ZI,ZS,ZQ,ZV> &z);

void ac_log2_cordic( const ac_fixed<AW,AI,false,AQ,AV> &x,
                     ac_fixed<ZW,ZI,ZS,ZQ,ZV> &z);
```

# 3.6. Power (ac_pow_cordic)

The *ac_pow_cordic* function evaluates the result of raising a fixed point number *x* to the power of another fixed point number *y*.

The following code shows the function declarations for the AC Numerical Datatypes (*ac_fixed*). There are no SystemC variants.

```
void ac_pow_cordic( const ac_fixed<AW,AI,AS,AQ,AO> &a,
                    const ac_fixed<BW,BI,BS,BQ,BO> &b,
                    ac_fixed<ZW,ZI,false,ZQ,ZV> &z);
```

Note: The logarithm and exponential cordic functions have limitations on the types they can accept for their input/output. The logarithm functions will only accept unsigned inputs, as the actual logarithm functions operate on the positive real-valued domain, while the exponential functions only accept unsigned types for the output, as they output values in the range of positive real values. This is different than the standard C math functions for logarithm and exponential functions, which accept signed doubles at their input and output.

# Section 4.  Lookup Table (LUT) Functions

The *ac_math* package includes the following Lookup Table based functions:

- Sine/Cosine (ac_sincos_lut)

The following subsections describes the implementation and usage of these functions in more detail.

## 4.1.  Sine/Cosine (ac_sincos_lut)

The *ac_sincos_lut* function is designed to provide the sine and cosine values using a lookup table (LUT).

## 4.1.1. The ac_sincos_lut Implementation

The *ac_sincos_lut* function accepts *ac_fixed* datatypes as input. The domain for the input is (0,1) radians/2PI. The function returns an *ac_complex<ac_fixed>* variable where the real part represents the *cosine* value and the imaginary part represents the *sine* value. The number of table entries is 512 values. A prototype of the ac_*sincos_lut* function is shown below:

```
template<class T_in, class T_out>
void ac_sincos_lut(const T_in &input, ac_complex<T_out> &output)
```

In the above prototype, the class *T_in* represents the datatype of the input which should be an *ac_fixed* datatype. The class *T_out* represents the datatype of the output which should also be an *ac_fixed* datatype. Here, *input* represents the input to the *ac_sincos_lut* function, the *output* is an *ac_complex* variable whose real part is the output *cosine* value and whose imaginary part is the output *sine* value.

### Algorithm

All the entries for the sine and cosine lookup tables were generated using the C++ math library sin() and cos() functions and the values are represented as constant doubles. The table contains only those values for the range of angles from (0-PI/4] radians and symmetry is used to adjust the output accordingly. By leveraging the symmetry the table can have more entries for greater precision while still covering the full range of angle inputs. To compute the index into the table, the first 3 bits are extracted from the MSB side of the input to determine in which octant the input lies. The octants can be defined as:

| Input Angle (1/2PI radians) | Octant |
|---|---|
| [ 0.000 – 0.125 ) | 0 |
| [ 0.125 – 0.250 ) | 1 |
| [ 0.250 – 0.375 ) | 2 |
| [ 0.375 – 0.500 ) | 3 |
| [ 0.500 – 0.625 ) | 4 |
| [ 0.625 – 0.750 ) | 5 |
| [ 0.750 – 0.875 ) | 6 |
| [ 0.875 – 1.000 ) | 7 |

The look up table index in this implementation uses the input bits (MSB-3 : LSB). When the input bitwidth is exactly 12 bits, this results in 512 possible indexes. If the input bitwidth is greater than 12 bits, then the closest table entry is found. If the input bitwidth is less than 12 bits, then a stride is implemented.

### Handling negative inputs

If the input angle is a negative value (for example -0.125) it would be represented in 2's complement form as 1.1110 0000. Ignoring the integer portion of the result, the fractional portion is 0.875 which is identical to -0.125.

### Returning by Value

The implementations of the ac_sincos_lut() function can return the output by value as well as by reference. See the sample code below for examples of each form. The prototype for the function to return by value is shown below:

```
template<class T_out, class T_in>
T_out ac_sincos_lut(const T_in &input);
```

# 4.1.2. Example Function Call

An example of a function call is as follows where x represents the input angle and y represents a complex variable whose real part is the cos value and imaginary part is the sin value:

```
typedef ac_fixed<12, 1, true, AC_RND, AC_SAT> input_type;
typedef ac_complex<ac_fixed<23, 1, true, AC_RND, AC_SAT> > output_type;

input_type x = 0.37;
output_type y;

 // returning value by reference
ac_sincos_lut(x, y);
cout << "Sine: " << y.i() << endl;
cout << "Cosine: " << y.r() << endl;

// returning by value
y = ac_sincos_lut<output_type>(x);
```

# 4.1.3. Increasing look up table entries

In the current implementation, the number of lookup table entries for the sine and cosine lookup tables is 512. This means that the implementation gives an accurate output for the cases when input bitwidth is less than or equal to 12 bits. So for cases when input bitwidth is greater than 12 bits, the closest table entry is found (as mentioned above) which might incur some error (maximum error is approximately equal to the largest difference between consecutive lookup table entries) and therefore if a more accurate implementation is desired, then the current lookup table entries can be replaced. The formula for the number of lookup table entries for sine and cosine is as follows:-

Number of lookup table entries = 2 ^ (input bitwidth - 3)

For example, the number of lookup table entries to get an accurate output for the following input bitwidths are:-

- 13 bits – 1024

- 4 bits – 2048

- 15 bits – 4096

The *lutgensincos.cpp* file under the *.../src/examples/Math/ac_sincos_lut* directory can be referred for generating lookup table entries. Also the *ac_sincos_lut.h* library header file has to be accordingly modified if the number of lookup table entries are changed.

## Synthesizing the ac_sincos_lut function

In High Level Synthesis, the *ac_sincos_lut* function can be completely pipelined as there are no data dependencies thus ensuring high throughput and low latency.

# Section 5.  Linear Algebra Functions

The *ac_math* package includes the following linear algebra functions:

- Cholesky Decomposition (ac_chol_d)
- Cholesky Inverse (ac_cholinv)
- Determinant (ac_determinant)
- Matrix Multiplication (ac_matrixmul)
- QR Decomposition (ac_qrd)

The following subsections describes the implementation and usage of these functions in more detail.

## 5.1.  Cholesky Decomposition (ac_chol_d)

The *ac_chol_d* library is designed to provide a Cholesky Decomposition of a square, positive definite input matrix using the Cholesky-Crout algorithm. The user can utilize either accurate math functions or the piece-wise linear (pwl) math library for the internal calculations involved. Cholesky Decomposition is an important linear algebra operation that has applications in solving linear equations and in Monte Carlo simulations.

### 5.1.1. The ac_chol_d Implementation

The *ac_chol_d* library provides four overloaded functions for computing the Cholesky Decomposition of real and complex matrices, and returns the lower triangular matrix result of the Cholesky Decomposition. Combined, the four functions handle two datatypes, which are (a) ac_fixed and (b) ac_complex<ac_fixed>. The matrix of data elements of each type can either be passed as standard two-dimensional C++ arrays, or can be packaged in the ac_matrix class.

#### Algorithm

As discussed earlier, the Cholesky-Crout algorithm is used to compute the Cholesky Decomposition. The computation for this algorithm is done in a column-wise manner. We first compute the diagonal element for each matrix either using the accurate or the approximate (PWL) *sqrt* function. After we calculate the diagonal element, we store its inverse in a separate variable, which is then reused for the computation of the non-diagonal elements below the diagonal. This inverse square root value can be calculated using the accurate *ac_div* function, or the approximate PWL version from the *ac_inverse_sqrt_pwl* library.

The Cholesky-Crout algorithm is used due to its simplicity and the reusability of the reciprocal value in calculating the non-diagonal elements.

#### Accurate Math Functions vs. PWL Approximations

The user has the option of being able to choose the accurate versions of the reciprocal and the sqrt functions, or their PWL approximations, as mentioned earlier, depending upon how much accuracy they may desire. Both have their advantages and disadvantages. The accurate math functions, while providing high precision,

can also add a lot of overhead in terms of throughput/area. The PWL functions, while providing higher throughput at a lesser cost in area, are also imprecise.

To use PWL approximations, the user has to override a default template parameter (*use_pwl*) for the *ac_chol_d* function call. An example of doing so is giving in the Example Function Calls section.

## Type of Intermediate Variables

Intermediate variables are used within the function to store the result of repeated subtractions in the process of calculating each element, as well as the reciprocal of diagonal elements. The default precision of these intermediate variables is set to be equal to the precision of the output variables. However, the user can choose to add any number of bits to these default values for word width as well as the integer width of the intermediate variables (in case of complex intermediate variables, these extra bits are used for the real/imaginary parts). They can do this by overriding the default template parameters (*delta_w* and *delta_i* for word and integer width, respectively). Furthermore, the user can also choose the rounding and saturation modes for the intermediate variables by overriding the appropriate template parameters (int_Q for rounding mode and int_O for saturation). By default, rounding and saturation is turned on for intermediate variables. An example of how to do so is giving in the Example Function Calls section below.

## Input Checking

As explained earlier, the input matrix must be positive definite. While calculating the values for the diagonal elements of the matrix, a square root operation is performed, either using a PWL approximation or the accurate math function. If the input to the *sqrt* function is negative/zero for any value, that means that the input matrix is not positive definite, and a macro-enabled *AC_ASSERT* is provided that will throw a run-time error in such a case. In case the *AC_ASSERT* does not kick in, additional, synthesizable functionality is also provided to output a matrix of zeros.

For certain input matrices when the PWL implementations are used for calculation, the inaccuracy incurred during computations might be large enough to result in the input matrix being wrongly perceived as not positive definite even when it is such. This particularly occurs when the diagonal values of the output matrix are very small and hence, as a result, even a small absolute error in the calculation of the reciprocal value of this diagonal quickly blows up and results in a large error when the remaining elements of that column are calculated. As this error builds up during the calculation of the value for the next diagonal element, the value calculated as the input to the square root function can turn out to be negative, hence resulting in the input checking failing for this particular case even though the input matrix is positive definite.

## 5.1.2. Function Prototypes

The following are the overloaded function prototypes for the *ac_chol_d* function to handle different datatypes:

```
template<bool use_pwl = false,
         int delta_w = 0, int delta_i = 0,
         ac_q_mode int_Q = AC_RND, ac_o_mode int_O = AC_SAT,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, bool outS, ac_q_mode outQ,
         ac_o_mode outO,
         unsigned M>
void ac_chol_d(
  const ac_fixed<W, I, S, Q, O> A[M][M],
  ac_fixed<outW, outI, outS, outQ, outO> L[M][M]
```

```
)
```

```
template<bool use_pwl = false,
         int delta_w = 0, int delta_i = 0,
         ac_q_mode int_Q = AC_RND, ac_o_mode int_O = AC_SAT,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, bool outS, ac_q_mode outQ, ac_o_mode outO,
         unsigned M>
void ac_chol_d(
  const ac_complex<ac_fixed<W, I, S, Q, O> > A[M][M],
  ac_complex<ac_fixed<outW, outI, outS, outQ, outO> > L[M][M]
)
```

```
template<bool use_pwl = false,
         int delta_w = 0, int delta_i = 0,
         ac_q_mode int_Q = AC_RND, ac_o_mode int_O = AC_SAT,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, bool outS, ac_q_mode outQ, ac_o_mode outO,
         unsigned M>
void ac_chol_d(
  const ac_matrix<ac_fixed<W, I, S, Q, O>, M, M> &A,
  ac_matrix<ac_fixed<outW, outI, outS, outQ, outO>, M, M> &L
)
```

```
template<bool use_pwl = false,
         int delta_w = 0, int delta_i = 0,
         ac_q_mode int_Q = AC_RND, ac_o_mode int_O = AC_SAT,
         int W, int I, bool S, ac_q_mode Q, ac_o_mode O,
         int outW, int outI, bool outS, ac_q_mode outQ, ac_o_mode outO,
         unsigned M>
void ac_chol_d(
  const ac_matrix<ac_complex<ac_fixed<W, I, S, Q, O> >, M, M> &A,
  ac_matrix<ac_complex<ac_fixed<outW, outI, outS, outQ, outO> >, M, M> &L
)
```

# 5.1.3. C++ Compiler

The functions use default template arguments. This requires using a C++ compiler that supports the C++11 standard. Failing to use such a compiler will result in compilation errors.

# 5.1.4. Example Function Calls

An example of a function call to store the lower triangular Cholesky Decomposition matrix of a sample 8x8 positive definite matrix A in a matrix L is shown below:

```
typedef ac_matrix<ac_fixed<20,11,true,AC_RND,AC_SAT>,8,8> i_type;
typedef ac_matrix<ac_fixed<30,15,true,AC_RND,AC_SAT>,8,8> o_type;
i_type A;
o_type L;
//Hypothetical function that generates a positive definite matrix
//and stores the output in A.
gen_pos_def_matrix(A);
ac_chol_d(A, L);
```

The ac_matrix L hereafter stores the lower triangular matrix result of the Cholesky decomposition.

## Choosing PWL approximation

By default, accurate math functions are used for the calculation of internal variables. However, as mentioned earlier, the user can override the default. They can use PWL approximation functions by giving a Boolean template argument, as follows:

```
ac_chol_d<true>(A, L);
```

## Configuring the Type for Temporary Variables

The user can add extra bits to the default precision of the temporary variables in the function, by overriding the default template parameters. They can also add other rounding and saturation modes for the temporary variables. They can do it as follows:

```
ac_chol_d<false, 10, 5, AC_TRN, AC_WRAP>(A, L);
```

By doing this, the user will add 10 bits to the default value for the temporary variable word width, 5 bits to the default value for the temporary variable integer width and switch off rounding and saturation for the temporary variables. Note that the user must also explicitly specify whether they want to use PWL approximation or the accurate math functions in such a case. (In this case, the user specifies that they want to use the accurate math functions)

Similarly, the user can also subtract bits from the default precision values. To do so, they merely need to pass negative parameters for the same. If the user wants to subtract 3 bits from word width and 2 bits from integer width, they can pass the following template parameters:

```
ac_chol_d<false, -3, -2>(A, L);
```

# 5.2. Cholesky Inverse (ac_cholinv)

The *ac_cholinv* function provides matrix inversion of a positive definite matrix using forward substitution and Cholesky Decomposition which uses the Cholesky-Crout algorithm. The user can utilize either accurate math functions or the PWL math library for the internal calculations involved.

# 5.2.1. The ac_cholinv Implementation

The *ac_cholinv* library provides two overloaded functions for computing the inverse of input matrices and returns the inverted matrix. Each overloaded function handles a different input datatype. The two datatypes hence handled are *ac_fixed* and *ac_complex<ac_fixed>*. The matrix of data elements of each type are packaged in the *ac_matrix* class.

## Algorithm

As mentioned earlier, the Cholesky-Crout algorithm is used to compute the Cholesky Decomposition which is done in the *ac_chol_d.h* library and it returns a lower triangular matrix. Please refer to the documentation of the library *ac_chol_d.h* for Cholesky Decomposition. Then the inverse of the lower triangular matrix is computed using forward substitution and then this inverse is multiplied with its conjugate transpose to obtain the inverse matrix.

## Accurate Math Functions vs. PWL Approximations

The user has the option of being able to choose the accurate versions of the reciprocal and the *sqrt* functions, or their PWL approximations, as mentioned earlier, depending upon how much accuracy they may desire. Both have their advantages and disadvantages. The accurate math functions, while accurate, can also add a lot of overhead in terms of throughput/area. The PWL functions, while providing higher throughput, are a bit imprecise.

To use the pwl functions, the user has to override a default template parameter (*use_pwl1*) for the *ac_chol_d* function call in the *ac_cholinv* function and (use_pwl2) for the *ac_cholinv* function call. An example of doing so is giving in the *Example Function Calls* section.

## Type of Intermediate Variables

Intermediate variables are used within the function to store the result of repeated additions in the process of calculating each element, as well as the reciprocal of diagonal elements. The default precision of these temporary variables is set to be equal to the precision of the output variables. However, the user can choose to add any number of bits to these default values for word width as well as the integer width of the temporary variables (in case of complex temporary variables, these extra bits are used for the real/imaginary parts). They can do this by overriding the default template parameters (add2w and add2i for word and integer width, respectively). Furthermore, the user can also choose the rounding and saturation modes for the temporary variables by overriding the appropriate template parameters (temp_Q for rounding mode and temp_O for saturation). An example of how to do so is giving in the Example Function Calls section.

# 5.2.2. Function Prototypes

The following are the overloaded function prototypes for the ac_chol_d function to handle different datatypes:

```
template<bool use_pwl1 = false,
        bool use_pwl2 = false,
        int add2w = 0,
        int add2i = 0,
        ac_q_mode temp_Q = AC_RND,
        ac_o_mode temp_O = AC_SAT,
        unsigned M,
        class T_in,
         class T_out>
 void ac_chol_d(
   const ac_matrix<T_in, M, M> &A,
   ac_matrix<T_out, M, M> &L
 )
 template<bool use_pwl1 = false, bool use_pwl2 = false, int      add2w = 0, int add2i =
0, ac_q_mode temp_Q = AC_RND, ac_o_mode temp_O = AC_SAT, unsigned M, class T_in, class
T_out>
  void ac_chol_d(
   const ac_matrix<ac_complex<T_in>, M, M> &A,
   ac_matrix<ac_complex<T_out>, M, M> &L
 )
```

# 5.2.3. C++ Compiler

The functions use default template arguments. In order to use a C++ compiler that supports this functionality, the user must use C++11 or a later standard as the standard for their compilation, failing which a compile-time error is thrown.

# 5.2.4. Example Function Calls

An example of a function call to store the lower triangular Cholesky Decomposition matrix of a sample 8x8 positive definite matrix A in a matrix L is shown below:

```
typedef ac_matrix<ac_fixed<20,11,true,AC_RND,AC_SAT>,8,8> i_type;
typedef ac_matrix<ac_fixed<30,15,true,AC_RND,AC_SAT>,8,8> o_type;
i_type A;
o_type Ainv;
//Hypothetical function that generates a positive definite matrix
//and stores the output in A.
gen_pos_def_matrix(A);
ac_cholinv(A, Ainv);
```

## Choosing PWL approximation

By default, accurate math functions are used for the calculation of internal variables. However, as mentioned earlier, the user can override the default. The first Boolean template parameter if true uses PWL approximation functions for calculation of internal variables and if false uses accurate math functions in the ac_chol_d library and similarly if the second Boolean template parameter if true uses PWL approximation functions for calculation of internal variables and if false uses accurate math functions in the ac_cholinv library.

1. If the user wants to use the PWL approximation functions for both the ac_chol_d and ac_chol_inv libraries, then the following is the way the function should be called

   ```
   ac_chol_inv<true, true>(A, L);
   ```

2. If the user wants to use the PWL approximation functions for the ac_chol_d library and accurate math functions for the ac_chol_inv library, then the following is the way the function should be called.

   ```
   ac_chol_inv<true, false>(A, L);
   ```

   Note: If the user wants to specify the template parameter for choosing the PWL approximation functions only for the ac_cholinv library, then the user has to explicitly specify whether they want to use PWL approximation or the accurate math functions for the ac_chol_d library as well.

## Configuring the Type For Temporary Variables

The user can add extra bits to the default precision of the temporary variables in the function, by overriding the default template parameters. They can also add their own quantization and overflow modes for the temporary variables. For example, they can do it as follows:

```
ac_chol_inv<false, false, 10, 5, AC_RND, AC_SAT>(A, Ainv);
```

By doing this, the user will add 10 bits to the default value for the temporary variable word width, 5 bits to the default value for the temporary variable integer width and switch on rounding and saturation for the temporary variables. Note that the user must also explicitly specify whether they want to use PWL approximation or the accurate math functions in such a case. (In this case, the user specifies that they want to use the accurate math functions)

Similarly, the user can also subtract bits from the default precision values. To do so, they merely need to pass negative parameters for the same. If the user wants to subtract 3 bits from word width and 2 bits from integer width, they can pass the following template parameters:

```
ac_chol_d<false, false, -3, -2>(A, Ainv);
```

# 5.3. Determinant (ac_determinant)

The *ac_determinant* library implementation is a fully parallelized and scalable implementation for determinant computation using template recursion functionality, with the user being able to choose between using an internally determined datatype or adding their own parameters for intermediate type precision, signedness, rounding and saturation.

## 5.3.1. The ac_determinant Implementation

The *ac_determinant* library provides four overloaded functions for computing the determinant of real and complex matrices. Combined, the four functions handle two datatypes, which are (a) ac_fixed and (b) ac_complex<ac_fixed>. The matrix of data elements of each type can either be passed as standard two-dimensional C++ arrays, or can be packaged in the ac_matrix class.

### Algorithm

Determinant computation is a recursive process. This implementation leverages the recursive tendency of this function to design a hardware-efficient implementation that is scalable and is highly parallelizable. Hence, higher order matrices are reduced to 2x2 matrices by computing minors recursively on them, with each recursion reducing the row and column size by 1. A specialization is also defined for a 1x1 matrix, in which case, we just return the sole value stored in the matrix.
In order to implement a fully parallelizable and synthesizable architecture, template recursion is used, with the 2x2 matrix being the specialized case for recursion.

### Internal Precision Adjustment

The ac_determinant implementation is such that it maintains full internal precision. To make sure that it does that, bitwidths are computed at every step of the coding.

Using full internal precision can result in a very large bitwidth, however. In case full internal precision might not be required, the user can set their own internal precision by overriding the default template parameters and providing their own internal types.

## 5.3.2. Function headers

The following are the overloaded function prototypes for the ac_determinant function to handle different datatypes:

```
template <bool override = false,
          int internal_width = 16, int internal_int = 8,
          bool internal_sign = true, ac_q_mode internal_rnd = AC_RND,
          ac_o_mode internal_sat = AC_SAT,
          unsigned M,
          int W1, int I1, bool S1, ac_q_mode q1, ac_o_mode o1,
```

```
            int W2, int I2, ac_q_mode q2, ac_o_mode o2>
void ac_determinant(
  ac_matrix <ac_fixed <W1, I1, S1, q1, o1>, M, M> &input,
  ac_fixed<W2, I2, true, q2, o2> &result
)
```

```
template <bool override = false,
          int internal_width = 16, int internal_int = 8,
          bool internal_sign = true, ac_q_mode internal_rnd = AC_RND,
          ac_o_mode internal_sat = AC_SAT,
          unsigned M,
          int W1, int I1, bool S1, ac_q_mode q1, ac_o_mode o1,
          int W2, int I2, ac_q_mode q2, ac_o_mode o2>
void ac_determinant(
  const ac_matrix<ac_complex<ac_fixed<W1,I1,S1,q1,o1> >, M, M> &input,
  ac_complex<ac_fixed<W2, I2, true, q2, o2> > &result
)
```

```
template <bool override = false,
          int internal_width = 16, int internal_int = 8,
          bool internal_sign = true, ac_q_mode internal_rnd = AC_RND,
          ac_o_mode internal_sat = AC_SAT,
          unsigned M,
          int W1, int I1, bool S1, ac_q_mode q1, ac_o_mode o1,
          int W2, int I2, ac_q_mode q2, ac_o_mode o2>
void ac_determinant(
  const ac_fixed<W1, I1, S1, q1, o1> a[M][M],
  ac_fixed<W2, I2, true, q2, o2> &result
)
```

```
template <bool override = false,
          int internal_width = 16, int internal_int = 8,
          bool internal_sign = true, ac_q_mode internal_rnd = AC_RND,
          ac_o_mode internal_sat = AC_SAT,
          unsigned M,
          int W1, int I1, bool S1, ac_q_mode q1, ac_o_mode o1,
          int W2, int I2, ac_q_mode q2, ac_o_mode o2>
void ac_determinant(
  const ac_complex<ac_fixed<W1, I1, S1, q1, o1> > a[M][M],
  ac_complex<ac_fixed<W2, I2, true, q2, o2> > &result
)
```

The library also provides the following two functions to allow the user to return by value.

```
template<class T_out,
         bool override = false,
         int internal_width = 16, int internal_int = 8,
         bool internal_sign = true, ac_q_mode internal_rnd = AC_RND,
         ac_o_mode internal_sat = AC_SAT,
         unsigned M, class T_in>
T_out ac_determinant(const ac_matrix <T_in, M, M> &input)
```

```
template<class T_out,
         bool override = false,
         int internal_width = 16, int internal_int = 8,
         bool internal_sign = true, ac_q_mode internal_rnd = AC_RND,
         ac_o_mode internal_sat = AC_SAT,
         unsigned M, class T_in>
```

```
T_out ac_determinant(const T_in input[M][M])
```

## 5.3.3. C++ Compiler

The functions use default template arguments. This requires using a C++ compiler that supports the C++11 standard, or a later standard. Failing to use such a compiler will result in compilation errors.

## 5.3.4. Example Function Call

The following give examples of using the determinant function with both ac_matrix and C style array inputs.

```
// Using the ac_matrix wrapper class
ac_matrix <ac_fixed <8, 5, true, AC_RND, AC_SAT>, 2, 2> input;
typedef ac_fixed <27, 18, true, AC_RND, AC_SAT> output_type;
output_type output;
// Assign elements to input
input(0, 0) = 1;
input(0, 1) = 1.5;
input(1, 0) = 0.5;
input(1, 1) = 1.75;
ac_determinant (input, output);
// The above function call can also be replaced by:
output = ac_determinant <output_type> (input);
```

```
// Using C-style arrays
ac_fixed <8, 5, true, AC_RND, AC_SAT> input[2][2];
typedef ac_fixed <27, 18, true, AC_RND, AC_SAT> output_type;
output_type output;
// Assign elements to input
input[0][0] = 1;
input[0][1] = 1.5;
input[1][0] = 0.5;
input[1][1] = 1.75;
ac_determinant (input, output);
// The above function call can also be replaced by:
output = ac_determinant <output_type> (input);
```

# 5.4.  Matrix Multiplication (ac_matrixmul)

The ac_matrixmul function is designed to provide the multiplication of two matrices with full precision by default and if user wants to specify his own precision then a provision is provided for this purpose.

## 5.4.1. The ac_matrixmul Implementation

 The ac_matrixmul library provides implementation for ac_fixed and    ac_complex<ac_fixed> datatypes. There is one generalized function which handles the  computation for both the datatypes. A prototype of the ac_matrixmul function is shown  below:

```
template<int mw = 0, int mi = 0, ac_q_mode mq = AC_RND,
         ac_o_mode mo= AC_SAT,
         int sw = 0, int si = 0, ac_q_mode sq = AC_RND,
         ac_o_mode so = AC_SAT,
```

```
        unsigned M, unsigned N, unsigned P, class T_in_A,
        class T_in_B, class T_op>
void
ac_matrixmul(
    const ac_matrix<T_in_A, M, N> &A,
    const ac_matrix<T_in_B, N, P> &A,
    ac_matrix<T_op, M, P> &C
)
```

## Definition of matrices:

In the above prototype, ac_matrix is a container class that helps to define two dimensional matrices. The following is the way to define a matrix using the ac_matrix class for the ac_fixed datatype:

```
ac_matrix<ac_fixed<20, 10, true, AC_RND, AC_SAT>, M, N> A;
```

In the definition above, A is a matrix of dimension M x N where its elements are of ac_fixed datatype.

Similarly, the following is the way to define a matrix using the ac_matrix class for the ac_complex<ac_fixed> datatype:

```
ac_matrix<ac_complex<ac_fixed<40, 20, true, AC_RND, AC_SAT> >, M, N> A;
```

In the definition above, A is a matrix of dimension M x N where its elements are of ac_fixed datatype.

## Default Template Arguments:

By default, the internal variables have full precision turned on. There are two internal variables mult and sum used in the computation of matrix multiplication. The code snippet where the sum and mult variables are used in the matrix multiplication block is as follows:

```
for (unsigned i=0; i<M; i++) {
  for (unsigned j=0; j<P; j++) {
     mult = 0;
     sum = 0;
     for (unsigned k=0; k<N; k++) {
         mult = A(i,k) * B(k,j);
         sum = sum + mult;
     }
     C(i,j) = sum;
  }
}
```

The default template arguments in the order which they have been mentioned in the list of template parameters in the function prototype for this block are:

- Width of the mult variable – mw

- Integer width of the mult variable – mi

- Quantization mode of the mult variable – mq

- Overflow mode of the mult variable – mo

- Width of the sum variable – sw

- Integer width of the sum variable – si

- Quantization mode of the sum variable – sq

- Overflow mode of the sum variable - so

## For the ac_fixed datatype

From the above code snippet, in order to achieve full precision for the mult variable its width and integer width is calculated as the sum of the widths and the sum of the integer widths of the elements of A and B matrices respectively. For the sum variable, the full precision is dependent on the number of times it is accumulated by the mult variable which is in turn decided by the number of iterations in the for loop in which it is accumulated. As the number of iterations are N in this for loop, the width and the integer width for the sum variable is calculated as the width of the mult variable + log2_ceil(N) and the integer width of the mult variable + log2_ceil(N) respectively. The default quantization and overflow modes for the mult and the sum variables are defined as AC_RND and AC_SAT respectively.

## For the ac_complex<ac_fixed> datatype

Everything above applies for the ac_complex<ac_fixed> datatype except that the width and integer of the mult variable is calculated as the sum of the widths + 1 and the sum of the integer widths + 1 of the elements of A and B matrices respectively. This is because when two complex numbers are multiplied for example:

```
(a1+b1i) x (a2+b2i) = a1a2 – b1b2 + (a1b2 + a2b1)i
```

The computation of the real part and the imaginary part involves an extra addition or subtraction and thus 1 is added. The user must keep in mind that in order to allow the usage of default template arguments for functions, they must use a version of C++ that allows this usage, such as C++ 11, in case they wish to compile and execute the code for the matrixmul library, failing which an error will be thrown by the compiler.

If the user wants to specify his/her own precision for the internal variables instead of the default, then please refer to the "Example Function Calls" section below for the same.

## Wrong Datatypes

If any of the two inputs and the output have different datatypes from each other, then a static assert is thrown stating that '*Both arguments need to be of the same datatype*' in case of C++11 version of the compiler and in versions prior to that a compiler error is generated.

# 5.4.2. Example Function Call

An example of a function call to store the result matrix of the two input matrices A and B who have elements of ac_fixed datatype is shown below:

```
typedef ac_matrix<ac_fixed<20, 10, true, AC_RND, AC_SAT>, M, N> input_type1;
typedef ac_matrix<ac_fixed<20, 10, true, AC_RND, AC_SAT>, N, P> input_type2;
typedef ac_matrix<ac_fixed<61, 31, true, AC_RND, AC_SAT>, M, P> output_type;
input_type1 A;
input_type2 B;
output_type C;
ac_matrixmul(A, B, C);    //Function call
```

The two input matrices here are A and B and the output matrix here is C.

## Changing default template parameters through function call:

- If the user wants a different precision from the default precision, then the default

template parameters have to be modified. If the user wants to modify all the template parameters then the user has to modify the above function call in the following way:

```
ac_matrixmul<mw, mi, mq, mo, sw, si, sq, so> (A, B, C);
```

The above order in which the template arguments are specified has to be maintained.

- Suppose the user only wants to modify the width of the mult variable i.e mw and keep all the other template arguments as it is, then the function call can be written as:

```
ac_matrixmul<mw> (A, B, C);
```

- But if the user wants to modify only sw for instance, then in order to do this the user has to make a function call like this:

```
ac_matrixmul<mw, mi, mq, mo, sw> (A, B, C);
```

## 5.4.3. Debug

Here is a provision to enable debug mode in the *ac_matrixmul.h* file. It can be done by defining the macro as

```
#define MATRIXMUL_DEBUG
```

# 5.5.  QR Decomposition (ac_qrd)

The ac_qrd is a function that provides QR decomposition of input matrix A. In linear algebra, QR decomposition is a decomposition of matrix A into product A = QR where Q is an orthogonal matrix and R is an upper triangular matrix.

This function uses systolic array based implementation where each iteration performs a Given's rotation algorithm to convert given matrix into R matrix, parallelly Q matrix is computed by performing the same rotations on the identity matrix. Note that this function uses ac_matrix container class and input and output matrices are either supplied in the form of ac_matrix objects or via standard AC datatype 2D arrays. This QR decomposition implementation supports ac_fixed and ac_complex <ac_fixed> datatypes and implementation remain same for both. If the function is called with other datatypes, compilation error is thrown.

## 5.5.1. The ac_qrd Implementation

The ac_qrd implementation used Givens rotation algorithm to convert any given matrix into an upper triangular matrix which is R. Computation of Q requires same set of sequence of rotations and can be computed simultaneously. Givens rotation has a property of implementing the QRD decomposition in systolic manner, which basically makes it possible to synthesize fully parallelized architecture, making throughput of 1 achievable, saving time and space complexity.

### Givens rotation algorithm

Givens rotation basically applies series of orthogonal rotations on the original input matrix, A and zeroes the lower triangular submatrix one element at a time.

For this purpose, first a 2x2 rotation matrix is computed (call it rot_mat(i,k)), which is:

```
rot_mat (i,k) = [  c   s  ]
                [ -s   c  ]
```

 where, i = [2,m], k=[1,n], such that, i > k

Here c and s stands for cosine and sine parameters respectively and $c^2 + s^2 = 1$ always holds true.

Each such (i,k) is computed and then is applied to the input matrix, which zeroes out element A(i,k) of the input matrix. This rotations applied in terms of rotation matrix are applied from element A(4,1) and applied all the way to the up until the diagonal. After that we move to the next column and so on and so forth.

Note that the givens rotation in each iteration is applied to only two rows (the row at which the element which is to be zeroed contains and a row above that), this makes it possible to the architecture in the Systolic manner (will be explained in later sections).

The two rotational matrix parameters c and s are computed using, following two formula:

```
c = x/sqrt(mod) and s = -y/sqrt(mod)
```

where mod = x^x+y^y

y is basically the element to be zeroed and x is the element in the row above that but present in the same column.

These values of c and s are then applied to the rows R1 (one above R2) and R2 (one containing the element to be zeroed).

Again, as mentioned earlier since this computation only involves two matrix rows of element at a time and zeroes one element at a time, multiple rotations can be applied in parallel if they work on different rows, leading to possibility of massive parallel architecture.

## Systolic Array Structure of QRD

QRD can be massively parallelized if only it performs rotation on two different set of rows. This algorithm uses two different 1D arrays two store and the operate on two different rows at a time.

The above process of applying rotations and computing rotational parameters c and s can be divided into two different functions (processing elements) that are called as off-diagonal and diagonal processing elements (PEs). Diagonal PEs are used to compute c and s, where as off-diagonal PEs are used to apply them to the rows and compute the new values of the rows in the final R matrix.

- **Diagonal Processing Elements** are expressed in the form of function, diagonal_PE which takes element to be zeroed (y) and element in the row just above that in the same column (x) as inputs and returns the rotational matrix parameters. It also takes the default Boolean parameter ispwl, which allows user to do space-time exploration by switching between pwl version of reciprocal and divider. By default, divider is used to make sure that the error is minimum, although this increases the area of the Diagonal Processing elements.

- **Off-Diagonal Processing Elements** are expressed using function offdiagonal_PE which takes the input matrix (which is combined matrix in the implementation), pivot (which acts as an indicator variable of which row element is to zeroed) and orthogonal parameters.

The computation for rows using c and s works on following formula:

```
r1'(i) = c*r1(i) + s*r2(i)
and r2'(i) = c*r2(i) – s*r1(i)
```

where i varies from 0 to the matrix size M.

Note that above algorithm although has a dependency which is broken by using temporary space and hence the processing element can be fully pipelined.

The diagonal and off-diagonal elements can be separated from rest of the algorithm and can be put as CCOREs leading to further chance of improvement from system perspective.

## Computation of Q

Q is an orthogonal matrix and computation of Q in Givens rotations is computed by taking transpose of multiplication of matrices that are obtained after each Givens rotation in input matrix.

In other words,

```
Q = (G1*G2*G3)'
```

Where G1, G2, G3 are outputs of Givens rotation after each iteration on input matrix.

This calculation of Q requires matrix multiplication and hence takes cubic time complexity. To avoid this, implementation uses the Systolic properties of QRD to apply rotations on the identity matrix. Since QRD has systolic nature and since only two rows are changing in every rotation, applying rotations in systolic manner to the previously obtained matrix is equivalent to the matrix multiplication.

Hence, first the given matrix is combined with the identity matrix to obtain a higher order matrix. Identity matrix used is of size M*M and it is placed after the input matrix to form a new higher order matrix of size 2M*M. After passing this matrix through the off-diagonal and diagonal processing elements we obtain final matrix that is basically of 2M*M size, from which then Q and R can be separated, by making sure that the place at which input matrix was present is new R and place at which identity matrix was present is now, Q'. Then transpose is taken of Q' to obtain Q. This whole process of forming higher order matrix (performed by initialize matrix function) and dividing the higher order matrix into final answers, Q and R (taken care by qr_divide function), requires two for loops each, which can be fully unrolled.

## PWL and non-PWL performance analysis

PWL functions are known for their efficiency in terms of area, while downside of being less accurate as compare to their accurate versions. Although in the applications where, area is more of point of concern, PWL implementation might prove more useful than their accurate counterparts. Hence, this function gives an option to the user to choose between inverse square root pwl and accurate divider and square root pwl, using a Boolean parameter, 'ispwl'.

If set to true, this implementation switches to the PWL implementation. In this case, the error is higher as compared to the non-PWL functions and in fact accumulates due to the for loops present in the function. As a result, we get higher error percentage, although significantly lesser area (it is observed that area is reduced as high as ten times, where as error percentage increases to the 10 times the error in case of accurate functions).

# 5.5.2. Function prototypes

## For ac_matrix of fixed point:

```
template<bool ispwl = false,
         unsigned M,
         int W1,
         int I1,
         ac_q_mode q1,
         ac_o_mode o1,
         int W2,
         int I2,
         ac_q_mode q2,
         ac_o_mode o2>
void ac_qrd(
```

```
ac_matrix <ac_fixed <W1, I1, true, q1, o1>, M, M> &A,
          ac_matrix <ac_fixed <W2, I2, true, q2, o2>, M, M> &Q,
          ac_matrix <ac_fixed <W2, I2, true, q2, o2>, M, M> &R
);
```

## For ac_matrix of complex of fixed point:

```
template<bool ispwl = false, unsigned M, int W1, int I1, ac_q_mode q1, ac_o_mode o1,
int W2, int I2, ac_q_mode q2, ac_o_mode o2>
void ac_qrd(
  ac_matrix <ac_complex <ac_fixed <W1, I1, true, q1, o1> >, M, M> &A, ac_matrix
<ac_complex <ac_fixed <W2, I2, true, q2, o2> >, M, M> &Q, ac_matrix <ac_complex
<ac_fixed <W2, I2, true, q2, o2> >, M, M> &R
);
```

## For array of fixed point:

```
template<bool ispwl = false, unsigned M, int W1, int I1, ac_q_mode q1, ac_o_mode o1,
int W2, int I2, ac_q_mode q2, ac_o_mode o2>
void ac_qrd(
  ac_fixed <W1, I1, true, q1, o1> A[M][M], ac_fixed <W2, I2, true, q2, o2>, M, M> Q[M]
[M], ac_fixed <W2, I2, true, q2, o2> R[M][M]
);
```

# Section 6.  Miscellaneous Functions

The *ac_math* package also provides the following functions, in addition to the PWL, LUT and CORDIC functions:

- Absolute Value (ac_abs)
- Division (ac_div)
- Square Root (ac_sqrt)
- Shifts (ac_shift_left/ac_shift_right)

The following subsections describes the implementation and usage of these functions in more detail.

## 6.1.  Absolute Value (ac_abs)

The absolute value (ac_abs) operation produces a positive result by negating values less than zero.

- Integer Signed

```
void ac_abs(ac_int<XW,true> x, ac_int<YW,false> &y)
void ac_abs(ac_int<XW,true> x, ac_int<YW,true> &y)
```

- Fixed Point Signed

```
void ac_abs(ac_fixed<XW,XI,true,XQ,XO> x, ac_fixed<YW,YI,false,YQ,YO> &y)
void ac_abs(ac_fixed<XW,XI,true,XQ,XO> x, ac_fixed<YW,YI,true,YQ,YO> &y)
```

- Float

```
void ac_abs(ac_float<XW,XI,XE,XQ> x, ac_float<YW,YI,YE,YQ> &y)
```

## 6.2.  Division (ac_div)

The division functions are supported for integer, fixed-point, complex, and float data types. The division functions takes two inputs: dividend and divisor and computes the quotient. If the inputs are integer there is a version of the function that also computes the remainder of the division. In all cases the div function returns true if the computed remainder is nonzero.

*NOTE: Division by zero triggers an assertion failure during simulation.*

### 6.2.1. Integer Division

There are several integer division functions defined. Some of the functions compute just the quotient, some compute both the quotient and the remainder. All return a bool flag to indicate whether the remainder is non zero.

The first argument is the dividend (or numerator), the second is the divider (or denominator), the third is the quotient (output argument) and the fourth (optional) argument is the remainder (output argument). All the arguments are either all signed or all unsigned. The computed quotient is equivalent to the resulted computed by the operator `'/'`:

```
ac_int<8,false> n;
```

```
ac_int<5,false> d;
ac_int<6,false> q;
ac_div(n, d, q);
ac_int<6,false> q1 = n/d;     // q == q1
```

While the behavior of the last two lines is identical, Catapult will produce different hardware for each case. Catapult will allocate a component from the library for the '/' whereas calling the div function will inline the function.

| | |
|---|---|
| ```bool ac_div(`<br>`  ac_int<NW,false> dividend,`<br>`  ac_int<DW,false> divisor,`<br>`  ac_int<QW,false> &quotient,`<br>`  ac_int<RW,false> &remainder`<br>`);``` | ```bool ac_div(`<br>`  ac_int<NW,true> dividend,`<br>`  ac_int<DW,true> divisor,`<br>`  ac_int<QW,true> &quotient,`<br>`  ac_int<RW,true> &remainder`<br>`);``` |
| ```bool ac_div(`<br>`  ac_int<NW,false> dividend,`<br>`  ac_int<DW,false> divisor,`<br>`  ac_int<QW,false> &quotient`<br>`);``` | ```bool ac_div(`<br>`  ac_int<NW,true> dividend,`<br>`  ac_int<DW,true> divisor,`<br>`  ac_int<QW,true> &quotient`<br>`);``` |

**Table 3: Functions for Integer Division that Return Remainder != 0**

For the div functions, dividend == divisor*quotient + remainder provided that both the quotient and the remainder have sufficient precision so that they don't overflow.

## 6.2.2. Fixed-point Division

There are several fixed-point division functions defined. Each function returns a *bool* flag to indicate whether the remainder is non zero (assumes that the quotient has enough precision that it does not overflow).

The first argument is the *dividend* (or numerator), the second is the *divider* (or denominator), the third is the *quotient* (output argument). All the arguments are either all signed or all unsigned. The computed quotient takes into account the target bitwidth, integer bitwidth, quantization and overflow modes. All quantization modes work correctly since they are based on the computation of the remainder. The operator '/', on the other hand, is forced to truncate the result without knowing the data type of the target. For example, the operator '/' for *ac_fixed* returns a result that is dependent on the type of both dividend and divisor and the bits (possibly infinite bits) to the right of the result are truncated before the quantization mode of the target is known.

The hardware produced by Catapult for a call to the '/' and '/=' operators and for a call to the div function are different. In the first case, Catapult allocates a component from the library, whereas in the second case the div function is inlined.

| | |
|---|---|
| ```bool ac_div(`<br>`  ac_fixed<NW,NI,false,NQ,NO> dividend,`<br>`  ac_fixed<DW,DI,false,DQ,DO> divisor,`<br>`  ac_fixed<QW,QI,false,QQ,QO> &quotient`<br>`);``` | ```bool ac_div(`<br>`  ac_fixed<NW,NI,true,NQ,NO> dividend,`<br>`  ac_fixed<DW,DI,true,DQ,DO> divisor,`<br>`  ac_fixed<QW,QI,true,QQ,QO> &quotient`<br>`);``` |

**Table 4: Functions for Fixed-Point Division that Return Remainder != 0**

The *div* functions return the *bool* value remainder != 0 which is equivalent to (divisor*quotient != dividend, provided that the quotient does not overflow) and may be dependent on the width and integer width of the quotient. For example:

```
ac_fixed<3,3,false> q; ac_fixed<4,3> q2;
ac_fixed<2,2,false> a = 3; ac_fixed<2,2,false> b = 2;
bool nonzero_rem = ac_div(a, b, q);    // nonzero_rem == true,  q == 1
bool nonzero_rem2 = ac_div(a, b, q2); // nonzero_rem2 == false, q2 == 1.5
```

## 6.2.3. Float Division

Inlined *ac_float* division is implemented using *ac_fixed* division and has an equivalent return value. The quotient may lose precision as neither the dividend or divisor are normalized. The resulting quotient will overflow if the result of division cannot be represented in the quotient type, otherwise the result is representable and normalized.

```
bool ac_div(
   ac_float<NW,NI,NE,NQ> dividend,
   ac_float<DW,DI,DE,DQ> divisor,
   ac_float<QW,QI,QE,QQ> &quotient
);
```

## 6.2.4. Complex Division

The *ac_div* library also provides for the division of complex inputs.

```
void ac_div(
   ac_complex<NT> dividend,
   ac_complex<DT> divisor,
   ac_complex<QT> &quotient
)
```

# 6.3.  Square Root (ac_sqrt)

The square root function has only input argument and one output argument both of which are unsigned. There is an integral version and a fixed-point version.

## 6.3.1. Integer square root

The integer square root computes the largest integer value such that when squared it is equal or less than the argument. The prototype of the function is given as follows:

```
void ac_sqrt(
   ac_int<XW,false> x,
   ac_int<OW,false> &sqrt
)
```

## 6.3.2. Fixed-point square root

The fixed-point *ac_sqrt* function allows for full flexibility on both the input and the output precision. The bitwidth, integer bitwidth, quantization and overflow modes of the target (second argument) determine the bits of precision for computing the square root and the quantization and overflow that is performed on the

result. All quantization modes are computed exactly based on the remainder of the square root computation. The prototype of the function is given as follows:

```
void ac_sqrt(
   ac_fixed<XW,XI,false,XQ,XO> x,
   ac_fixed<OW,OI,false,OQ,OO> &sqrt
)
```

An example of how the square root functions is used, consider the case where the input value is *ac_fixed<4,4, false>* (4 bit unsigned integer), and the required precision for the result is *ac_fixed<8,3, false>* (3 bit integral, 5 fractional: xxx.xxxxx), the call code would look like:

```
ac_fixed<4,4,false> x = 13;
ac_fixed<8,3,false> sqrt_x;
ac_sqrt(x, sqrt_x); // sqrt_x == (ac_fixed<8,3,false>) sqrt( x.to_double() );
```

# 6.4.  Shifts (ac_shift_left/ac_shift_right)

The *ac_shift* functions allow the specification of precision and quantization and overflow modes of the target. This gives a simple way to get around the issue of the fixed-point shift operations >> and << for *ac_fixed* returning the type of the first operand. The shifting provided by these functions is "arithmetic" in nature, that is, these functions provide the same result as multiplying the input by $2^{(shift\_count)}$, all the while keeping in mind the quantization and overflow modes of the target. In this sense, they allow for saturation and rounding, something that is not provided by the >> and << operators.

## 6.4.1. Bidirectional shifts

For ac_fixed data-types, both the right shift >> and the left shift << operators shift in the opposite direction when the shift value is negative. The equivalent functionality is provided by the following functions:

| | |
|---|---|
| <pre>void ac_shift_right(<br>   ac_fixed&lt;XW,XI,false,XQ,XO&gt; x,<br>   int n,<br>   ac_fixed&lt;OW,OI,false,OQ,OO&gt; &sr<br>)</pre> | <pre>void ac_shift_right(<br>   ac_fixed&lt;XW,XI,true,XQ,XO&gt; x,<br>   int n,<br>   ac_fixed&lt;OW,OI,true,OQ,OO&gt; &sr<br>)</pre> |
| <pre>void ac_shift_left(<br>   ac_fixed&lt;XW,XI,false,XQ,XO&gt; x,<br>   int n,<br>   ac_fixed&lt;OW,OI,false,OQ,OO&gt; &sl<br>)</pre> | <pre>void ac_shift_left(<br>   ac_fixed&lt;XW,XI,true,XQ,XO&gt; x,<br>   int n,<br>   ac_fixed&lt;OW,OI,true,OQ,OO&gt; &sl<br>)</pre> |

**Table 5: Functions for Fixed-Point Bidirectional Shifts**

## 6.4.2. Unidirectional shifts

If the shift value is know to be non-negative, it is best to cast it to (unsigned int) so that the following functions are inlined. These functions will deliver better quality of results during Catapult synthesis.

```
void ac_shift_right(
  ac_fixed<XW,XI,false,XQ,XO> x,
  unsigned int n,
  ac_fixed<OW,OI,false,OQ,OO> &sr
)
```
```
void ac_shift_right(
  ac_fixed<XW,XI,true,XQ,XO> x,
  unsigned int n,
  ac_fixed<OW,OI,true,OQ,OO> &sr
)
```
```
void ac_shift_left(
  ac_fixed<XW,XI,false,XQ,XO> x,
  unsigned int n,
  ac_fixed<OW,OI,false,OQ,OO> &sl
)
```
```
void ac_shift_left(
  ac_fixed<XW,XI,true,XQ,XO> x,
  unsigned int n,
  ac_fixed<OW,OI,true,OQ,OO> &sl
)
```

**Table 6: Functions for Fixed-Point Unidirectional Shift**

# 6.4.3. Complex shifts

Wrapper functions for the unidirectional *ac_shift_left* and *ac_shift_right* are provided for complex types to perform the corresponding shift for the real and imaginary parts of the underlying type.

```
void ac_shift_right(
  ac_complex<XT> x,
  unsigned int n,
  ac_complex<OT> &sr
)
```
```
void ac_shift_left(
  ac_complex<XT> x,
  unsigned int n,
  ac_complex<OT> &sl
)
```

**Table 7: Functions for Complex Unidirectional Shift**