



Magical Math

Adam@adiuvoengineering.com

Objective

The objective of this session are:

- Introduction to applications using maths in programmable logic
- The difference between fixed and floating point maths
- Why programmable logic is often more suited to fixed point maths
- How do we work with fixed point maths in programmable logic
- What are the rules of fixed point maths in programmable logic
- How can we implement complex algorithms and filters using fixed point maths
- The benefits of using AMD Vitis™ High Level Synthesis to implement algorithms
- Leveraging Matlab / Simulink and AMD Vitis™ Model Composer to implement math solutions.

Why We Need to Do Math in FPGA

Programmable Logic enables accelerated applications

- **Image Processing** – Noise removal, edge detection, filtering
- **Radar Processing** – Signal generation, reply processing
- **Signal Processing** – Signal filtering & manipulation
- **Robotics** – End Effector positioning, Navigation
- **Control Systems** – Kalman Filtering, PID Loops
- **Motor Control** – Servo Motor, DC Motor Control

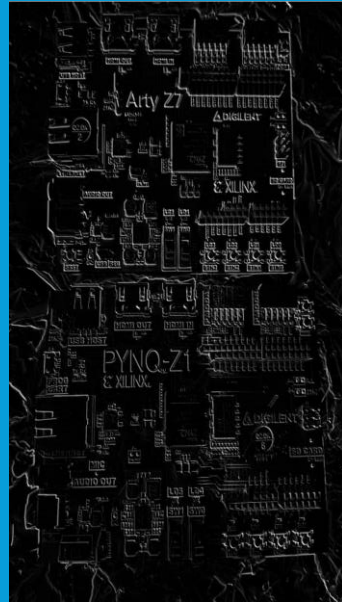
All these applications require the ability of FPGA to do maths and implement algorithms

Example Applications

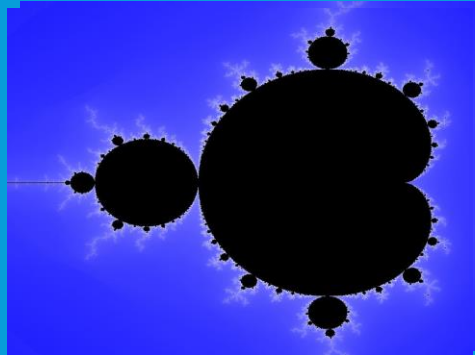
Background Removal & Substitution



Edge Detection



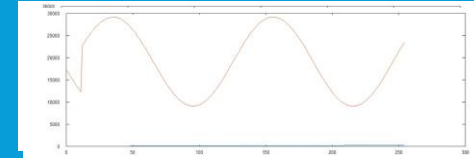
Fractals



AI/ML



Signal Processing



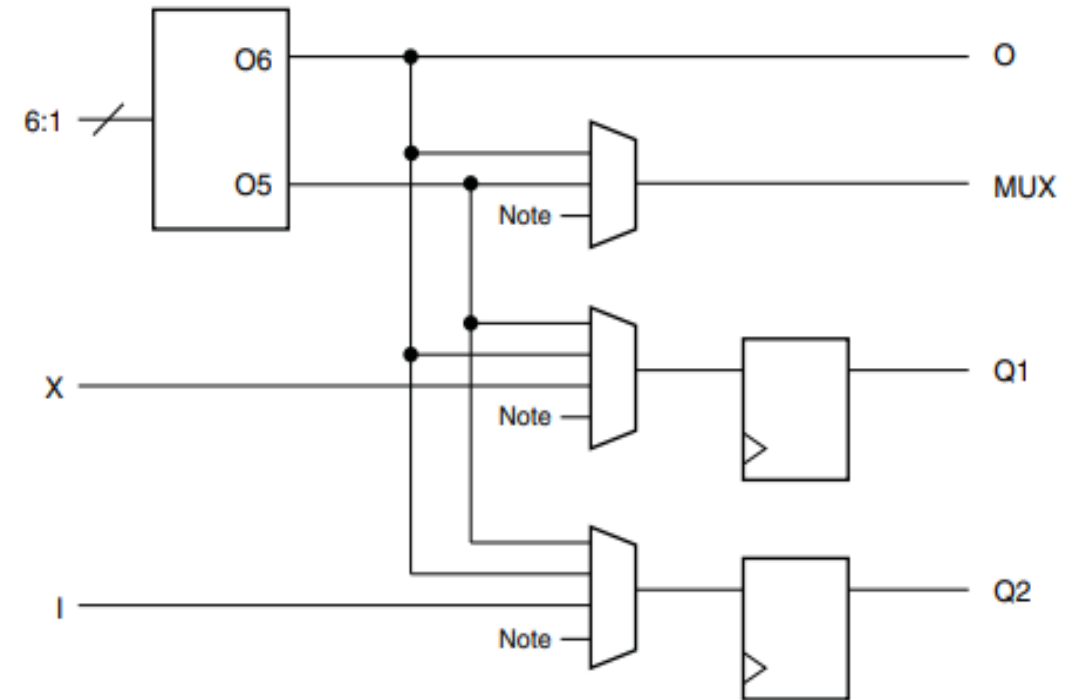
FPGA Architecture

FPGA are register and logic rich

Configurable Logic Blocks contain:

- Registers
- Look Up Table
- Distributed RAM
- Carry Mux

Logic resources are the basic building blocks of our algorithms. It is where we implement our mathematical algorithms.



DSP Elements

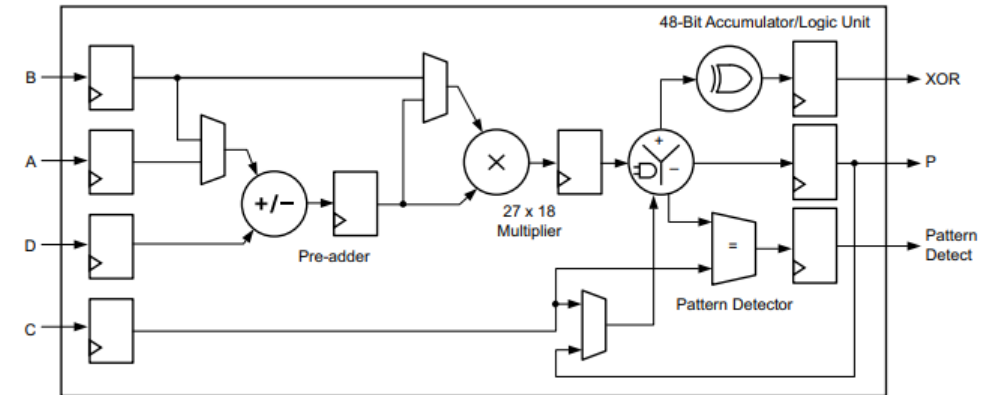
Implementing math directly in logic is costly

- Resources – increased resources
- Performance – reduced performance

Manufactures to address this include dedicated DSP elements (e.g., DSP48)

Capable of doing 48 bit Multiply accumulator

- Pre-Adder
- Multiply
- Accumulator
- Can do advance things SIMD – More later!



FPGA Maths

So far, we have looked at Logic & DSP elements, but they all have one thing in common ?

They are all ideal for the implementation of fixed-point solutions.

What is the difference between fixed point and floating-point numbers ?

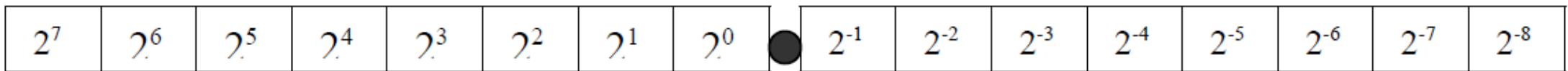
Fixed Point Number

Fixed-point representation maintains the decimal point within a fixed position allowing for straight forward arithmetic operations.

The major drawback of fixed-point representation is that to represent larger numbers or to achieve a more accurate result with fractional numbers, a larger number of bits are required.

A fixed-point number consists of two parts called the integer and fractional parts.

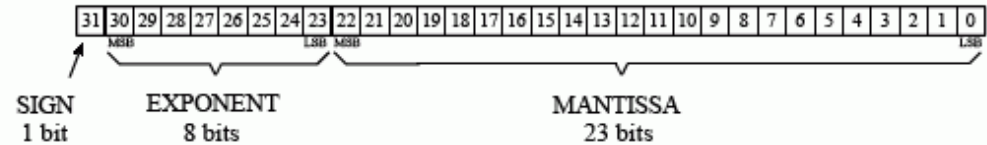
But in programmable logic Fixed Point Maths can be very fast



Floating Point Number

Floating point representation allows the decimal point to float to different places within the number depending upon the magnitude.

The floating-point number is standardized by an IEEE / ANSI Standard 754-1985 the basic IEEE floating point number



Example 1

0	00000111	110000000000000000000000	
↓	↓	↓	
+	7	0.75	$+ 1.75 \times 2^{(7-127)} = + 1.316554 \times 10^{-36}$

Example 2

1	10000001	011000000000000000000000	
↓	↓	↓	
-	129	0.375	$- 1.375 \times 2^{(129-127)} = - 5.500000$

FIGURE 4-2

Single precision floating point storage format. The 32 bits are broken into three separate parts, the sign bit, the exponent and the mantissa. Equations 4-1 and 4-2 shows how the represented number is found from these three parts. MSB and LSB refer to "most significant bit" and "least significant bit," respectively.

Why Fixed Point

Less complex to implement in logic

Enables a faster solution

Can be more power efficient

FPGA are register rich!

No standard for implementation but, Q format is popular

- Q15 – 15 fractional bits
- M,N – M integer bits and N fractional bits

Number Schemes

Fixed point numbers need to represent positive and negative numbers

- **Sign and Magnitude** - Utilises the left most bit to represent the sign of the number (0 = positive, 1 = negative) the remainder of the bits represent the magnitude. BUT! Positive and Negative Numbers.
- **Ones Complement** – Same unsigned representation for positive numbers as Sign and Magnitude representation. However, for negative numbers the inversion (ones complement) of the positive number are used. Requires end around carry for subtraction – Added complexity.
- **Twos Complement** – Positive are represented in the same manner as an unsigned numbers. While negative numbers are represented as the binary number you add to a positive number of the same magnitude to get zero

Twos Complement

A negative twos complement number is calculated by first taking the ones complement (inversion) of the positive number and then adding one to it. The twos complement number system allows subtraction of one number from another by performing an addition of the two numbers. The range a twos complement number can represent is given by:

$$(2^{n-1}) \text{ to } -(2^{n-1} - 1)$$

One method we can use to convert a number to its twos complement format is to work right to left leaving the number the same until the first one is encountered, after this each bit is inverted.

Fixed Point

How many bits do we need to represent my value ?

$$\text{Integer Bits Required} = \text{Ceil} \left(\frac{\text{LOG}_{10} \text{Integer_Maximum}}{\text{LOG}_{10} 2} \right)$$

For example, the number of integer bits required to represent a value between 0.0 and 423.0

$$9 = \text{Ceil} \left(\frac{\text{LOG}_{10} 423}{\text{LOG}_{10} 2} \right)$$

Fixed Point

How do we work out fractional bit ? Trade off between bit length and accuracy

To store the number $1.45309806319 \times 10^{-4}$

Multiply by 2^{16} $1.45309806319 \times 10^{-4} * 65536 = 9.523023$

Can only store 9 in the FPGA registers.

$9/65536 = 1.37329101563 \times 10^{-4}$

Significant loss of accuracy, how can we address this?

Fixed Point

We can obtain a more accurate result by scaling the number up by a factor of 2 that produces a result of between 32768 and 65535 therefore still allowing storage in a 16-bit number

$$268435456 * 1.45309806319 \times 10^{-4} = 39006.3041205$$

Stored number therefore $1.45308673382 \times 10^{-4}$ (39006/ 268435456)

Number is formatted as Q28 or 1,28

Fixed Point Rules

Fixed Point Arithmetic does have some rules which must be followed.

- **Addition** – Decimal points must be aligned
- **Subtraction** – Decimal points must be aligned
- **Division** – Decimal Points must be aligned
- **Multiplication** – Decimal points do not need to be aligned

Fixed Point Result Sizes

Operation	
$A + B$	Max(A'left, B'left) + 1 downto Min (A'right, B'right)
$A - B$	Max(A'left, B'left) + 1 downto Min (A'right, B'right)
$A * B$	(A'left + B'left) + 1 downto A'right + B'right
A / B - Unsigned	A'left - B'left downto (A'right + B'right) -1
A / B - Signed	(A'left - B'left) +1 downto A'right + B'right

Implementing in VHDL

Two options:

Numeric Standard (pre VHDL 2008)

- Unsigned
- Signed
- Need to keep track of decimal point – Range of number from X downto 0
- Quantisation needs to be performed by developer
- No in-built checking, requires design to check correct sizing / overflow etc

Fixed Point (VHDL 2008)

- Ufixed
- Sfixed
- Decimal point located between the 0 and -1 bit
- Inbuilt checking to ensure correct sizing of results

Fixed Package

Integer bits are represented in the range MSB down to 0

Fractional bits are represented in the range -1 down to LSB

SIGNAL example : `ufixed(2 DOWNT0 -3);`

Which represents the vector of 000.000 allowing for a range of 0.0 to 7.875

To help initialise signals, variables and constants in our algorithm we can use the `to_ufixed` and `to_sfixed`, these can be used with integers, real, `ufixed`, `sfixed` and `std_logic_vectors`.



Fixed Point Example



Implementing Complex Functions

What about more complex math

How would I implement the following functions

- Sine / Cosine / ArcTan
- SineH / CosH / ArcTanH
- Square Root
- Exponential
- Ln

Taylor / Maclaurin Series? Look Up Table ?

But how do we achieve performance ?

CORDIC Algorithm

CORDIC (COordinate Rotation Digital Computer) algorithm invented by Jack Volder for B58 Program

Deployed in first scientific calculator HP35

Shift and Add algorithm which can be used to implement transcendental functions.

No dedicated Multiplier required



CORDIC Algorithm

Three configurations - Linear /
Hyperbolic / Circular

Each mode has two modes – Rotation /
Vectoring

Enable a range of complex math's
functions to be implemented.

Configuration	Rotation	Vectoring
Linear	Op Y = X * Y	Op Z = X / Y
Hyperbolic	Op X = CosH(X) Op Y = SinH(Y)	Op Z = ArcTanH
Circular	Op X = Cos(X) Op Y = Sin(Y)	Op Z = ArcTan(Y) Op X = SQR(X ² + Y ²)

Additional Functions

$$\text{Tan} = \text{Sin} / \text{Cos}$$

$$\text{TanH} = \text{Sinh} / \text{Cosh}$$

$$\text{Exponential} = \text{Sinh} + \text{Cosh}$$

$$\text{Natural Logarithm} = 2 * \text{ArcTanH} \text{ note 1}$$

$$\text{SQR} = (X^2 - Y^2)^{0.5}$$

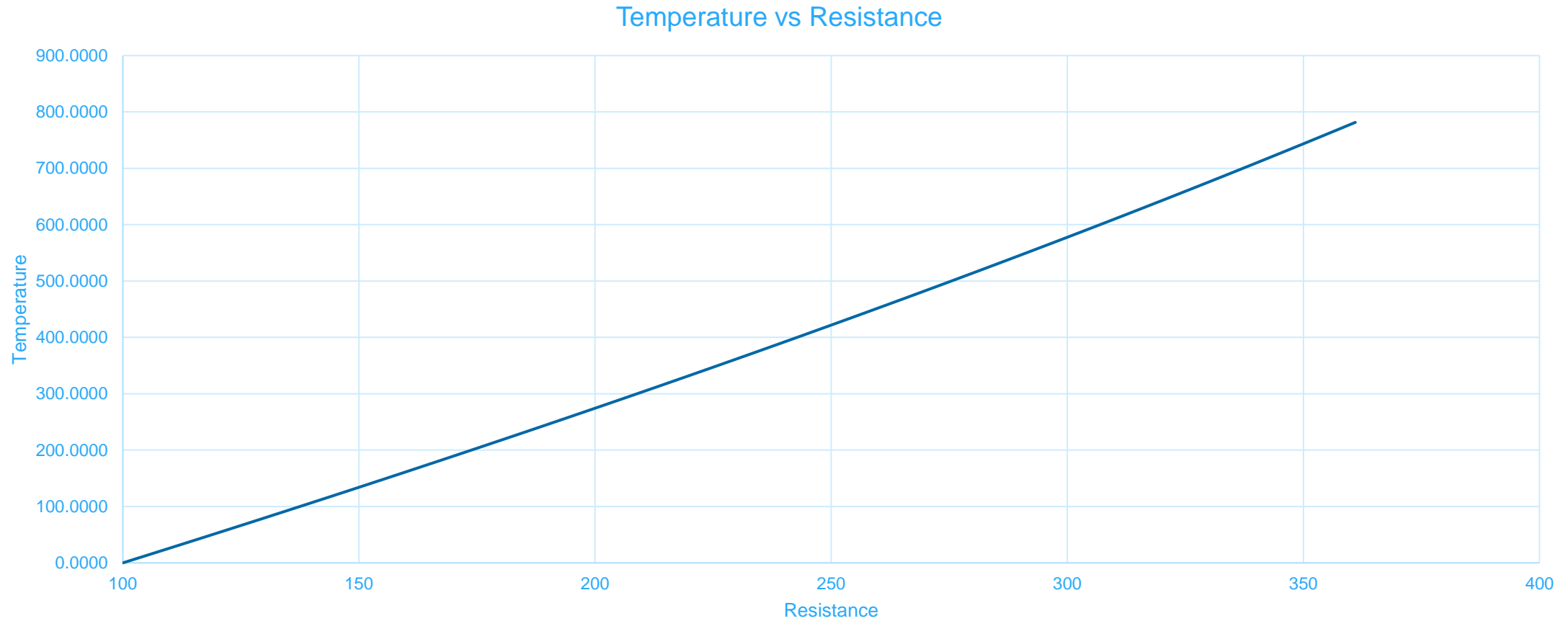
Complex Algorithm

How would you implement the following algorithm

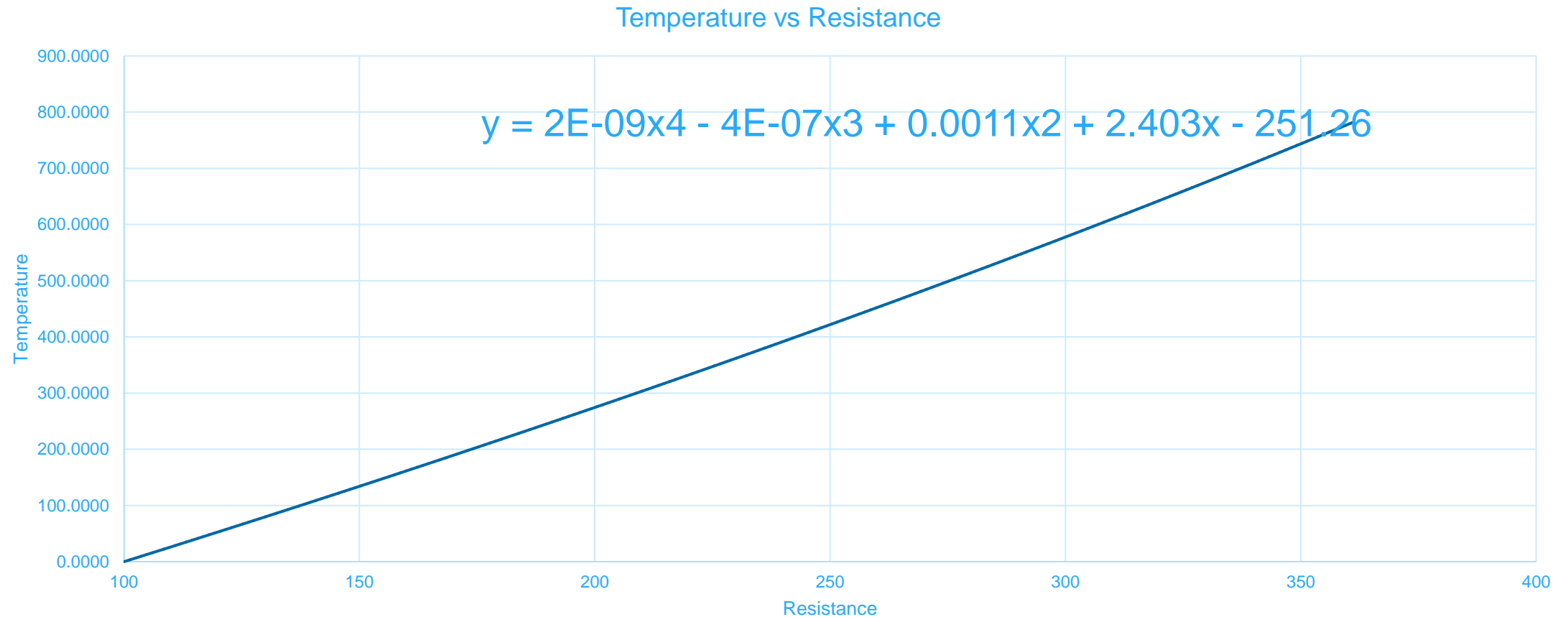
$$t = \frac{-R_0 \times a + \sqrt{R_0^2 \times a^2 - 4 \times R_0 \times b \times (R_0 - R)}}{2 \times R_0 \times b}$$

Typical Platinum Resistance Thermometer conversion equation used in industrial applications

Complex Algorithm



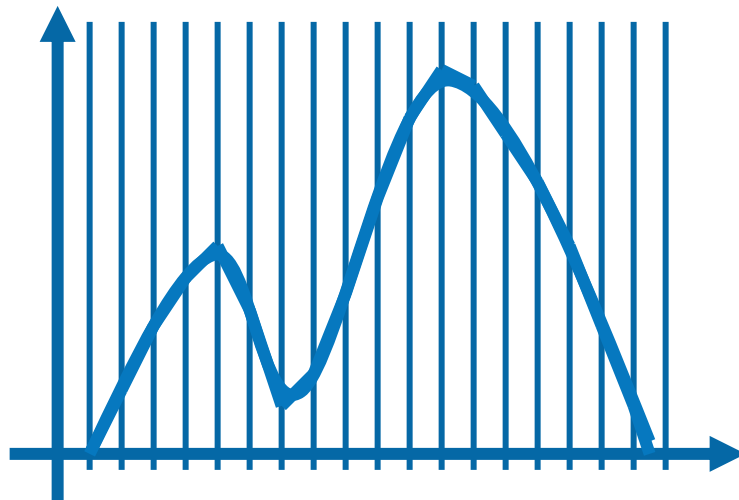
Polynomial Approximation



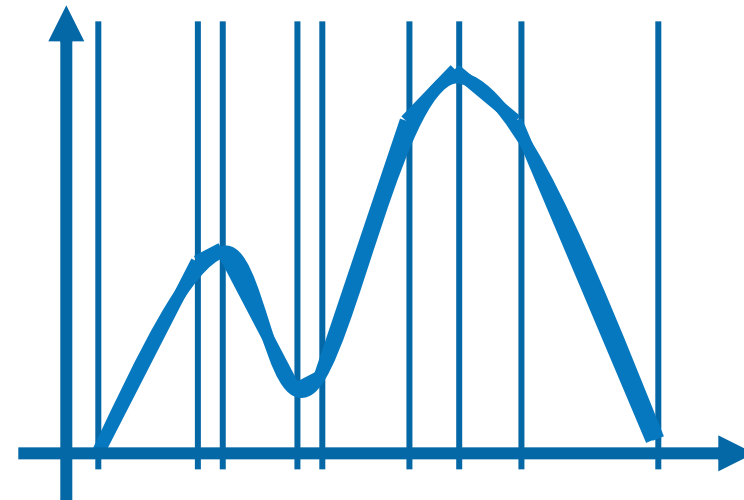
Polynomial Approximation

Leverage FPGA DSP rich environment – Addition and Multiplication easy to do in FPGA especially as we now know how!

If Accuracy is difficult with one overall polynomial equation – Segment it to several elements



Uniform segmentation



Non-uniform segmentation

Complex fixed example

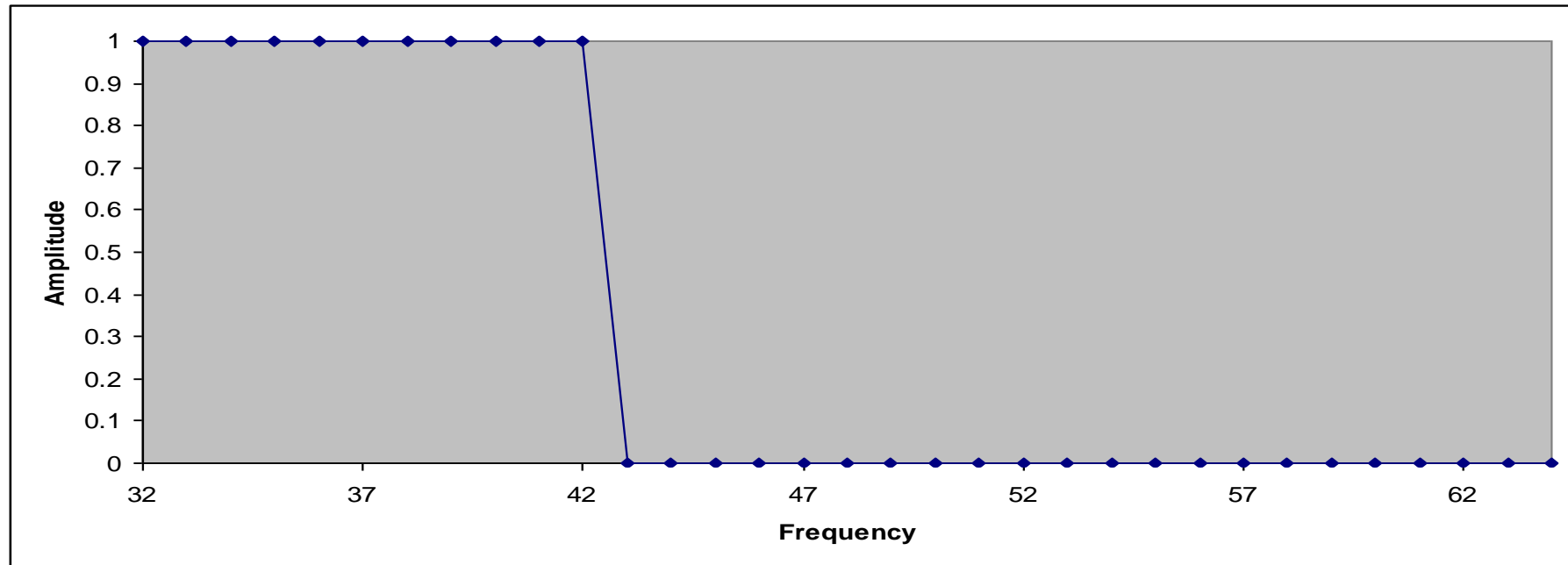


Filters

What about signal processing

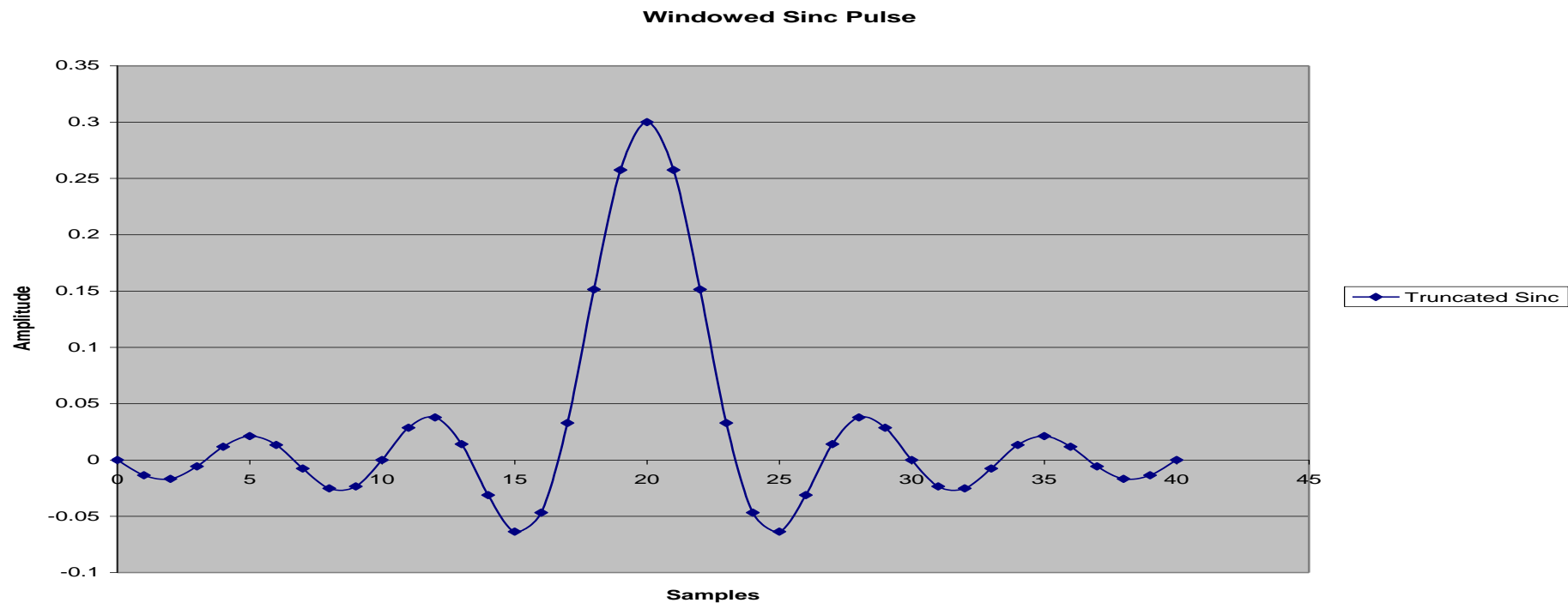
Finite Impulse Response Filters – Leverage the Multiple Accumulate Capability

Assume an ideal filter in the frequency domain – Brick Wall



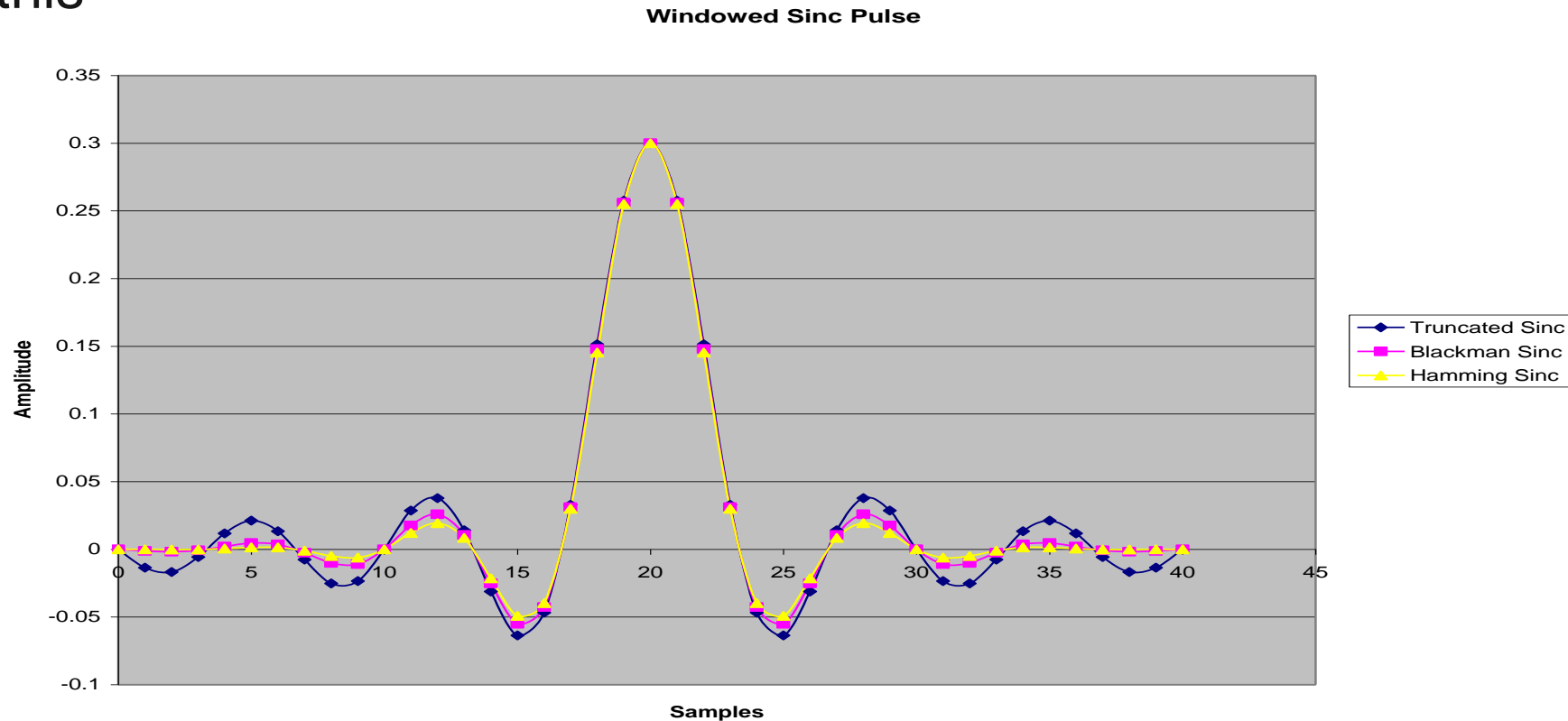
FIR Filter

IFFT of the brick wall filter gives us the Windowed Sync Pulse
The ripples extend to infinity and never settle to zero



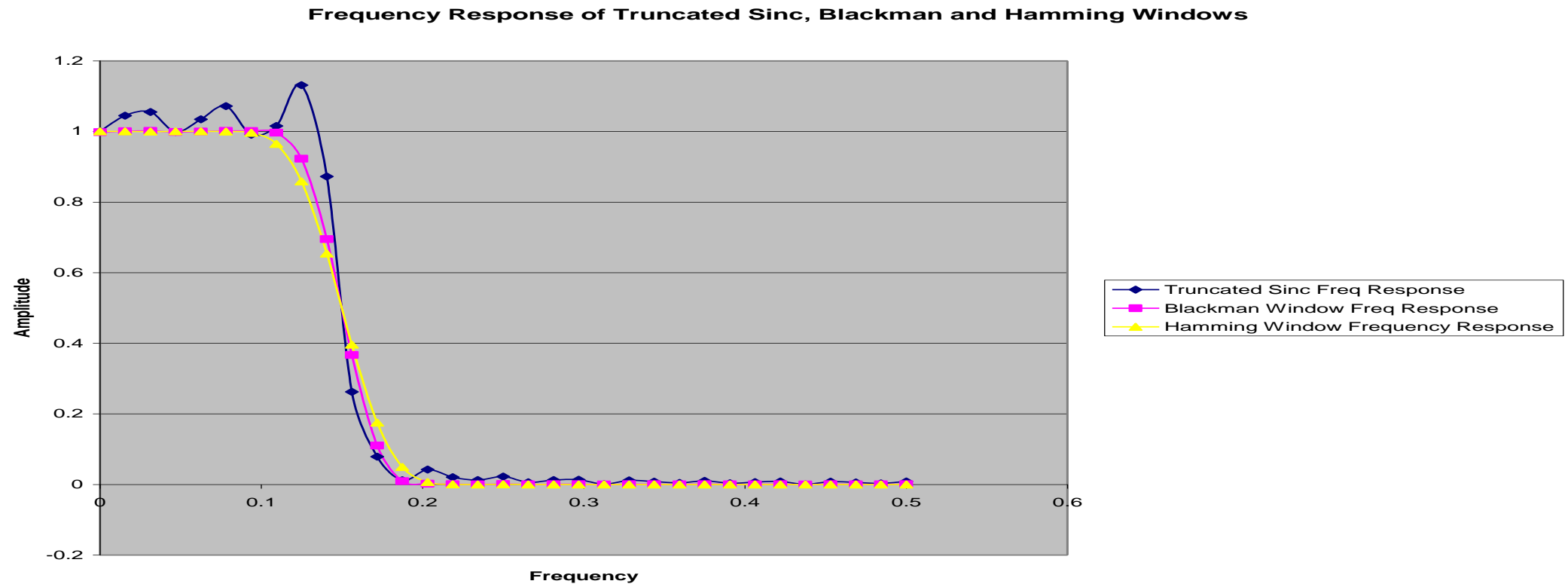
FIR Filter

Truncating the impulse response gives us ripples – Windowing helps address this

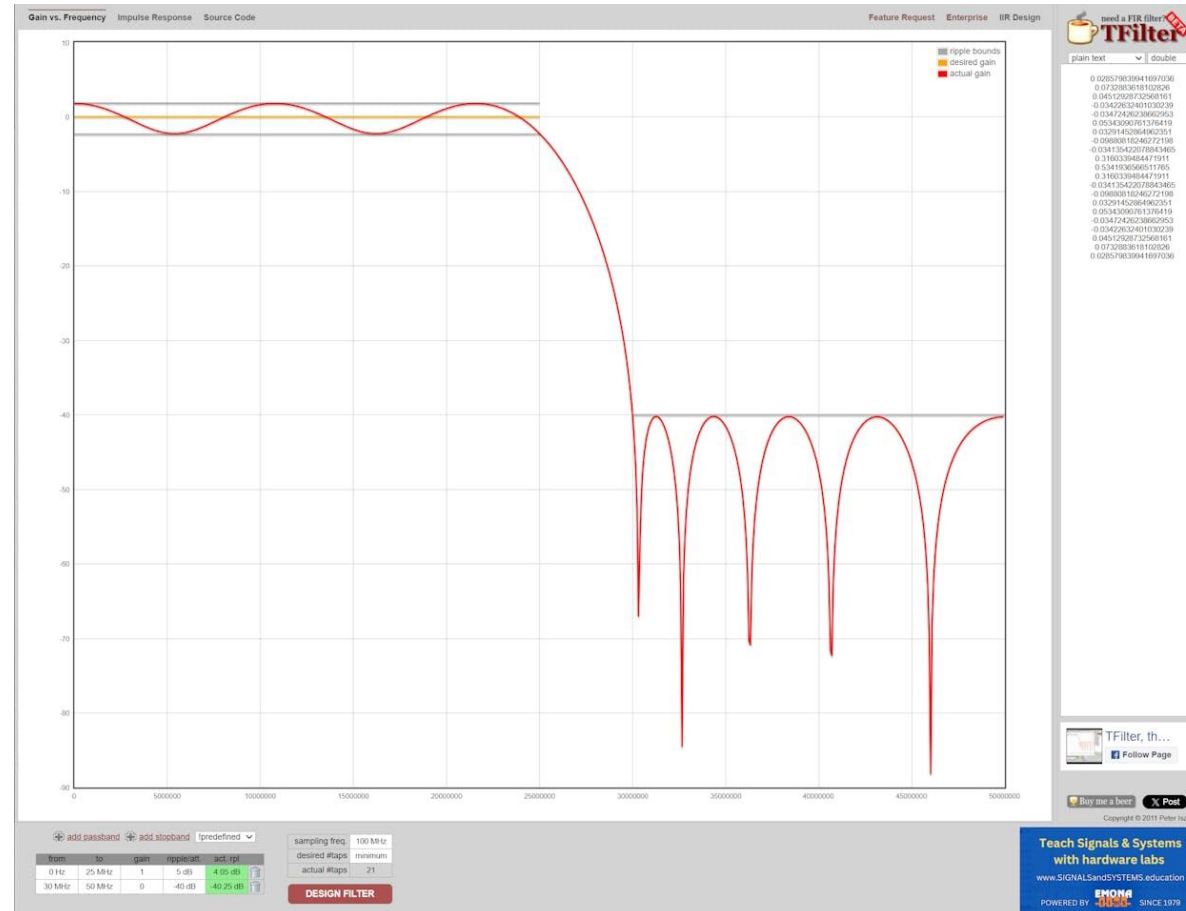


FIR Filter

Filter Response is improved with window



Filter Design Tools



Implementing

Re-customize IP

FIR Compiler (7.2)

Documentation IP Location

IP Symbol Freq. Response Implementation Details Coefficients

☐ Show disabled ports

Component Name: fir_compiler_0

Filter Options Channel Specification Implementation Detailed Implementation Interface Summary

Filter Coefficients

Select Source: Vector

Coefficient Vector: 937, 2402, 1479, -1122, -1138, 1751, 1079, -3238, -1119, 10356, 17504, 10356, -1119, -3238, 1079, 1751, -1138, -1122, 1479, 2402, 937

Coefficient File: no_coe_file_loaded

Number of Coefficient Sets: 1 [1 - 1024]

Number of Coefficients (per set): 21

☐ Use Reloadable Coefficients

Filter Specification

Filter Type: Single Rate

Inferred Coefficient Structure(s): Symmetric or Non Symmetric

Rate Change Type: Integer

Interpolation Rate Value: 1 [1 - 1]

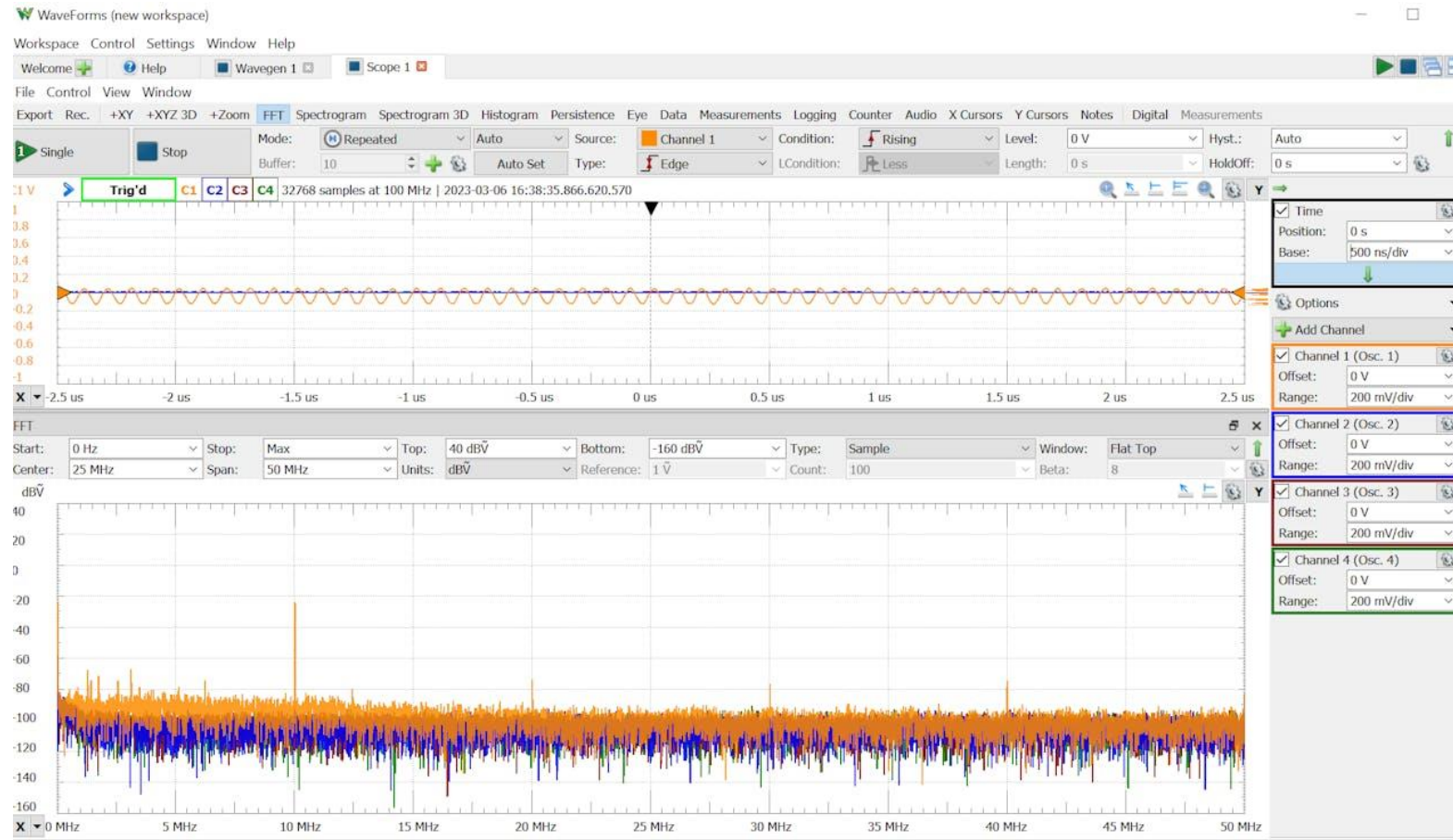
Decimation Rate Value: 1 [1 - 1]

Zero Pack Factor: 1 [1 - 1]

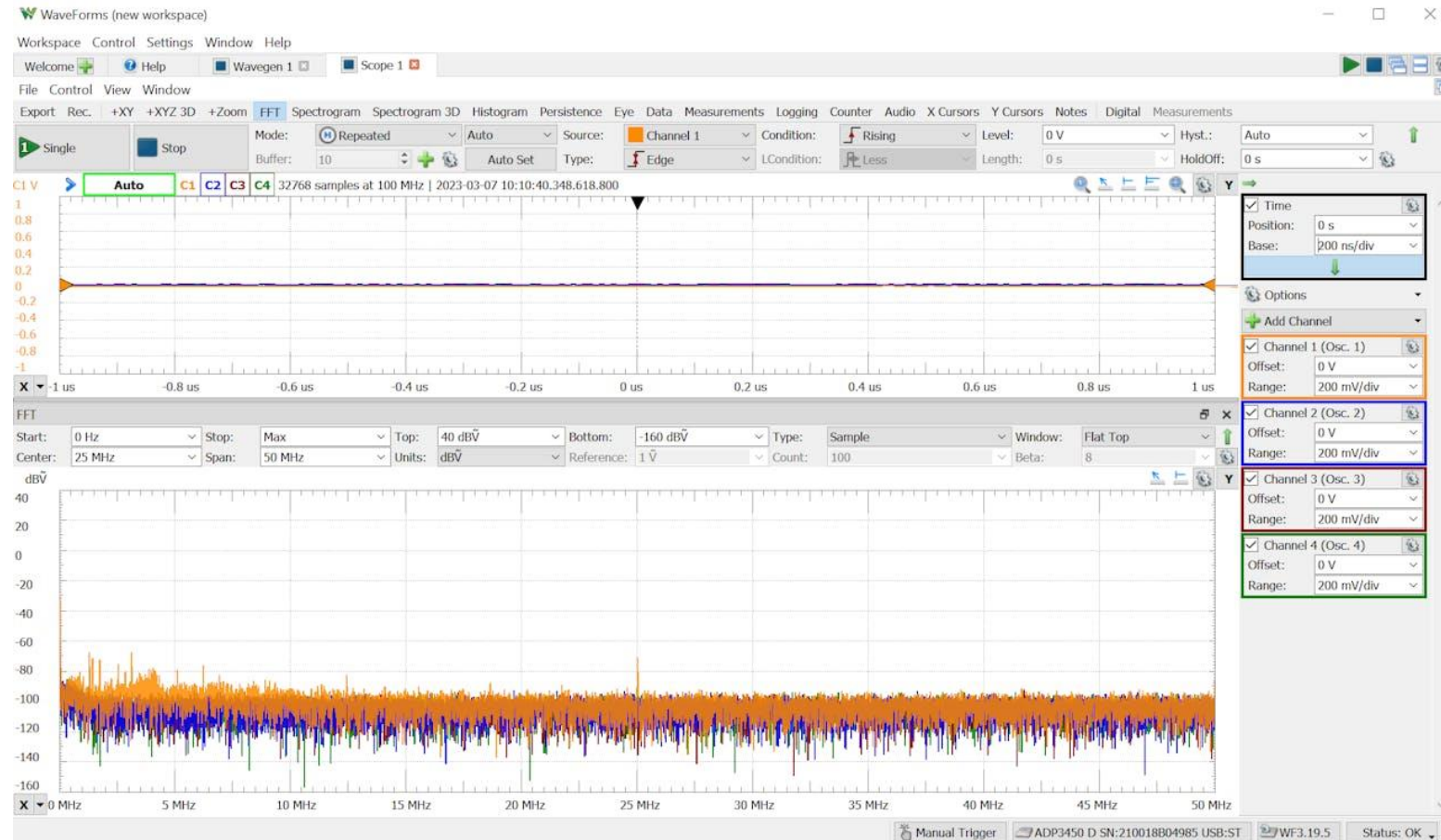
OK Cancel

+ S_AXIS_DATA
- aclk
M_AXIS_DATA +

Results 10MHz



Results 25 MHz





High Level Synthesis

What is HLS

High Level Synthesis (HLS) enables generation of RTL modules from higher level language such as C, C++, OpenCL

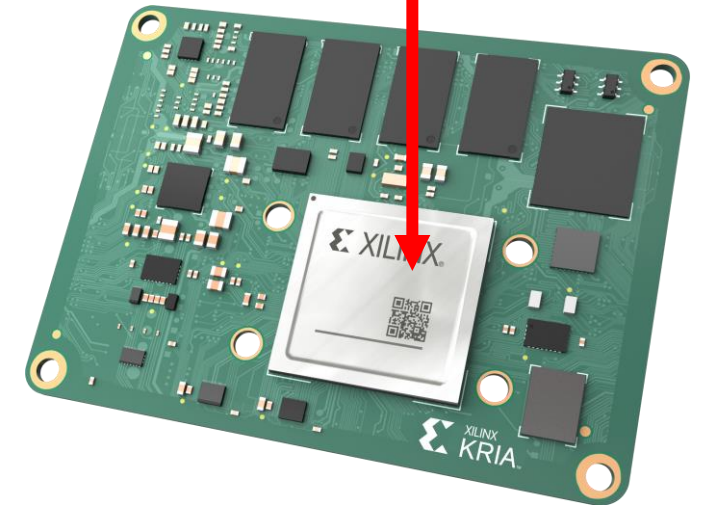
Of course, SW engineers still consider these low-level languages.

HLS offers several benefits for the development of Signal / Data / Image processing algorithms

```
error = set_point - sample;

p = error * KP;
i = i_prev + (error * ts * KI);
d = KD * ((error - error_prev) / ts);

op = p+i+d;
error_prev = error;
if (op > pmax) {
    i_prev = i_prev;
    op = pmax;
}else{
    i_prev = i;
}
return op;
}
```



HLS Benefits

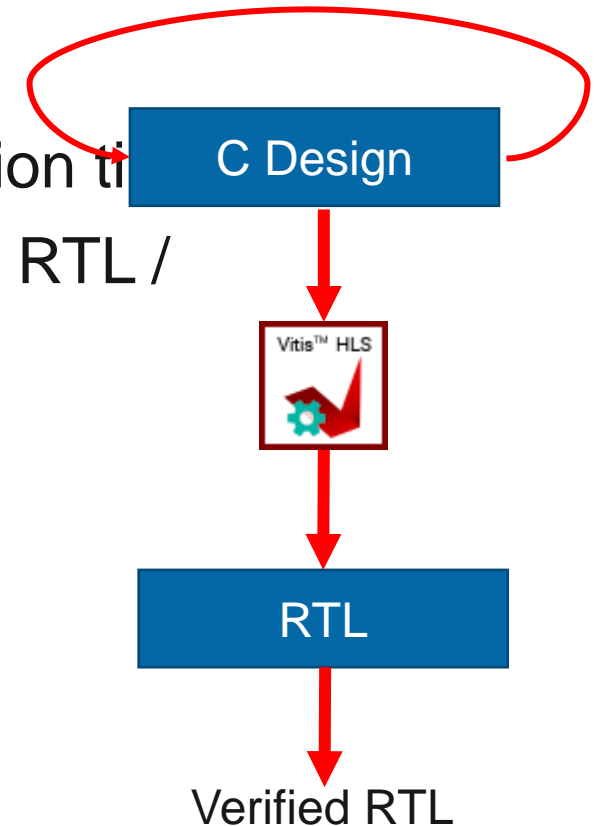
Developing in Higher Level language enables a faster iteration time

- Development Time decreased as untimed language – No RTL / Behavioral level
- Increased level of abstraction – accelerates development
- Verification time is reduced as untimed Simulation

Hours / Days Iteration



Seconds / Minutes Iteration

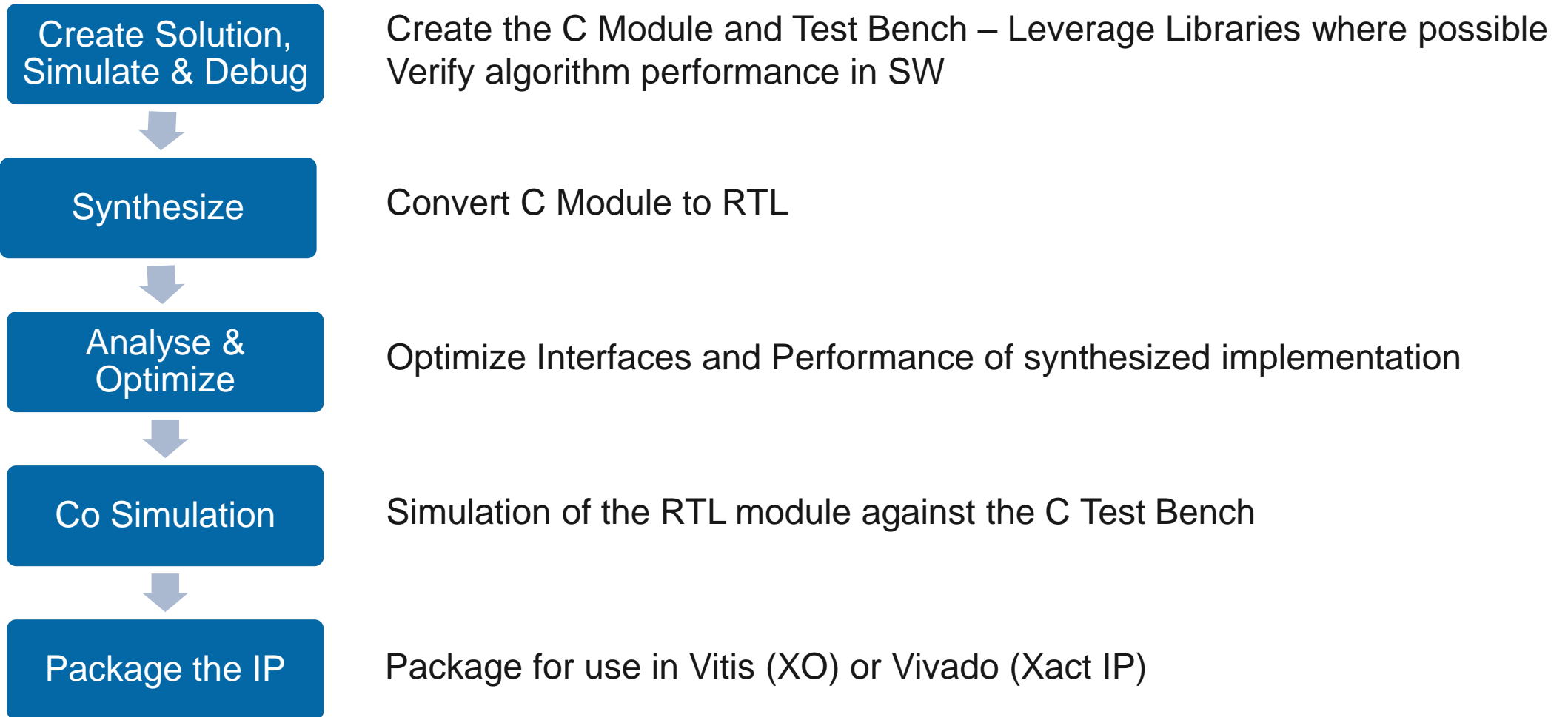


Creating HLS Solutions

Software written for CPUs and software written for FPGAs is fundamentally different

1. Not all C constructs can be synthesised
2. Learn about synthesizable C/C++ coding styles.
3. Need to focus on correct micro architecture
 1. Understand the producers and consumers
 2. Decompose the algorithm into small section which interconnect
 3. Understand throughput required for each element to achieve overall performance goals
4. Learn how to interpret the design reports

HLS Flow



Untimed to Timed

Scheduling

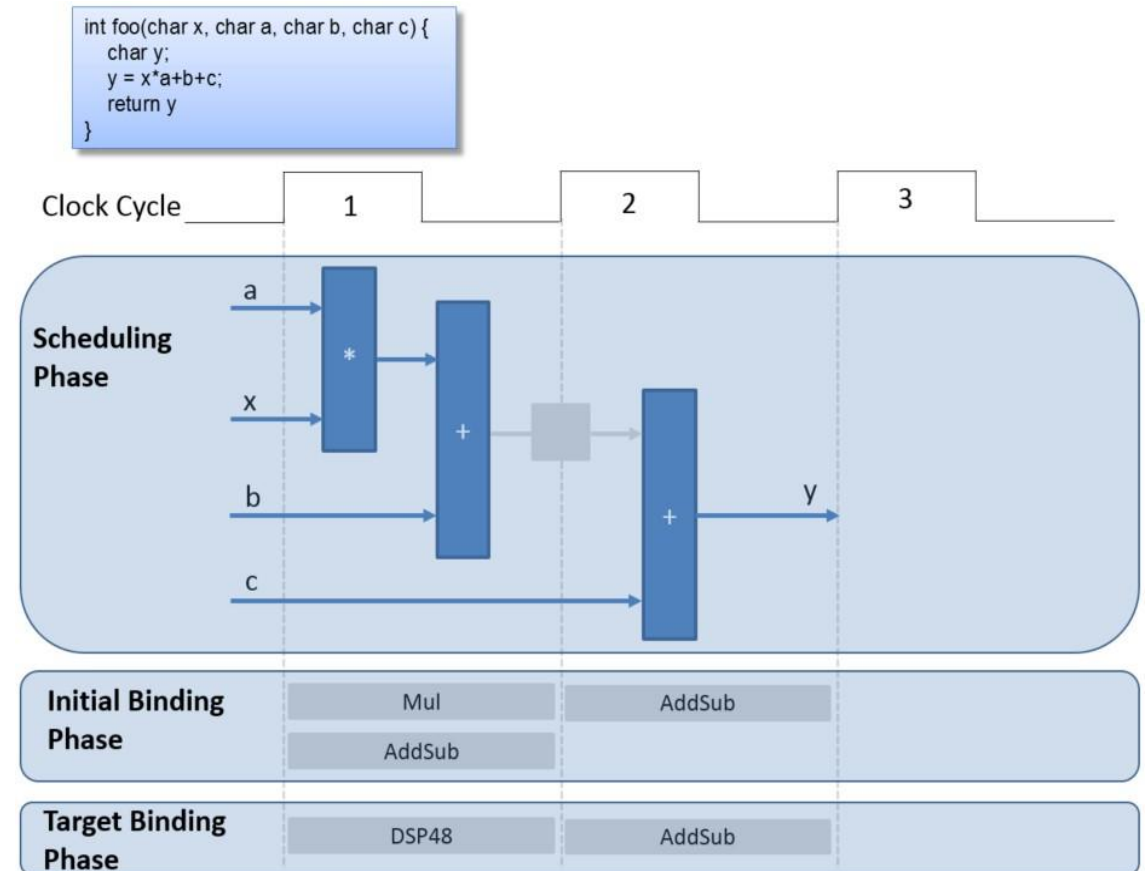
- Determines which operations occur during each clock cycle

Binding

- Determines which hardware resource implements each scheduled operation

Control logic extraction

- Extracts the control logic to create a finite state machine (FSM) that sequences the operations in the RTL design



HLS Math

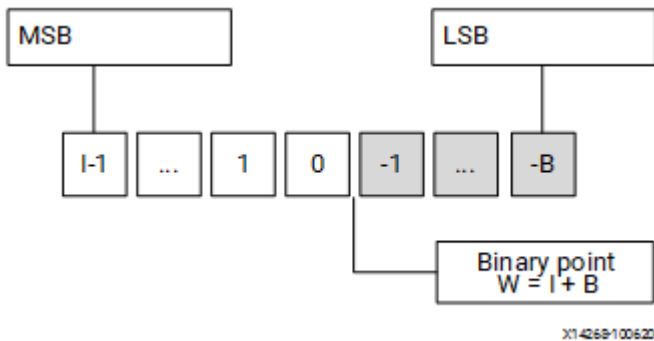
HLS is excellent for accelerating development time of algorithms.

As such there are several features we can leverage including

- Arbitrary precision data types – integer and fixed-point data types provided in `ap_[u]int` / `ap_[u]fixed`
- `HLS_Math.h` - library which provides functions implemented in `cmath`, provides floating- and fixed-point support.
- Domain specific HLS libraries including DSP, and Vision etc

HLS Math

ap_[u]fixed format



Identifier	Description																
W	Word length in bits																
I	<p>The number of bits used to represent the integer value, that is, the number of integer bits to the <i>left</i> of the binary point, including the sign bit.</p> <p>When I is negative, as shown in the example below, it represents the number of <i>implicit</i> sign bits (for signed representation), or the number of <i>implicit</i> zero bits (for unsigned representation) to the <i>right</i> of the binary point. For example,</p> <pre> ap_fixed<2, 0> a = -0.5; // a can be -0.5, ap_ufixed<1, 0> x = 0.5; // 1-bit representation. x can be 0 or 0.5 ap_ufixed<1, -1> y = 0.25; // 1-bit representation. y can be 0 or 0.25 const ap_fixed<1, -7> z = 1.0/256; // 1-bit representation for z = 2^-8 </pre>																
Q	<p>Quantization mode: This dictates the behavior when greater precision is generated than can be defined by smallest fractional bit in the variable used to store the result.</p> <table> <tr> <th>ap_fixed Types</th><th>Description</th></tr> <tr> <td>AP_RND</td><td>Round to plus infinity</td></tr> <tr> <td>AP_RND_ZERO</td><td>Round to zero</td></tr> <tr> <td>AP_RND_MIN_INF</td><td>Round to minus infinity</td></tr> <tr> <td>AP_RND_INF</td><td>Round to infinity</td></tr> <tr> <td>AP_RND_CONV</td><td>Convergent rounding</td></tr> <tr> <td>AP_TRN</td><td>Truncation to minus infinity (default)</td></tr> <tr> <td>AP_TRN_ZERO</td><td>Truncation to zero</td></tr> </table>	ap_fixed Types	Description	AP_RND	Round to plus infinity	AP_RND_ZERO	Round to zero	AP_RND_MIN_INF	Round to minus infinity	AP_RND_INF	Round to infinity	AP_RND_CONV	Convergent rounding	AP_TRN	Truncation to minus infinity (default)	AP_TRN_ZERO	Truncation to zero
ap_fixed Types	Description																
AP_RND	Round to plus infinity																
AP_RND_ZERO	Round to zero																
AP_RND_MIN_INF	Round to minus infinity																
AP_RND_INF	Round to infinity																
AP_RND_CONV	Convergent rounding																
AP_TRN	Truncation to minus infinity (default)																
AP_TRN_ZERO	Truncation to zero																
O	<p>Overflow mode: This dictates the behavior when the result of an operation exceeds the maximum (or minimum in the case of negative numbers) possible value that can be stored in the variable used to store the result.</p> <table> <tr> <th>ap_fixed Types</th><th>Description</th></tr> <tr> <td>AP_SAT¹</td><td>Saturation</td></tr> <tr> <td>AP_SAT_ZERO¹</td><td>Saturation to zero</td></tr> <tr> <td>AP_SAT_SYM¹</td><td>Symmetrical saturation</td></tr> <tr> <td>AP_WRAP</td><td>Wrap around (default)</td></tr> <tr> <td>AP_WRAP_SM</td><td>Sign magnitude wrap around</td></tr> </table>	ap_fixed Types	Description	AP_SAT ¹	Saturation	AP_SAT_ZERO ¹	Saturation to zero	AP_SAT_SYM ¹	Symmetrical saturation	AP_WRAP	Wrap around (default)	AP_WRAP_SM	Sign magnitude wrap around				
ap_fixed Types	Description																
AP_SAT ¹	Saturation																
AP_SAT_ZERO ¹	Saturation to zero																
AP_SAT_SYM ¹	Symmetrical saturation																
AP_WRAP	Wrap around (default)																
AP_WRAP_SM	Sign magnitude wrap around																
N	This defines the number of saturation bits in overflow wrap modes.																



HLS Example



Matlab Simulink

Matlab / Simulink Example

Example to convert voltage measured by negative Temperature Coefficient Thermistor into temperature.

Development will use Simulink HDL coder modules.

Two Stages

First determine resistance of NTC thermistor

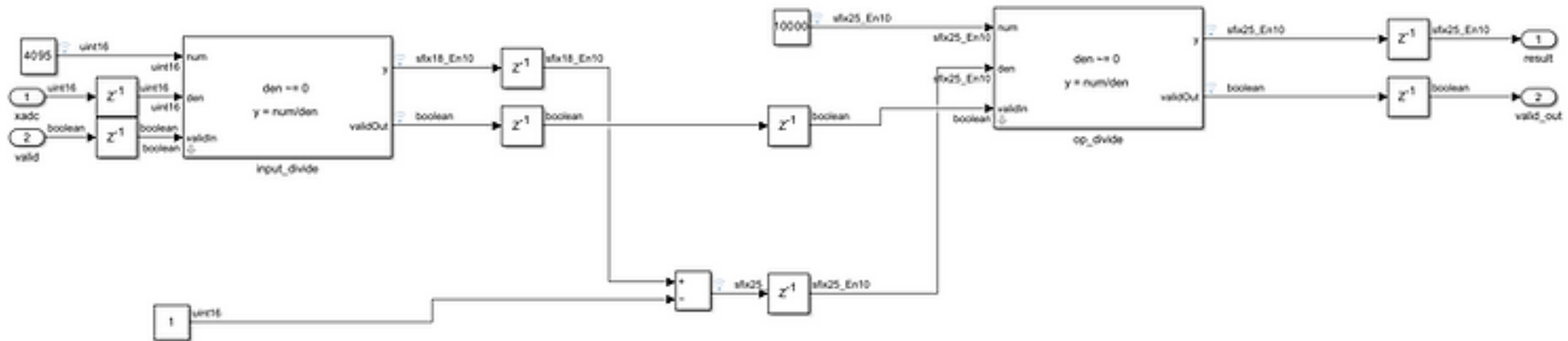
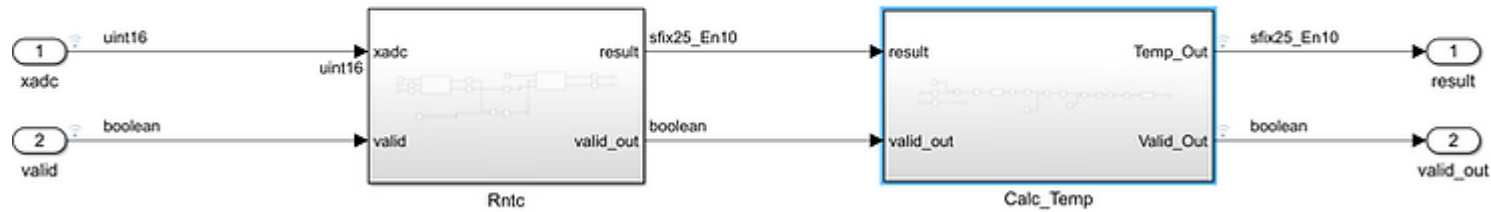
$$R_{ntc} = \frac{R_{series}}{\left(\frac{adc\ resolution}{adc\ value} - 1\right)}$$

Second determine the temperature

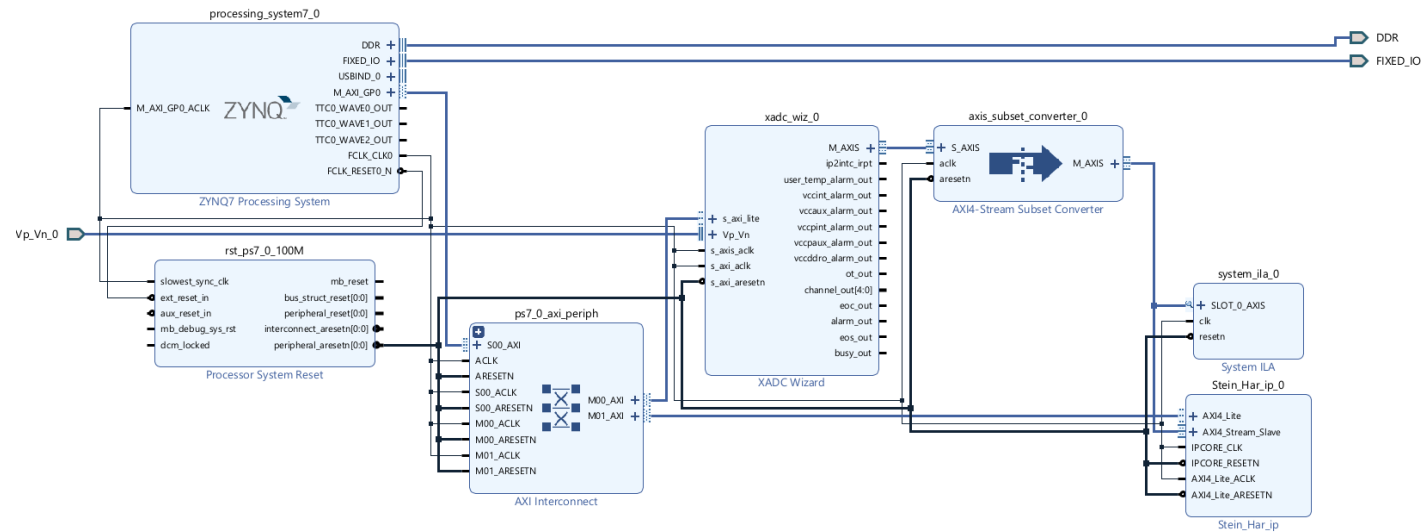
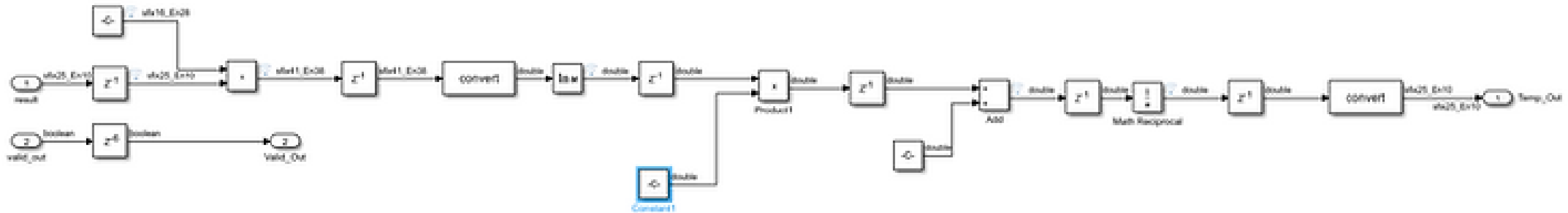
$$\frac{1}{T} = \frac{1}{T_0} + \frac{1}{Beta} * \ln \frac{R_{ntc}}{R_{nom}}$$

Matlab / Simulink Example

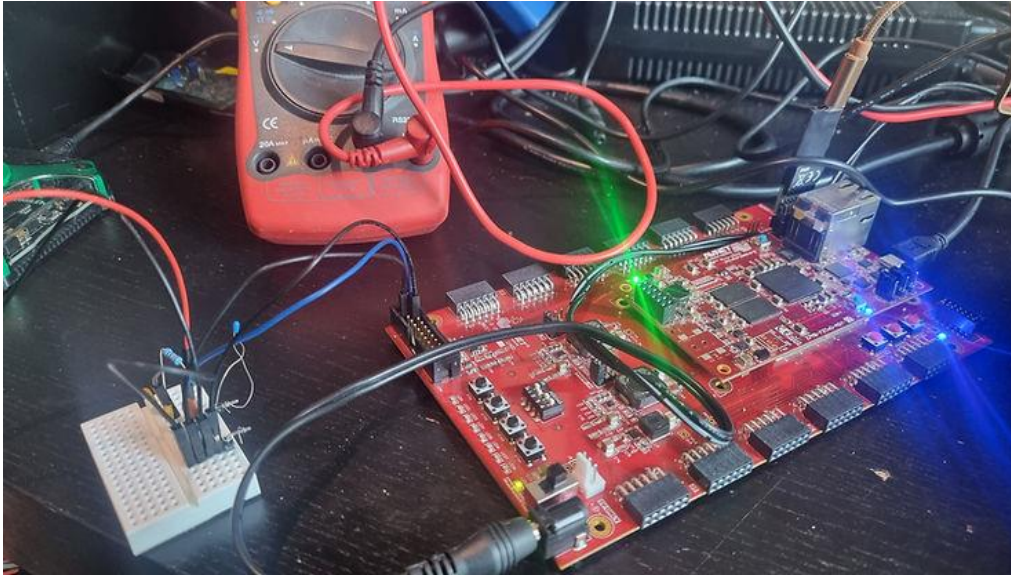
Create two modules as per the stages outlined previously



Matlab / Simulink Example



Matlab / Simulink Example



Hardware testing with NTC shows the functionality is as required the temperature is correctly shown in Kelvin

```
COM12 - Tera Term VT
File Edit Setup Control Window Help
295
002511
295
002635
295
009396
002
009934
002
010105
002
010293
003
010555
003
010275
003
009705
002
009377
002
009260
002
008867
001
008207
000
008086
000
007717
000
007489
000
007572
000
007469
000
```

Summary

- Implementing Math functions and algorithms in FPGA is not as challenging as it might appear.
- There are several different ways to implement algorithms as presented
 - HDL
 - HLS
 - Matlab / Simulink
- Each has their own pros/cons – HLS is excellent for rapid acceleration of algorithm development and aligns often with C based models



Questions



ADIUVO

ENGINEERING AND TRAINING, LTD.

www.adiuvoengineering.com



info@adiuvoengineering.com