

ODE Solver Design Reference

Team #5

April 22, 2020

Contents

I	Introduction	1
0.1	Abstract	2
0.2	Interfaces and Actions	2
0.2.1	Ports	3
0.3	Pipelining	5
0.3.1	Possible Pitfalls	5
0.4	Parallelism	5
0.5	Inter-module Communication	5
0.5.1	IO and CPU	5
0.5.2	IO and Solver/Interpolator	6
0.5.3	Solver and Interpolator	6
0.5.4	Communication and Clock Edges	6
II	Simulation Workflow	7
0.6	Input Preprocessing	8
0.7	Instantiating Main Module	9
0.8	Input Loading	9
0.9	Processing	10
0.10	Output Extraction	10
III	Specifications	11
0.11	Memory Mapping	12
0.11.1	Solver Memory Mapping	12
0.11.2	Interpolator Memory Mapping	12
0.12	Header Data Structure	13
0.13	Input Format	13
0.14	Modules	14

0.14.1	IO	14
0.14.2	Solver	17
0.14.3	Interpolator	19
0.14.4	Fixed/Floating Point Unit (FPU)	21
0.15	Compression	24
0.15.1	Input	25
0.15.2	Output	25
0.15.3	Example	25
0.15.4	Algorithm	25
0.16	Decompression	26

IV Bibliography

27

List of Figures

1	Overall System Design	2
2	I/O Internal Design	14
3	Solver Internal Design	16
4	Interpolator Internal Design	19
5	FP Unit	22

List of Tables

1	Solver Memory Mapping	12
2	Interpolator Memory Mapping	12
3	Header Data Structure	13
4	Inputs Specifications	13
5	Compression Output	25
6	Compression Example	25

Part I

Introduction

0.1 Abstract

This document is intended to be used as a comprehensive reference for the coming phases, and to be up-to-date with all future changes in the design.

0.2 Interfaces and Actions

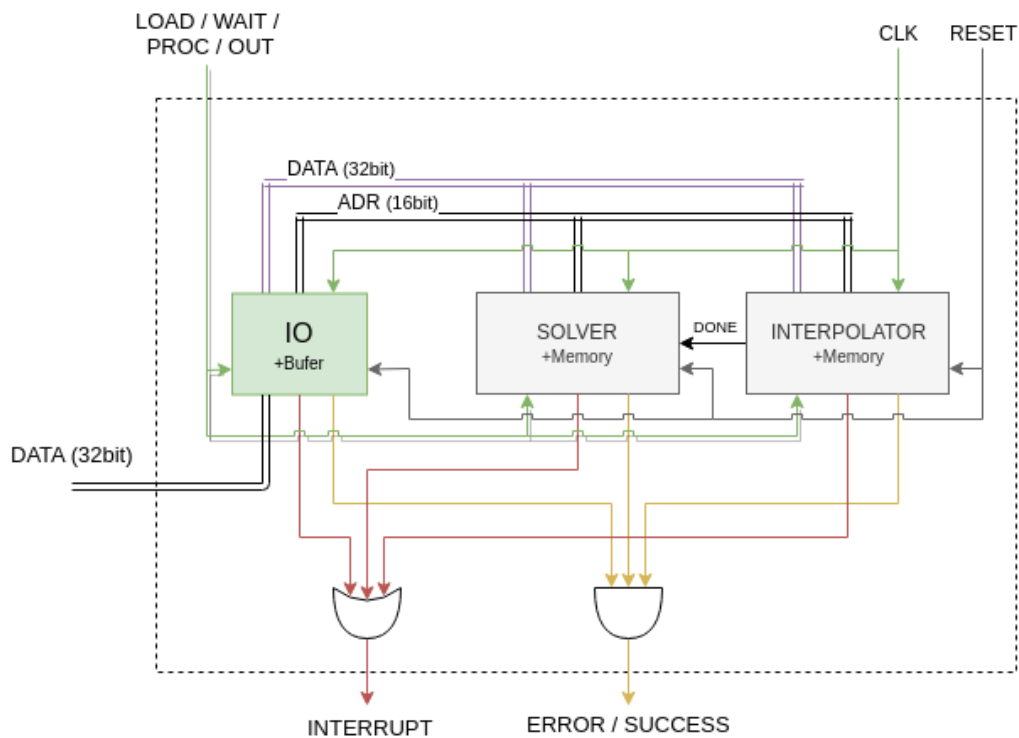


Figure 1: Overall System Design

The main module, as shown in Figure 1, has the ports listed in Subsection 0.2.1 that triggers some actions.

Ports and their actions are summarized below and detailed in the rest of the document.

0.2.1 Ports

CLK: IN

RESET: ASYNC IN

- Called at each time cpu wants to start/restart hardware.
- After enough time for reset, CPU is expected to turn the LOAD / WAIT / PROC / OUT into LOAD state and we start loading input.
- Clears all internal states of all modules:
 - IO internal buffer.
 - ERROR/SUCCESS of all modules resets to SUCCESS.
 - INTERRUPT resets to zero.
 - SOLVER and INTERPOLATORE clear all internal state and parallel-load their memories and latches with 0.
- Signal is async: not related to clock edges, could be called at any time and perform its action immediately with only the propagation delay.

LOAD / WAIT / PROC / OUT (2bit): IN

Set the current major state of the machine.

- LOAD(0):
 - IO receives **compressed** data from the CPU.
 - IO decompresses data into buffer.
 - Buffer is flushed into data bus with appropriate address.
 - Ends when cpu finishes its data loading and switches to WAIT state.
- WAIT(1):
 - Same state as LOAD, but IO doesn't receive data anymore from CPU.
 - Ends when IO flushes all its buffer and raises INTERRUPT with either ERROR or SUCCESS.

- PROC(2):
 - SOLVER sends time step to calculate U at.
 - SOLVER and INTERPOLATOR work concurrently to calculate their outputs.
 - INTERPOLATOR sends DONE signal to SOLVER when it finishes the interpolated U.
 - SOLVER can request to copy the interpolated U.
 - INTERPOLATOR waits for SOLVER to send next time step.
 - Ends when either SOLVER or INTERP raises INTERRUPT with either SUCCESS or ERROR.
- OUT(3):
 - IO connects internal bus with external bus.
 - SOLVER starts copying it's Xout, word each cycle.
 - Ends when SOLVER raises INTERRUPT with either SUCCESS or ERROR.

DATA (32bit): INOUT

- Data bus between cpu and io.

INTERRUPT: OUT

- IO or SOLVER or INTERPOLATOR set this signal when one/more finish their task.
- If task finished with success, the ERROR / SUCCESS should be set to SUCCESS, otherwise ERROR.

ERROR(0) / SUCCESS(1): OUT

- CPU should operate on this value ONLY when INTERRUPT is 1.
- Errors that could happen include: divide by zero, $H \neq 1$, incomplete input.

0.3 Pipelining

Solver's main equation to calculate is $X_{n+1} = X_n + hAX_n + hBU_n$, so we can divide it into two parts, part 1: $X_n + hAX_n$, and part 2: hBU_n , so our approach is that at the begining of operating that equation Solver sends 'h' to the interpolator, at which it (interpolator) will calculate U, meanwhile, solver is calculcating the first part (part 1), and when interpolator is finished, it sends 'done' signal to solver and raises the value at the bus, solver then proceed with the second part of the equation (part 2).

0.3.1 Possible Pitfalls

That whole assumption is built based on that part 1 of the solver's main equation is polynomial wrt. N (no. of elements at X), while the interpolator's whole equation is Linear in (M), that's why most of the time no delay will occur or be noticeable.

0.4 Parallelism

Parallelism: having multiple instances of module M as $m_0, m_1, m_2 \dots$, and splitting inputs $x_0, x_1, x_2 \dots$ between them, so each module m_i takes input x_i .

Based on the previous definition, this design doesn't implement parallelism. As the solver only solves one and only one X_{out} at a given time. And the main module has one instance of the solver and interpolator.

0.5 Inter-module Communication

0.5.1 IO and CPU

- CPU puts main module in LOAD state, IO reads 32bits of compressed data.
- When CPU finishes, it raises WAIT.
- IO flushes its buffer, and if succeeded, sets INTERRUPT with ERROR/SUCCESS.

- CPU switches to PROC state, so that all of other components starts calculating.
- When output data is ready, CPU is interrupted with error or success. In case of success, CPU switches to OUT state which implies that it's ready to take out the resulting output from IO.

0.5.2 IO and Solver/Interpolator

- LOAD/WAIT state: IO puts the data and address on their busses, each component of them two reads that address and check if this address is within its range, if so this component reads the data and stores it.
- OUT state: IO connects internal and external busses. for each cycle: SOLVER puts the next X_{out} value.

0.5.3 Solver and Interpolator

- Solver is working on X_h so it needs U, solver checks on DONE signal, if raised, reads the U from the data bus, Interpolator has already placed it there.
- Then Solver places h_{new} at the data bus, interpolator reads it and they both work...
- and repeat.

0.5.4 Communication and Clock Edges

Module X sends data and address on **only positive edge** of the clock, and module Y receives data on **either the following positive or negative edge** of the clock.

Part II

Simulation Workflow

0.6 Input Preprocessing

This stage is the responsibility of `preprocessor.py` that runs before the simulation:

Stage Input

JSON file that follows the format stated in the main project document [1].

Steps

- Create bit stream of the read data that follows the Input Data Structure specifications.

NOTE: if the mode is 16 bit (fixed), then each 16bit number must be padded with 16 bits of zeros, so that the number is in the lowest part (bits 0 ... 15 \rightarrow num, 16 ... 31 \rightarrow 0 padding.)

- Encode the bits following the Compression specifications.
- Collect encoding output in ASCII string, each byte in string is either '0' or '1' in ASCII format.
- When the string reaches the length of 32 bytes, push it to output file.
- If the last created string didn't reach the length of 32 bytes, complete the rest with '0' and push it to the output file.

Stage Output

- ASCII file that contains multiple lines of compressed data.
- each line has exactly 32 '0' or '1' ASCII characters.
- ONLY the ASCII characters 0 or 1 are permitted in the file and NOTHING ELSE.
- there is NO EMPTY LINE/s in the file or spaces.

0.7 Instantiating Main Module

This stage and all the next ones are the responsibility of the CPU simulation code.

CPU is a non-synthesisable HDL test-bench that:

- Instantiates the HW main module.
- Attaches the appropriate signals to the HW main module.
- Generates CLK with fixed frequency.
- Resets HW for 1 ps.
- Loads data into HW.
- Puts HW into PROCESS state.
- Load output out from the HW and into a file.

0.8 Input Loading

- Load the output of the former script into array of vectors each is 32bit wide that will hold one line in the file.
- RESET for 1 ps.
- Put HW at LOAD state.
- For each 32bit vector in the former array:
 - At the positive edge of CLK:
 - * Load vector into DATA bus.
- Load DATA with 0s.
- Wait for the positive edge of INTERRUPT signal.
- Check for ERROR / SUCCESS and only proceed if it is SUCCESS.

0.9 Processing

- Put HW at PROCESS state.
- Wait for the INTERRUPT positive edge.
- Check for ERROR / SUCCESS and only proceed if it is SUCCESS.

0.10 Output Extraction

- Put high impedance on DATA bus.
- Put HW at OUT state.
- Keep receiving data into array of vectors and outputting them into file in the same format of the input file.
- Wait for the positive edge of INTERRUPT signal.

Simulation is done!

You can turn the output into human-readable json using `output_formatting.py` script.

Part III

Specifications

0.11 Memory Mapping

The Following addresses are only meant for internal communicating between modules, and they don't need to resemble actual addresses stored at some memory.

The address loaded at the bus resembles what kind of data is on data bus.

If address bus is loaded with an address Adr that some module M is not assigned to, module M must ignore the data and address bus so the rest can use the bus.

This way communication is simplified.

0.11.1 Solver Memory Mapping

Solver module listens at the addresses listed in Table 1.

Table 1: Solver Memory Mapping

Address	Type	Words	Name	Description
0x0000	struct Header	1	Header	Includes Dimensions and modes
0x0001	f64	2	H	Timestep (variable step mode)
0x0003	f64	2	Error	Error Tolerance (variable step mode)
0x0005	f64[50][50]	5000	A	Matrix A
0x138D	f64[50][50]	5000	B	Matrix B
0x2715	f64[50]	100	X	Initial value of X
0x296D	f64[50]	100	U_0	Initial U vector

0.11.2 Interpolator Memory Mapping

Interpolator module listens at the addresses listed in Table 2.

Table 2: Interpolator Memory Mapping

Address	Type	Words	Name	Description
0x0000	struct Header	1	Header	Includes Dimensions and modes
0x296D	f64[50]	100	U_0	Initial U vector
0x29D1	f64[5]	10	T	Time points of solutions
0x29DB	f64[50][5]	500	U_s	U vector at required time steps
0x2C33	f64	2	h_{new}	Time to calculate U at
0x2C35	f64	2	h_{init}	Detect upcoming output point and switch ranges

0.12 Header Data Structure

Table 3 shows the structure of header data. Note the whole struct takes 2 words.

Table 3: Header Data Structure

Address	Bit	Name	Description	Datatype	Total Size
0x0000	31:26	N	Dimension of X	uint	6 bits
	25:20	M	Dimension of U	uint	6 bits
	19	Solver Mode	Fixed Step(0) or Variable Step(1)	enum	1 bit
	18:17	FPU Precision	fixed point(0), f64(1) or f32(2)	enum	2 bits
	16:14	T_{size}	Count of timesteps needed	uint	3 bits
	13:0	NOT USED	——	——	14 bits

0.13 Input Format

Table 4 shows the IO input structure. IO reads the header and figures out the addresses.

NOTE: if the mode is 16 bit (fixed), then each 16bit number is padded with 16 bits of zeros, so that the number is in the lowest part (bits 0 ... 15 \rightarrow num, 16 ... 31 \rightarrow 0 padding.)

Table 4: Inputs Specifications

Address	Type	Max Words	Name
0x0000	struct Header	1	Header
0x0001	f64	2	H
0x0003	f64	2	Error
0x0005	f64[50][50]	5000	A
0x138D	f64[50][50]	5000	B
0x2715	f64[50]	100	X
0x296D	f64[50]	100	U_0
0x29D1	f64[5]	10	T
0x29DB	f64[50][5]	500	U_s

0.14 Modules

0.14.1 IO

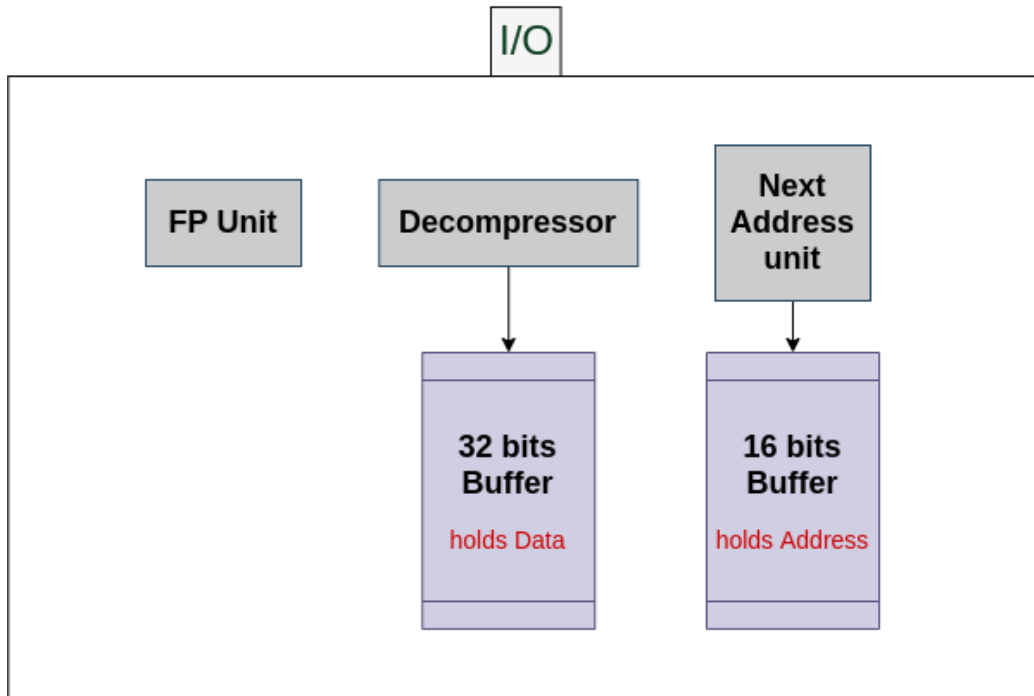


Figure 2: I/O Internal Design

Role

- In LOAD/WAIT state:
 - Receive packets of 32 bits from the CPU through *DATA* bus.
 - Decompress the data.
 - Send the decompressed data to the data-bus.
- In OUT state: just connect the internal data bus with external data bus.

Ports

- INOUT: 32bit data bus with other modules.
- INOUT: 32bit data bus with CPU.
- OUT: 16bit address bus.
- IN: CLOCK.
- ASYNC IN: Reset.
- IN: 2bit Load/Wait/Process/Out.
- OUT: Interrupt to CPU.
- OUT: Error/Success to CPU.

Sub-modules

- Decompressor: For decompression, see section 0.16.
- Next Address Unit: Calculates the next address to load at the address bus.

NAU has the header (the first 32 input data) that contains N, M and T_{size} . When reading data (e.g. matrix A) it calculates where it will end and proceeds with the next address.

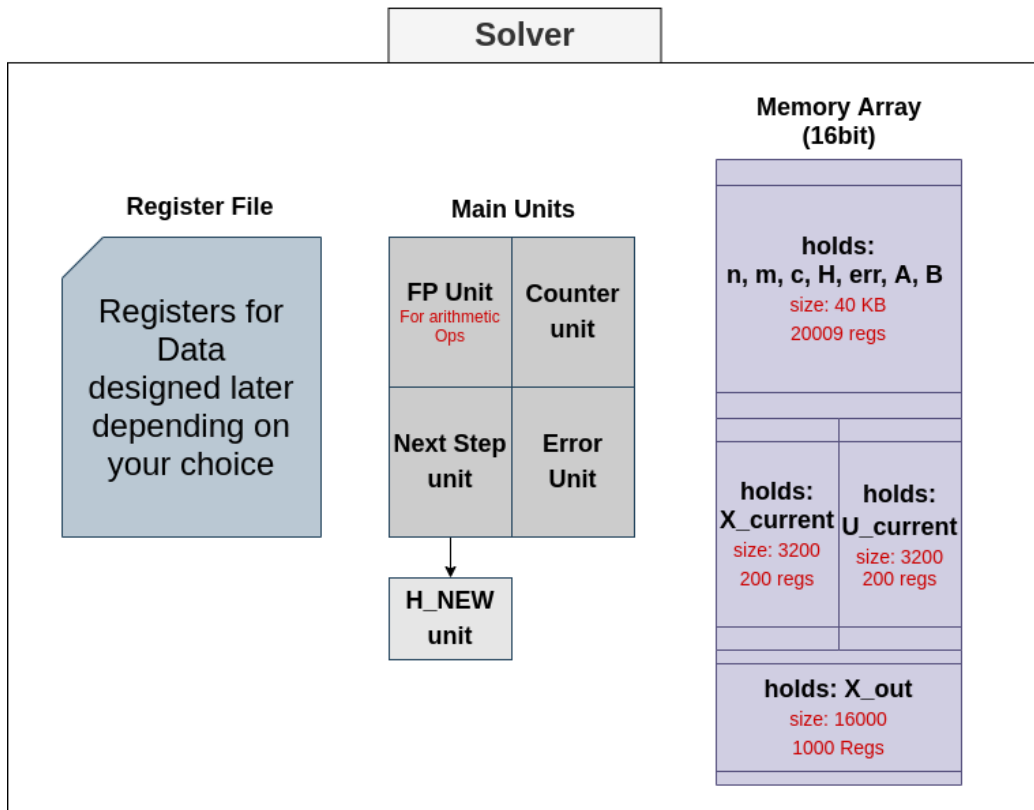


Figure 3: Solver Internal Design

0.14.2 Solver

Role

Receives U_h from interpolator, gives it another h to compute $U_{h_{new}}$ at, then computes X_h and decides to stop and flush output to I/O or continue.

At the beginning it receives its data from I/O such as $N, M, err, h \dots$ etc.

- Computes the upcoming X knowing h, the previous X and U.
- Counts the error difference and the new h.
- Checks for arithmetic errors that may occurs (e.g. div. by zero).
- Outputs the final X's at the desired times to the internal RAM.

Ports

- IN: Done signal from interpolator.
- INOUT: 32bit data bus with other modules.
- INOUT: 16bit address bus.
- IN: CLOCK.
- ASYNC IN: Reset.
- IN: 2bit Load/Wait/Process/Out.
- OUT: Interrupt to CPU.
- OUT: Error/Success to CPU.

Solver Memory

- Main Part: 40 KB \rightarrow 20009 (16 bits) Registers
 - N, M, C = 16 bits
 - h = 64 bits
 - err = 64 bits
 - A = $[50*50]*64 = 160000$ bits

- $B = [50*50]*64 = 160000$ bits
- $X_{current}$: 3200 bits \rightarrow 200 (16 bits) Registers
 - $X = 50*64 = 3200$ bits
 - 50: max of M
 - 64: max of numbers
- $U_{current}$: 3200 bits \rightarrow 200 (16 bits) Registers
 - $U = 50*64 = 3200$ bits
 - 50: max of M
 - 64: max of numbers
- X_{out} : 16000 bits \rightarrow 1000 (16 bits) Registers
 - $X = 5 * 50 * 64 = 16000$ bits
 - 50: max of M
 - 64: max of numbers
 - 5: max of times answer is required

Note: main module stores all the 5 output values and flushes them to the CPU all at once at the end.

Sub-modules

- FPU: see Section 0.14.4.
- Error Unit: to detect any error in sizes, h, numbers ... etc.
- Next Step Unit: helps create the upcoming h_{new} so that when solver is busy calculating X_h , interpolator is calculating $U_{h_{new}}$ ¹, this unit represents the stepper unit, holds the logic of calculating the adaptive h , and detects when to stop, in summation it calculates the next h , even if it was fixed step.
- Counter Unit: tells you when to calculate more, when to advance to next time (in T_s), and when to stop the whole operations.

¹See Section 0.3 to learn more about pipelining in design.

0.14.3 Interpolator

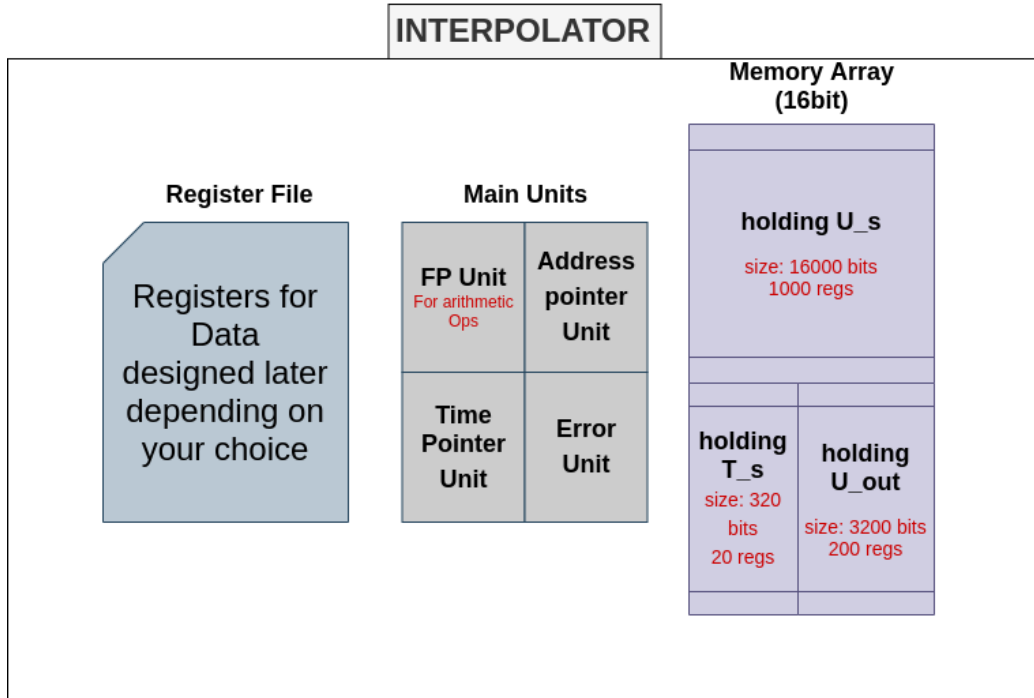


Figure 4: Interpolator Internal Design

Interpolator Job Description

- Only computes U at a specific time.
- At the beginning it receives from IO the U_s and T_s fixed variables/data, and stores them.
- Each U_s is at most of size $[64 \times 50] = 3200$ bits = 200 registers
- Each T_s is at most of size $[64]$ bits = 4 regs.
- T_s keeps hold of the time where each U_s represents, for example, $T_s = [1, 2, 3]$, there fore the first 200 regs. In U_s are the value of U at time 1, and so on...

- You will need an iterator to target the appropriate U_s from the array of memory, that's the **Address pointer Unit**.
- Time Pointer Unit: Identifies which T of the T_s array the unit is currently handling.
- Note: at the beginning of the program $T_{init} = 0$, $T_{final} = T_s[0]$, after a successful output, $T_{init} = T_{final}$, and $T_{final} = T_s[1]$, and so on...

Role

- Calculates the upcoming U knowing h , U initial and U final.

Ports

- OUT: Done signal to Solver with one more bit to indicate whether the calculated U_{int} is an output point.
- INOUT: 32bit data bus with other modules.
- IN: 16bit address bus.
- IN: CLOCK.
- ASYNC IN: Reset.
- IN: 2bit Load/Wait/Process/Out.
- OUT: Interrupt to CPU.
- OUT: Error/Success to CPU.

Interpolator Memory

- U_s : 16000 bits \rightarrow 1000 (16 bits) Registers
 - $U_s = 5 * [50] * 64 = 16000$ bits
 - 50: max of M
 - 64: max of numbers
 - 5: max of times answer is required

- T_s : 320 bits \rightarrow 20 (16 bits) Registers
 - $T_s = 5 \cdot 64 = 320$ bits
 - 5: max of times answer is required
 - 64: max of numbers
- U_{out} : 3200 bits \rightarrow 200 (16 bits) Registers
 - $U = 50 \cdot 64 = 3200$ bits
 - 50: max of M
 - 64: max of numbers

Sub-modules

- Error Unit: responsible for detecting when an arithmetic error may occur, like dividing by zero, **h** is getting bigger every time, different sizes....etc.
- FPU: see Section 0.14.4.

0.14.4 Fixed/Floating Point Unit (FPU)

Role

- Instantiated, in the rest modules, multiple times and for different purposes.
- perform the following operations given in port *OP*.
- perform them in different modes given in *MODE*.

Ports

- MODE: IN
 - FIXED(0): operate on LOWEST 16bits of both A and B following Fixed point specifications.

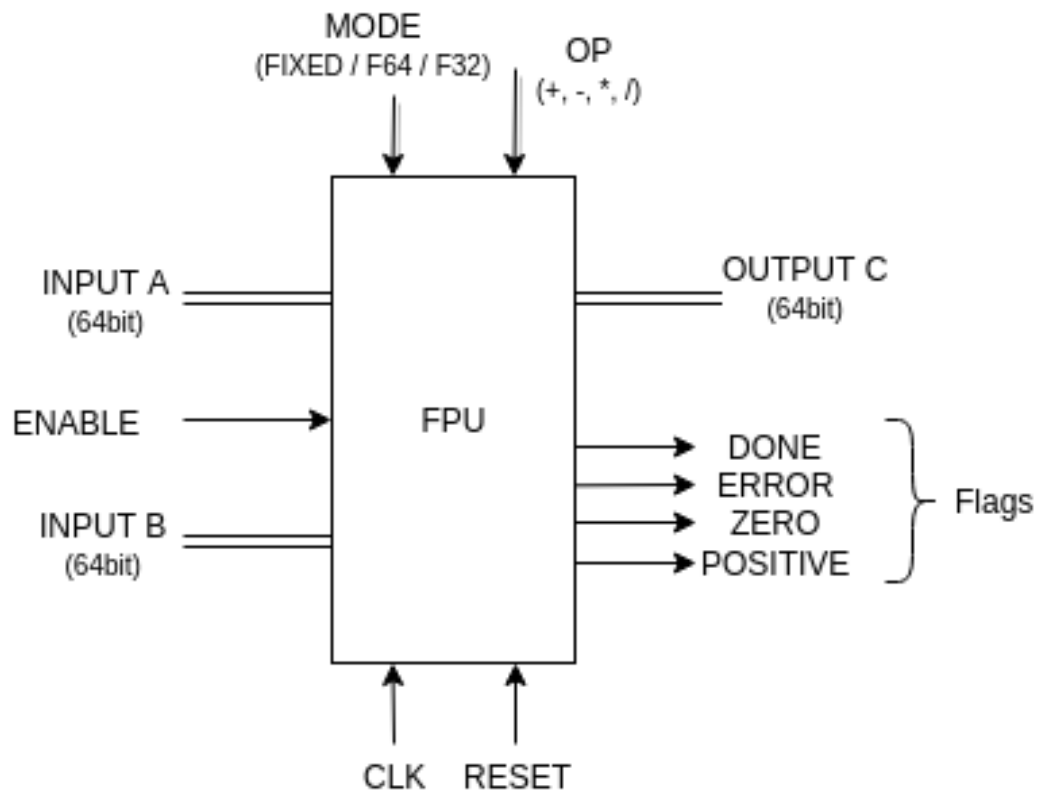


Figure 5: FP Unit

- F64(1): operate on ALL 64bits of both A and B following Floating point specifications.
 - F32(2): operate on LOWEST 32bits of both A and B following Floating point specifications.
- OP: IN
 - 0: $C = A + B$
 - 1: $C = A - B$
 - 2: $C = A * B$
 - 3: $C = \frac{A}{B}$
- CLK: IN

FPU operates on POSITIVE edge.
- DONE: OUT

one operation needs multiple clocks, so FPU must RAISE DONE to 1 when finished for EXACTLY ONE clock cycle.
- ENABLE: IN

FPU must only operate when ENABLE is set to 1 AND there is NO ERROR.
- A, B: IN

Input busses, all are 64bit wide.
- C: OUT

Output bus, 64bit wide.
- ERROR: OUT
 - FPU sets ERROR to 1 when an exception takes place.
 - Each mode has its possible set of exceptions, see the corresponding specifications.
 - MUST stay at 1 after setting it, until RESET is set to 1.
- RESET: ASYNC IN

clears ERROR state REGARDLESS of ENABLE input.

- ZERO: OUT
Set when output C is zero.
- POSITIVE: OUT
Set when output C is positive, otherwise cleared.

Fixed point specifications

- 16bit input
- 16bit output
- CONSTANT scale factor for both input and output = 7. see Subsubsection 0.14.4.
- in case of overflow or division by zero, ERROR MUST be set to 1 and MUST stay at 1 until RESET is set to 1.

Floating point specifications

- MUST adhere to IEEE-754 2019-revision [2] for both fp32 and fp64 modes.
- ERROR is set to 1 when ANY of the exceptions stated in the IEEE-754 takes place, and stay at 1 until RESET is set to 1.

Scale Factor Definition

Scale Factor is an integer used to obtain the real number from the fixed point number and vice versa using the following formula: $y = \frac{x}{2^s}$; where:

- $x \in \mathbb{N}$ is the fixed point number stored as an integer.
- $y \in \mathbb{R}$ is the real number that represents x with some error.
- $s \in \mathbb{N}$ is the scale factor.

0.15 Compression

Follow bit-level Run-length encoding to compress ram content before sending them, by taking each (one to eight) [1:8] repeating bits and compressing them into four bits, using RLE (Run length encoding) algorithm.

0.15.1 Input

Bit stream input data to IO.

0.15.2 Output

Y compressed 4bit packets, where $X \geq Y \geq \lceil \frac{X}{8} \rceil$ Each packet must follow the format in Table 5.

Table 5: Compression Output

Bit Index	Description	Size
3:1	Number of bits to generate - 1	3 bits
0	Bit to generate	1 bit

Note: the 3 bits must be incremented before decode.

0.15.3 Example

Table 6: Compression Example

Original	Compressed
11111111	1111
0000	0110

0.15.4 Algorithm

```

c = first bit in bit_stream
count = 0
for b in bit_stream:
    if c == b and count < 7:
        count++
    else:
        emit_packet(count-1, c)
        count = 1
        c = b

```

Compression Design Justification

Because the occurrence of more than 7 ones or zeros simultaneously is very rare.

Problems

This compression algorithm may not compress the data, rather than that it may increase the number of bits.

0.16 Decompression

Follow this simple algorithm to fill the buffer, once its full (has 32bits) flush it to data bus and update the address.

```
for packet in packets:
    # bit to repeat
    b = packet[0]

    # extract 3 bits
    n = packet[3:1]

    # increment by 1 to get number of repetitions
    n++

    # generate n bit of b
    for _ in range(n):
        buffer.push(b)

    if buffer.full():
        buffer.flush()
```

Part IV

Bibliography

Bibliography

- [1] *ODE Solver – System Description*. CMP Engineering Department, Cairo University. Spring 2020.
https://docs.google.com/document/d/1eHHQ9NC2HlrqLECSz6iQl5c6tR_y5AxtJBdxxkLAY5g4
- [2] Wikipedia contributors. (2020, March 15). IEEE 754. In Wikipedia, The Free Encyclopedia. Retrieved 19:00, March 21, 2020, from
https://en.wikipedia.org/w/index.php?title=IEEE_754&oldid=945680688