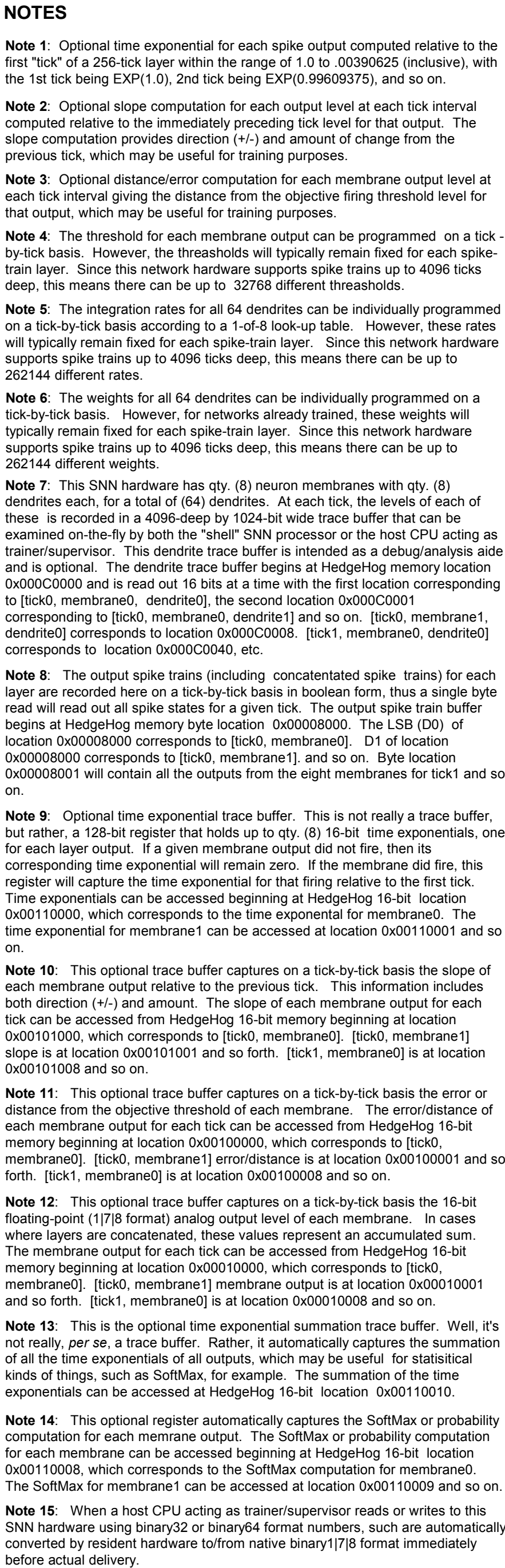


Dubbed the "HedgeHog" because it does this one thing so very well, it implements in hardware qty. (8) neuron membranes with qty. (8) programmable 128-tick input spike integrators each, 4096-tick real-time trace buffers for all dendrite output levels, all membrane output levels, slope computations, distance/error computations, time exponential computations, etc., such computations being done automatically in hardware using its qty. (72) multiply, (72) division, (88) addition, (72) exponential and (8) compare hardware (or equivalent) operators. Thus, clocking at 100MHz in a Xilinx Kintex Ultra Plus, it can perform roughly qty. (312) floating-point operations per clock cycle or roughly 31 billion floating-point operations per second. The Accumulate and Activate bits in the mover shell instruction set enable concatenation of these (8) neuron membranes to build layers with up to qty. (512) inputs and outputs each. The real-time monitor and data exchange feature of the mover shell enables a RISC-V or other CPU to read and write from/to the HedgeHog resources as if ordinary memory.

The Verilog test bench available at the SYMPL HedgeHog repository at GitHub automatically converts human-readable decimal character sequences for weights and thresholds, enabling you to create your test data using Google Sheets online spreadsheet. When your test run is complete, the test bench then automatically converts the binary formatted floating-point trace buffer results from each dendrite and membrane level to decimal character sequences before writing them to their respective output files, enabling you to view your results in numeric or graphic form using Google Sheets online spreadsheet. For examples of resulting graphs on data produced by the Verilog test bench simulations, refer to the reverse side of this document.

If you are exploring or experimenting with spiking neural networks for FPGA embeddable AI applications using RISC-V as host CPU, then the SYMPL HedgeHog is for you.

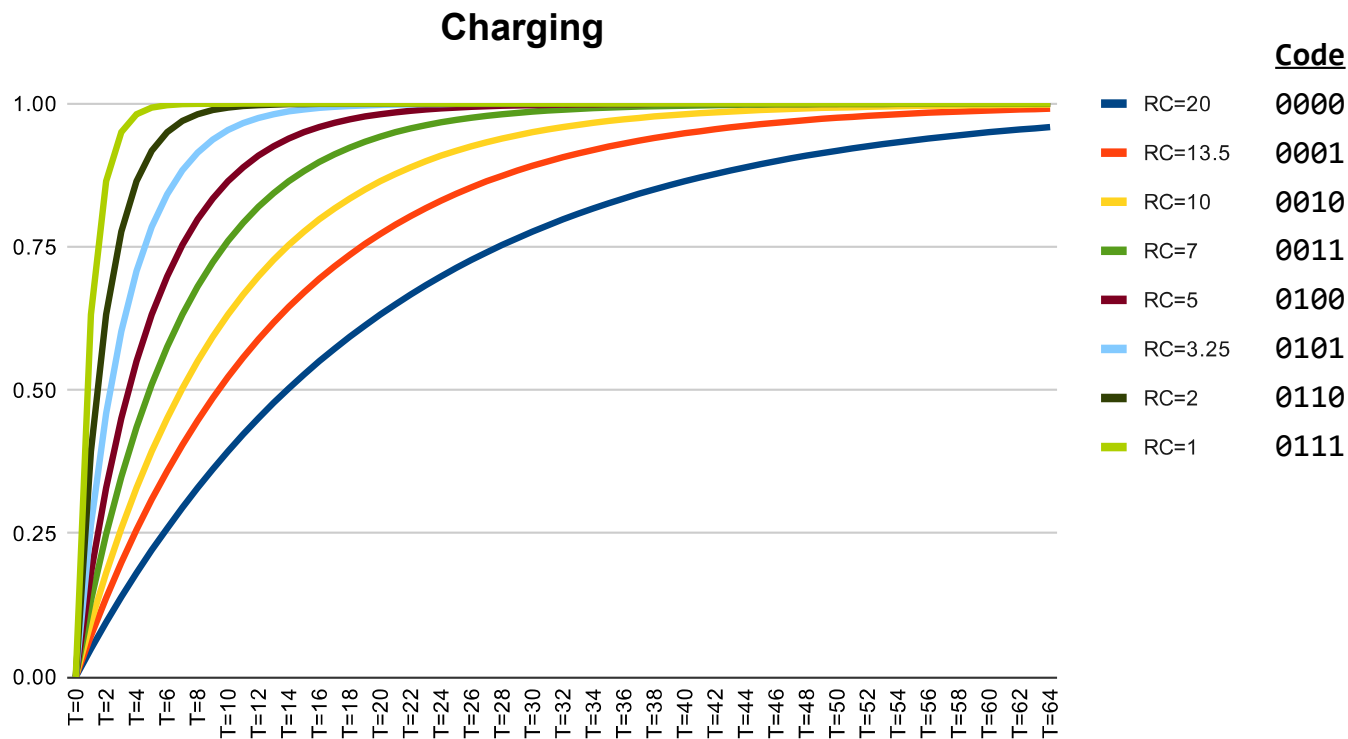


Dynamically Selectable Input Spike Integration Charge Rates

Each of the 64 dendrites comprises an input spike integrator that can be dynamically configured for one of eight charge rates. Once an input spike triggers charging, the dendrite will begin charging according to the selected rate specified for that dendrite and subsequent incoming spikes on that dendrite will be ignored during this charge period until the membrane the dendrite is connected to fires or the dendrite completely decays due to the membrane not firing because its threshold was never reached. There are 8 possible charge rates that can be assigned to a given dendrite.

For charging, the dendrite output level (V_{tick}) at any tick relative to an input spike is automatically computed according to the following equation and graph (note that the legend at the bottom of the graph is the "tick" number):

$$V_{tick}=V_{weight}(1-\exp(-tick/RC))$$

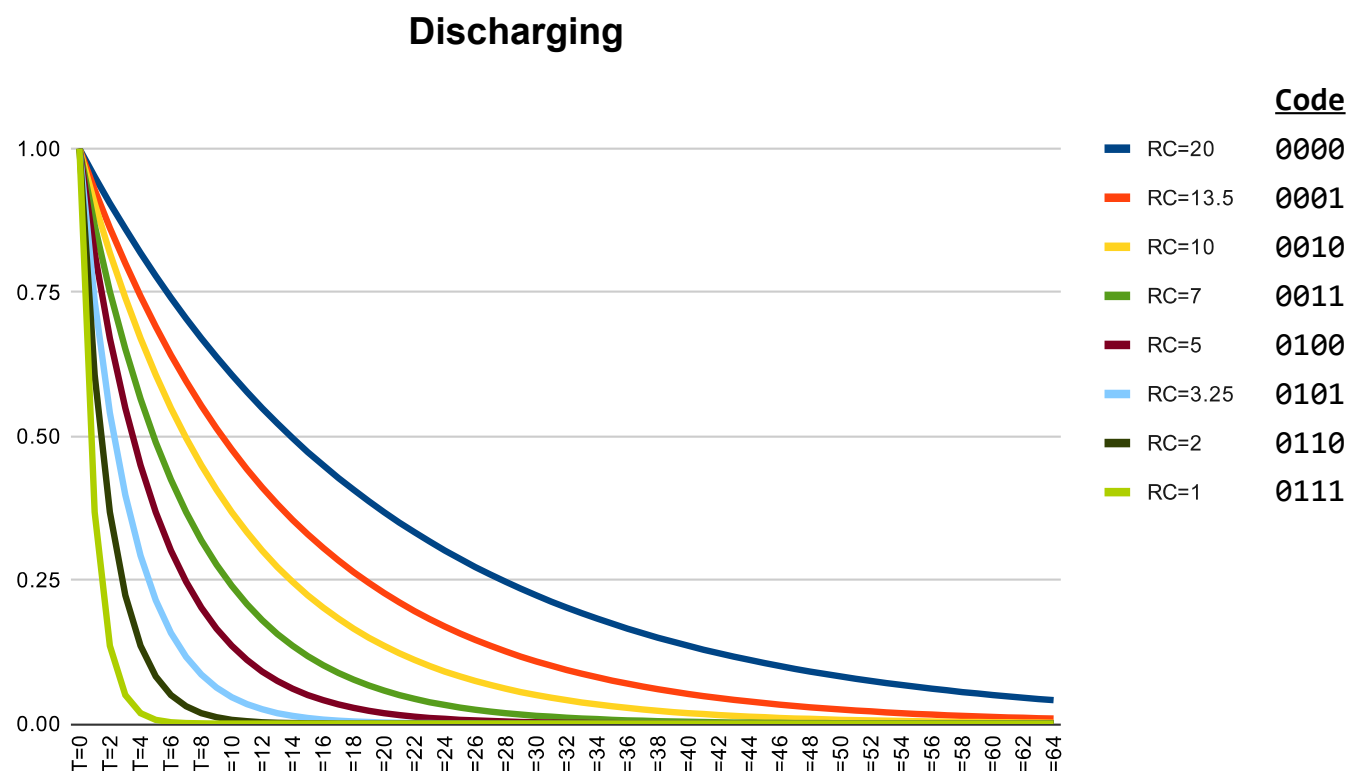


Once a dendrite charge level reaches 98% of its maximum, it is considered charged and on the next tick will automatically begin to discharge/decay according to one of eight possible rates assigned to that dendrite if the threshold of the membrane that dendrite is attached is not exceeded by the sum of all the membrane's dendrite output levels. Each rate has a 4-bit code, with code "0000" being the slowest charge rate and code "0111" being the fastest charge rate. With a code of "0111", a dendrite can be fully charged in just 4 ticks from spike input detection. With a code of "0000", the dendrite will not be fully charged until 64 ticks from spike input detection. Note that bit 3 of the code is always "0", since only 3 bits are ever used.

If the threshold of a membrane is exceeded by the sum of all of its eight dendrites, then the membrane fires, creating a spike on its output, which resets its eight dendrite levels to zero.

For discharging, the dendrite output level (V_{tick}) at any tick relative to the input spike is automatically computed according to the following equation and graph (note that the legend at the bottom of the graph is the "tick" number):

$$V_{tick}=V_{weight}(\exp(-tick/RC))$$



The charge rate codes form a 64-entry vector (one code for each of 64 dendrites), with each code residing in HedgeHog memory space on byte address boundaries starting at location 0x00080000. Thus the rate code for dendrite 00 (membrane0/dendrite0) is at location 0x00080000 and the rate code for dendrite 01 (membrane0/dendrite1) is at location 0x00080001 and so on. These locations can be read from or written to at any time by either the HedgeHog mover shell or host CPU. There are 4096 of these rate code vectors available. This means that it is possible to change the charge rate every tick. However, the usual practice is to use the same vector (and same dendrite charge rates) for all the ticks in a given layer when pushing a given spike train into it.

Dynamically Programmable Weight and Threshold Vectors

Associated with each charge rate vector are the weight and threshold vectors. Each of the 4096 weight vectors comprises qty. (64) 16-bit binary floating-point format 1/7/8 weights, one for each dendrite. Like the charge rate vectors, the weight vectors reside in asymmetrically accessible dual-port SRAM. On the "A" side of the SRAM, the weights can be individually accessed on 16-bit boundaries by the HedgeHog shell or host CPU for initialization and training purposes, while on the "B" side of the SRAM, all 64 weights are read out simultaneously by the HedgeHog shell as a 1024-bit vector during the spike-train push operation.

Weight vector memory resides in HedgeHog indirectly-addressable memory space starting at location 0x00040000. Thus the weight for dendrite 00 (membrane0/dendrite0) is at location 0x00040000 and the weight for dendrite 01 (membrane0/dendrite1) is at location 0x00040001 and so on. These locations can be read from or written to at any time by either the HedgeHog mover shell or host CPU. There are 4096 of these weight vectors available. This means that it is possible to change a given dendrite's weight every tick.

Similarly, each of the 4096 threshold vectors comprises qty. (8) 16-bit binary floating-point format 1/7/8 membrane thresholds, one for each membrane. These threshold values also reside in asymmetrically accessible dual-port SRAM. On the "A" side of the SRAM, the thresholds can be individually accessed on 16-bit boundaries by the HedgeHog shell or host CPU for initialization and training purposes, while on the "B" side of the SRAM all 8 of the thresholds are read out simultaneously by the HedgeHog shell as a 128-bit vector during the input spike-train push operation.

Threshold vector memory resides in HedgeHog indirectly-addressable memory space starting at location 0x00018000. Thus the threshold for membrane0 is at location 0x00018000 and the threshold for membrane1 is at location 0x00018001 and so on. These locations can be read from or written to at any time by either the HedgeHog mover shell or host CPU. There are 4096 of these threshold vectors available. This means that it is possible to change a given membrane's threshold every tick. But this is not the usual practice. The usual practice is to keep the thresholds constant during the entire spike train push into a given layer.

Theory of Operation

The theory of operation is really quite simple. The HedgeHog takes advantage of its instruction REPEAT capability along with its auto-post-modify indirect addressing mode. The basic format of its dual-operand instruction is as follows:

<activate/accumulate> <signal/size>:<destination> = (<signal/size>:<sourceA>, <signal/size>:<sourceB>)

Example 8-input, 8-membrane single-layer:

```
pushSpikes :
-      _4:AR2 = _4:#binWeightsMem      ;point to weights vector
-      _4:AR3 = _4:#spikeSrcMem          ;point to input spike train source buffer
-      _4:AR4 = _4:#SNNinput            ;point to SNN spike train target input
-      _2:REPEAT = _2:#31               ;input spike train is 32 ticks deep in
                                         ;this example

ACT     _1:*AR4++[1] = (_1:*AR3++[1], _128:*AR3++[0]) ;push the spike train and
                                         ;specified weights,rates, thresholds
                                         ;into SNN input
```

In the above example, Auxiliary Registers AR2, AR3 and AR4 are initialized to point to weights/rates/thresholds vector, spike train source buffer, and fused spiking neural network input (respectively). After the REPEAT counter is loaded with the number of ticks (minus 1), the immediately following instruction is executed 32 times. "ACT" means "activate" in the mode column. ACT is set here because this is a simple single-layer push with no concatenation. During execution of this last instruction, the 32-tick spike train, along with the qty. (64) weights, qty. (64) rates, and qty. (8) thresholds, are pushed into the FSSN, all requiring only qty. (32) system clocks.

Upon completion of this last instruction, the HedgeHog will have automatically performed qty. (2304) multiply, (2304) division, (2304) exponential, (2816) addition, and (256) compare floating-point computations and will have captured in their respective trace buffers qty. (32) ticks worth of all dendrite output levels, all membrane output levels, all membrane spike train outputs, time exponentials on all membrane firings, slopes on all membrane output levels, distancece/error on all membrane output levels, a summation of all time exponentials, and a SoftMax (probability) computation on all membrane firings.

Stated another way, with just a single instruction, the HedgeHog performed qty. (9984) floating-point operations in just qty. (32) clocks. The above example presumes of course that the weights, rates and thresholds vectors have been initialized beforehand. For an example of how to do the initialization of the vectors as well as an example of how to perform concatenation of inputs to create 16, 32, 64, 256 or 512 -input/output layers, refer to the demo assembly language listing file at the SYMPL HedgeHog source code repository at GitHub: www.github.com/jerry-D

Below is an example 32-tick, 8-input spike train pushed into the HedgeHog FSSN. Note that the last half of the train is all zeros. With a charge rate of "0110" or 0x6, this is the second fastest charge rate. Thus, these ticks with all zeros in the last half of the spike train provide time for the actual input spikes that precede them to charge up and maybe cause the membrane to fire as a result of the threshold being exceeded.

Input	Spike Train
0	- 0 0 0 0 0 0 0 0 1 0
1	- 0 0 1 0
2	- 0 0 0 1 0
3	- 0 0 0 0 0 0 0 0 0 1 0
4	- 0 0 0 0 0 0 1 0
5	- 0 0 0 1 0
6	- 0 0 0 0 1 0
7	- 0 1 0
	0 1 2 3 4 5 6 7 8 9 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 2 3 3
	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
	Tick Number

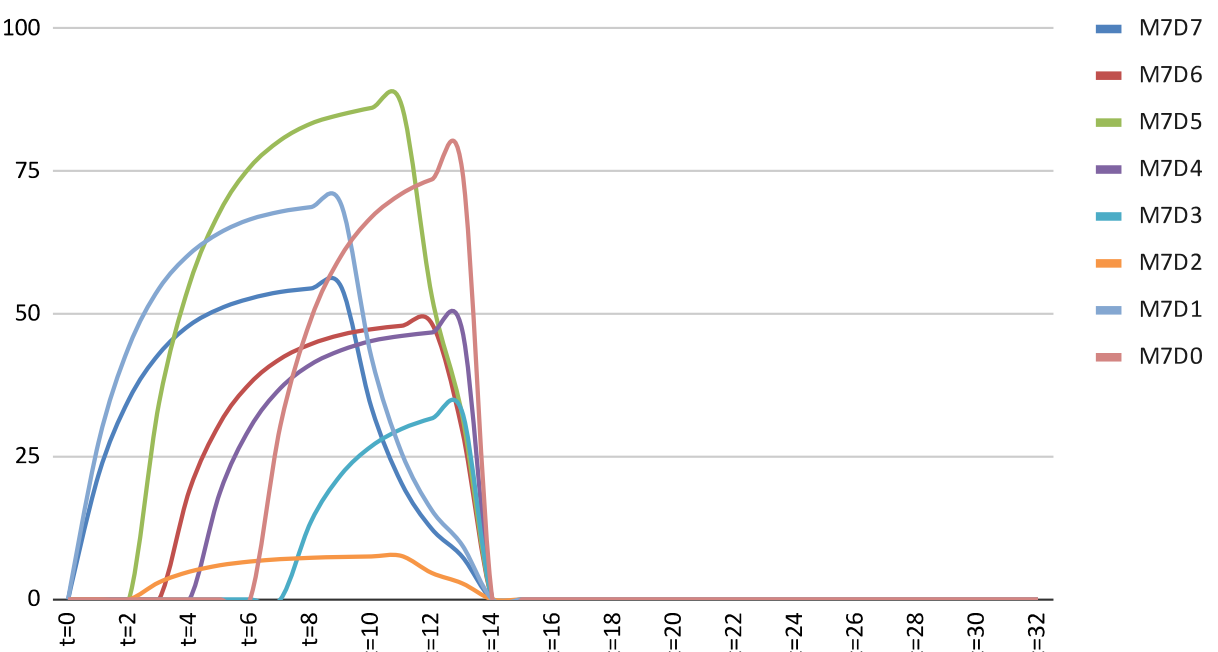
Dendrite Trace Buffer

The HedgeHog FSSN Emulator/Compute Engine includes an optional 1024-bit wide by 4096-tick deep real-time trace buffer that captures all the output levels of all qty. (64) dendrites, which may be useful for debug and analysis purposes. The levels of each dendrite are stored in binary floating-point 1/7/8 format, but the Verilog test bench that is included in the HedgeHog repository at GitHub, converts these to human-readable decimal character sequences before writing it to the "dendrites.txt" file. The textual contents of this file can then be copied and pasted directly into Google Sheets online spreadsheet and plotted. Below are plots generated by Google Sheets from the actual simulation run using the 32-tick spike train shown previously.

The dendrite trace buffer begins at HedgeHog memory location 0x000C0000 and is read out 16 bits at a time with the first location corresponding to [tick0, membrane0, dendrite0], the second location 0x000C0001 corresponding to [tick0, membrane0, dendrite1] and so on. [tick0, membrane1, dendrite0] corresponds to location 0x000C0008. [tick1, membrane0, dendrite0] corresponds to location 0x000C0040, etc.

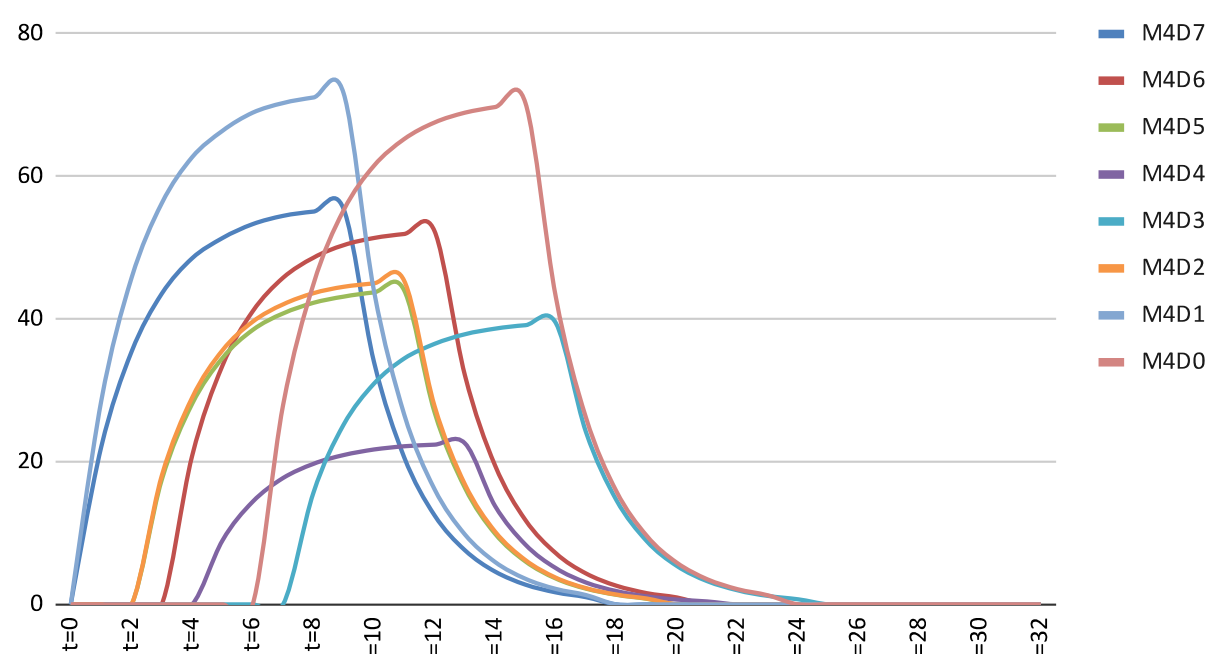
The plot below shows all the output levels of all 8 dendrites connected to Membrane 7, which fires at around tick# 14. Note how the dendrites simultaneously discharge to zero upon firing of the membrane they are connected to.

Dendrites Connected to Spiking Membrane 7



The plot below shows all the output levels of all 8 dendrites connected to Membrane 4, which never fires. Note how the dendrites gradually decay to zero over time according to the specified charge/decay rate.

Dendrites Connected to Non-Spiking Membrane 4



Membrane Output Level Trace Buffer

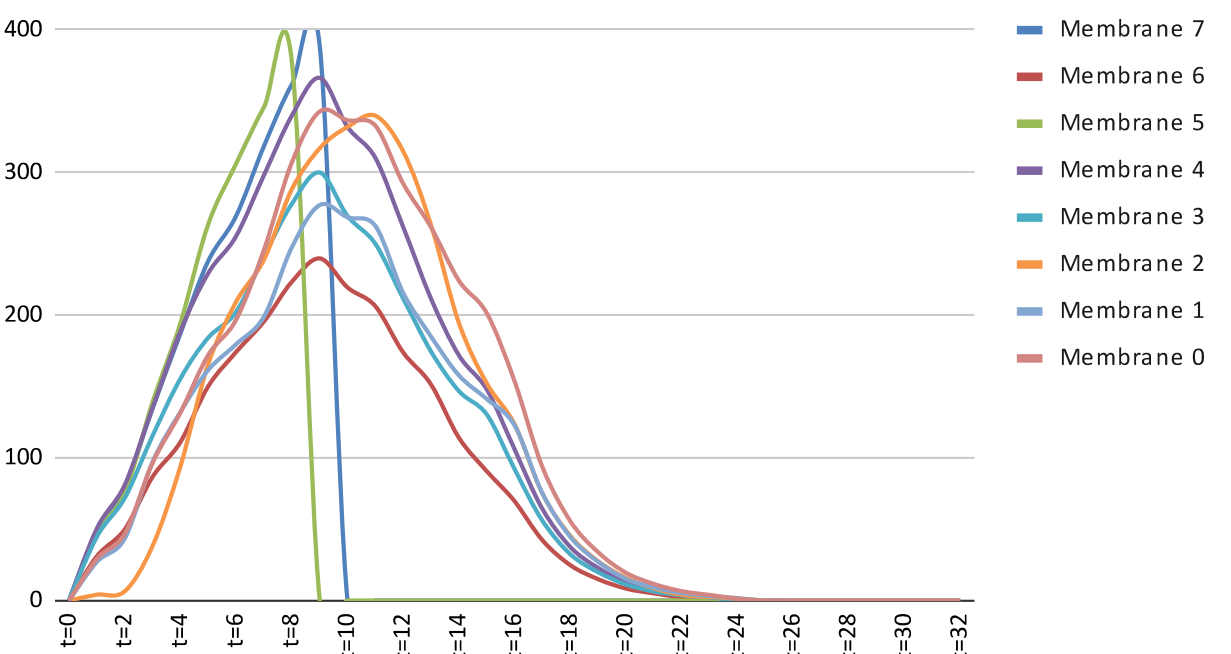
The HedgeHog FSSN Emulator/Compute Engine includes an optional 128-bit wide by 4096-tick deep real-time trace buffer that captures all the output levels of all qty. (8) membranes, which may be useful for debug, analysis and SNN training purposes. In cases where layers are concatenated, these values represent an accumulated sum.

The levels of each membrane are stored in binary floating-point 1/7/8 format, but the Verilog test bench that is included in the HedgeHog repository at GitHub, converts these to human-readable decimal character sequences before writing it to the "outputLevels.txt" file. The textual contents of this file can then be copied and pasted directly into Google Sheets online spreadsheet and plotted. Below are plots generated by Google Sheets from the actual simulation run using the 32-tick spike train shown previously.

The membrane output for each tick can be accessed from HedgeHog 16-bit memory beginning at location 0x00100000, which corresponds to [tick0, membrane0]. [tick0, membrane1] membrane output is at location 0x00100001 and so forth. [tick1, membrane0] is at location 0x00010008 and so on.

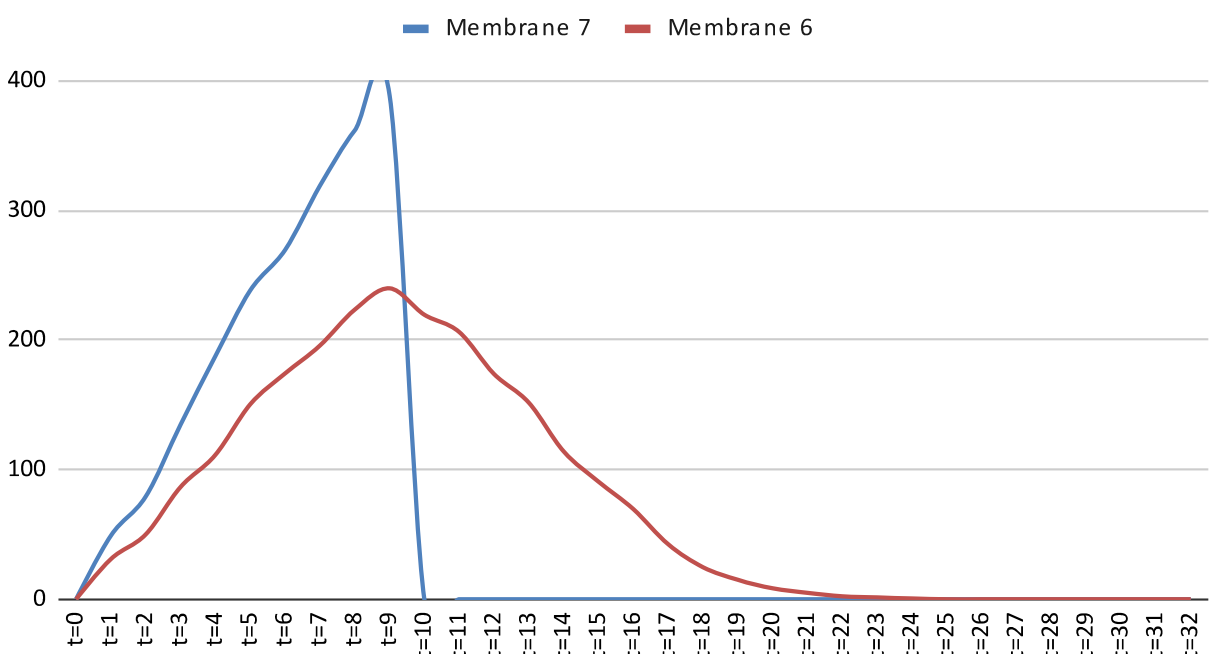
The plot below shows all the output levels of all 8 membranes, with Membrane 5 and Membrane 7 firing at tick# 9 and 10 respectively along with the other membranes gradually decaying over time because their thresholds were never reached.

Composite of all 8 Membrane Output Levels--2 Spiking



The plot below shows the distinction between a membrane spiking because its threshold was exceeded and a non-spiking membrane that gradually decays over time because its threshold was never exceeded.

Spiking Membrane 7 vs Non-Spiking Membrane 6



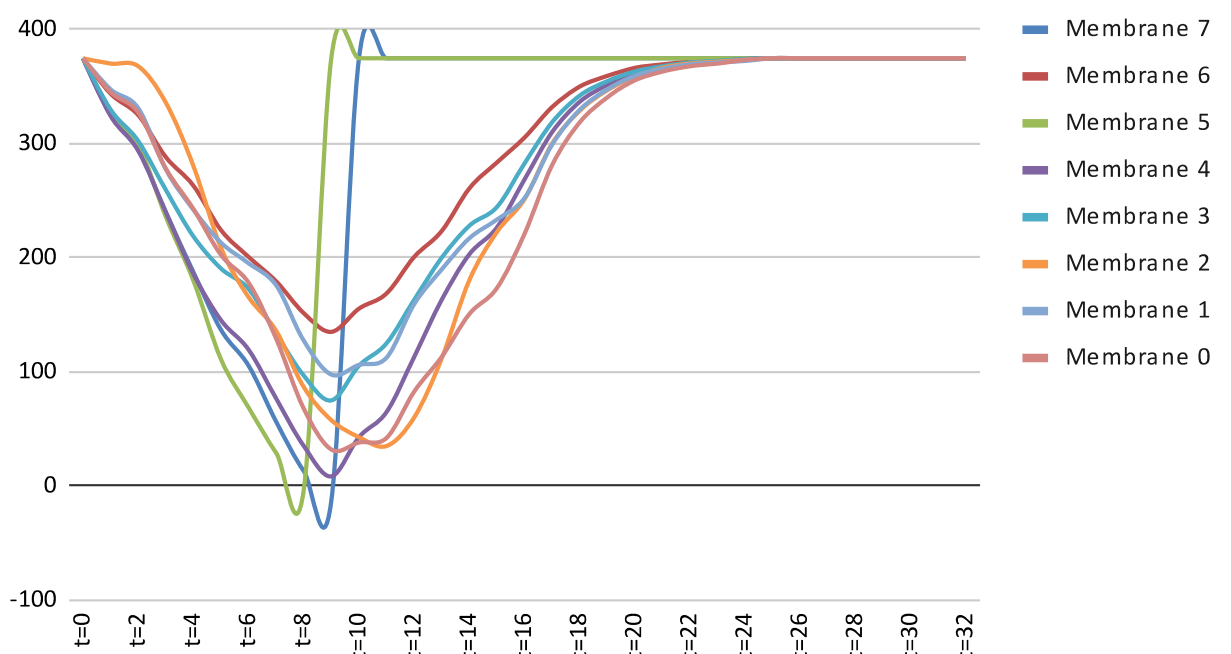
Distance/Error Trace Buffer

The HedgeHog FSSN Emulator/Compute Engine includes an optional 128-bit wide by 4096-tick deep real-time trace buffer that captures all the distance/error between of all 8 membrane output levels and their respective firing thresholds, which may be useful for debug, analysis and SNN training purposes. These distance values are stored in binary floating-point 1/7/8 format, but the Verilog test bench that is included in the HedgeHog repository at GitHub, converts these to human-readable decimal character sequences before writing it to the "errAmounts.txt" file. The textual contents of this file can then be copied and pasted directly into Google Sheets online spreadsheet and plotted. Below are plots generated by Google Sheets from the actual simulation run using the 32-tick spike train shown previously.

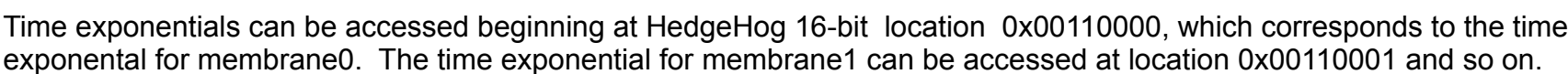
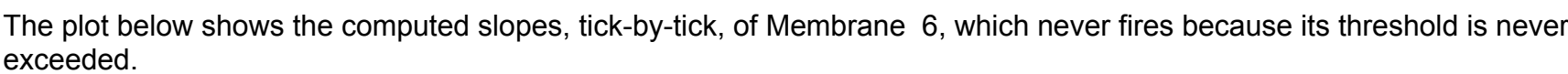
The error/distance of each membrane output for each tick can be accessed from HedgeHog 16-bit memory beginning at location 0x00100000, which corresponds to [tick0, membrane0]. [tick0, membrane1] error/distance is at location 0x00100001 and so forth. [tick1, membrane0] is at location 0x00100008 and so on.

The plot below shows a composite of all 8 membrane distance/error computations, which can be accessed on-the-fly by the HedgeHog mover shell or the host CPU.

Distance/Error Plot of all 8 Membrane Output Levels



The plot below shows a composite of all 8 membrane slope computations, which can be accessed on-the-fly by the HedgeHog mover shell or the host CPU.



The summation of the time exponentials can be accessed at HedgeHog 16-bit location 0x00110010.

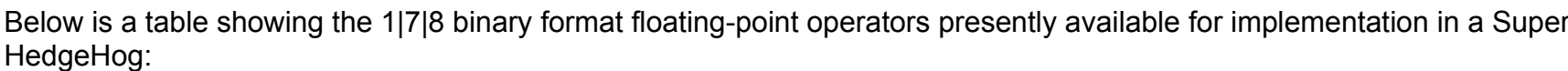
The SoftMax or probability computation for each membrane can be accessed beginning at HedgeHog 16-bit location 0x00110008, which corresponds to the SoftMax computation for membrane0. The SoftMax for membrane1 can be accessed at location 0x00110009 and so on.

Below are the contents of the membrane spike train output trace buffer of the actual simulation run using the 32-tick input spike train shown previously.

Below is the Verilog instantiation used in the test-bench included with the HedgeHog Verilog RTL source code, which emulates a host CPU like the RISC-V for these transactions.

In the table below are descriptions of the signals for the module ports of the above instantiation.

The block diagram below shows how simple a direct connection between the RISC-V CPU and the SYMPL HedgeHog can be. Other approaches such as the use of an AXI-4 or AXI-Lite master/slave arrangement can also be used. In such a configuration, the RISC-V can be employed as a training supervisor or simply be used to initialize the weight, rate and threshold vectors inside the HedgeHog in scenarios where the HedgeHog is being employed as a "PLAYER-only" SNN, wherein the vector data sets are already trained.



Or visit:
www.github.com/jerry-D