

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- config.vhd is a part of Innervator.
4  -- -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9      use ieee.fixed_float_types.all;
10
11 package constants is
12     /* Compile-Time Data File's Location */
13     -- NOTE: You could also use relative paths (../) here, but they
14     -- vary between simulators/synthesizers, defeating the purpose.
15     constant c_DAT_PATH : string :=
16         "C:/Users/Thrae/Desktop/Innervator/data";
17     /* FPGA Constrains & Configurations */
18     -- Clock
19     constant c_CLK_FREQ : positive := 100e6;
20     constant c_CLK_PERD : time := 1 sec / c_CLK_FREQ;
21     -- NOTE: Apparently, some FPGAs (e.g., Xilinx 7 series) work better
22     -- (internally) with active-high resets, because their flip-flops
23     -- were designed to take reset signals as so, and using active-low
24     -- would require an inverter before each flip-flop's SR port; yet,
25     -- board vendors might use negative/active-low signals for their
26     -- reset buttons, because electronics are designed easier that way.
27     -- If we have no control over the (external) polarity of our reset
28     -- signal, a solution is to place a single inverter in the top-
29     -- level hierarchy of the I/O pin logic ("IOB") & use that instead.
30     -- SEE:
31     --     ednasia.com/coding-consideration-for-pipeline-flip-flops
32     --     01signal.com/electronics/iob-registers
33     --
34     --     Also, you might want to "synchronize" (i.e., de-glitch, NOT
35     --     related to sync./async. types of reset) and possibly debounce,
36     --     if it is a button, your external reset signal prior to using it.
37     -- signal.
38     -- Reset (int./ext. = internal/external; neg. = negative)
39     constant c_RST_INVT : boolean := true; -- Invert ext. reset pin
40     constant c_RST_POLE : std_ulogic := '1'; -- '0' = int. neg. reset
41     constant c_RST_SYNC : boolean := true; -- false = async. reset
42     -- Input port synchronization (num. = number)
43     constant c_SYNC_NUM : natural := 3; -- Port sync./deglitch num.
44     -- Button/switch debouncing
45     constant c_DBNC_LIM : time := 30 ms; -- Debounce timeout
46     -- TODO: Have a constant that chooses rising_ or falling_edge
47     --constant c_EDG_RISE
48     /* Internal Fixed-Point Sizing */
49     constant c_WORD_INTG : natural := 4;
50     constant c_WORD_FRAC : natural := 4;
51     constant c_WORD_SIZE : positive := c_WORD_INTG + c_WORD_FRAC;
52     constant c_GUARD_BITS : natural := 0;
53     constant c_FIXED_ROUND : fixed_round_style_type :=
54         fixed_truncate;
55     constant c_FIXED_OFLOW : fixed_overflow_style_type :=
56         fixed_saturate;
57     /* Neuron Settings */
58     -- TODO: Make these arrays to configure each layer/neuron.
59     --
60     -- Number of data to be concurrently processed in a single neuron
61     -- (More = Faster network, at the expense of more FPGA logic usage)

```

## Innervator

```
62  constant c_BATCH_SIZE : positive := 1; -- < or = to data's length.
63  -- Number of pipelining registers in each neuron; this would be
64  -- the amount in clock cycles of latency in input --> output, too.
65  -- (Less = Faster network, at the expense of route timing failure)
66  constant c_PIPE_STAGE : natural := 3; -- 0 = Disable pipelining
67  -- TODO: Add options to select between executing the activation
68  -- function and/or setting the done signal inside or outside
69  -- the neurons' busy state. (Speed/Size trade-off)
70  /* UART Parameters */
71  -- NOTE: Baudrate = Baud, in the digital world
72  constant c_BIT_RATE : positive := 9_600;
73  constant c_BIT_PERD : time      := 1 sec / c_BIT_RATE;
74  end package constants;
75
76
77  -----
78  -- END OF FILE: config.vhd
79  -----
```

```

1  -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- debouncer.vhd is a part of Innervator.
4  -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 library config;
11     use config.constants.all;
12
13 -- Background: When you press a physical button, the metal contacts
14 -- don't make a perfect, clean contact instantly; instead, they might
15 -- "bounce" against each other several times, over a few milliseconds,
16 -- before settling into a closed state. Additionally, Microcontrollers
17 -- and FPGAs are incredibly fast, and they can detect each of those
18 -- tiny bounces as if they were separate button presses; this could
19 -- lead to a single button press being interpreted as multiple presses.
20 --     There are many ways to resolve this matter, and they could be
21 -- done using hardware approaches (e.g., using a resistor-capacitor)
22 -- or software-based ones. In a software approach, we could detect
23 -- a button transition and sample it again at a later point in time,
24 -- which is at least a few milliseconds long (like 10 ms); if the
25 -- button's state had remained the same (i.e., it was "stable"), we
26 -- output that the button was "pressed" once.
27 --     Be aware that other problems arising from external, wired
28 -- interfaces might still apply: we had better accounted for
29 -- metastability and asynchronorized clock domains.
30 entity debouncer is
31     generic (
32         -- Timeout in milliseconds
33         g_TIMEOUT_MS : time := 30 ms
34     );
35     port (
36         i_clk      : in  std_ulogic;
37         i_button    : in  std_ulogic;
38         o_button    : out std_ulogic
39     );
40 end entity debouncer;
41
42 architecture behavioral of debouncer is
43     -- Logically, we are supposed to divide the time by 1000,
44     -- but VHDL simulators just don't like it when you perform
45     -- integer math with physical units like time.
46     constant timeout_ticks : positive :=
47         (g_TIMEOUT_MS / ms) * (c_CLK_FREQ / 1000);
48     signal timeout_count   : natural range 0 to timeout_ticks-1 := 0;
49
50     signal previous_state : std_ulogic := '0';
51 begin
52
53     process (i_clk) is
54     begin
55         if rising_edge(i_clk) then
56             if (i_button /= previous_state and
57                 timeout_count < timeout_ticks-1)
58             then
59                 -- If there's been a change in the button's state, we
60                 -- begin to track it as long as it hasn't stayed
61                 -- "stable" (i.e., unchanged) over the given timeout.

```

## Innervator

```
62         timeout_count <= timeout_count + 1;
63     elsif timeout_count = timeout_ticks-1 then
64         -- Otherwise, for the duration of the timeout, the
65         -- button did not change and can be registered.
66         previous_state <= i_button;
67     else
68         -- No change in button state; keep waiting.
69         timeout_count <= 0;
70     end if;
71 end if;
72 end process;
73
74 o_button <= previous_state;
75
76 end architecture behavioral;
77
78 -- -----
79 -- END OF FILE: debouncer.vhd
80 -- -----
```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- neural_typedefs.vhd is a part of Innervator.
4  -- -----
5
6
7  -- NOTE: You NEED to re-declare used libraries AFTER you instantiate
8  -- the packages here; otherwise, basic types like "std_logic" will not
9  -- get recognized! This is because VHDL's and Ada's ancient compilers
10 -- were "one-pass" and did not keep track of contexes. See:
11 -- https://insights.sigasi.com/tech/use-and-library-vhdl/
12 --library work, std, ieee;
13
14 library ieee;
15
16 library config;
17     use config.constants.all;
18
19 -- TODO: Provide an option to use unsigned types, too.
20 package fixed_neural_pkg is
21     generic (
22         -- NOTE: these two generic numbers should NOT be taken as
23         -- generics, because Vivado 2023.2 has a bug where it throws
24         -- random error messages if types, which are defined in a
25         -- generic package based on generic parameters, are used
26         -- outside; this is not a problem with VHDL or the code but
27         -- with Vivado itself, as always. The solution is to
28         -- "hard-code" them as constants from a global config file.
29         --g_INTEGRAL_BITS, g_FRACTIONAL_BITS : natural;
30         package g_FIXED_PKG_INSTANCE is new ieee.fixed_generic_pkg
31             generic map ( <> ) -- VHDL-2008 Formal Generic Package
32     );
33     use g_FIXED_PKG_INSTANCE.all; -- Import our custom fixed_pkg
34
35     -- Vivado bug workaround (see the note above)
36     constant g_INTEGRAL_BITS    : natural := c_WORD_INTG;
37     constant g_FRACTIONAL_BITS : natural := c_WORD_FRAC;
38
39     -- NOTE: A word is the "primary" size of values that we are
40     -- going to use. Also, historically, a byte was the number
41     -- of bits used to encode a character of text in a computer.
42     -- SEE: https://en.wikipedia.org/wiki/Units_of_information
43     --
44     -- NOTE: In fixed-point arithmetic, the "Q" notation is used to
45     -- indicate the minimum number of bits required to represent a
46     -- range of values. For example, signed Q0.7 uses 1 bit for the
47     -- signedness, 0 bit for the integer part, and 7 bits for the
48     -- fractional part. Similarly, unsigned Q2.5 represents 2
49     -- integer and 5 fractional bits.
50     -- See: inst.eecs.berkeley.edu/~cs61c/sp06/handout/fixedpt.html
51     --
52     -- NOTE: Bit Width = |INTEGRAL_BITS| + |FRACTIONAL_BITS|; the
53     -- integral (i.e., whole) part comprises of [INTEGRAL_BITS-1,
54     -- 0] and the fractional (i.e., after the .decimal point)
55     -- comprises of [-1, -FRACTIONAL_BITS]
56
57
58     /*
59     Types derived from the generically given specifications
60     */
61

```

```

62  /* Scalars */
63  subtype neural_bit    is -- An "analogue" bit with the range [0, 1)
64      u_ufixed (-1 downto -(g_FRACTIONAL_BITS*2)); -- UNSIGNED and <1
65  subtype neural_nibble is -- Half-word
66      u_sfixed ((g_INTEGRAL_BITS/2)-1 downto -(g_FRACTIONAL_BITS/2));
67  subtype neural_word   is -- Full-word
68      u_sfixed (g_INTEGRAL_BITS-1 downto -g_FRACTIONAL_BITS);
69  subtype neural_dword  is -- Double-word
70      u_sfixed ((g_INTEGRAL_BITS*2)-1 downto -(g_FRACTIONAL_BITS*2));
71  subtype neural_qword  is -- Quadruple-word
72      u_sfixed ((g_INTEGRAL_BITS*4)-1 downto -(g_FRACTIONAL_BITS*4));
73  subtype neural_oword  is -- Octuple-word
74      u_sfixed ((g_INTEGRAL_BITS*8)-1 downto -(g_FRACTIONAL_BITS*8));
75  /* Vectors */
76  -- A single-row (1-D) array of neural_*word values.
77  type neural_vector    is array (natural range <>) of neural_word;
78  type neural_bvector   is array (natural range <>) of neural_bit;
79  type neural_nvector   is array (natural range <>) of neural_nibble;
80  alias neural_wvector  is neural_vector;
81  type neural_dvector   is array (natural range <>) of neural_dword;
82  type neural_qvector   is array (natural range <>) of neural_qword;
83  type neural_ovector   is array (natural range <>) of neural_oword;
84  /* Matrices */
85  -- A multi-row/nested (NOT 2-D) array of neural_*vector arrays.
86  type neural_matrix    is array (natural range <>) of neural_vector;
87  type neural_bmatrix   is array (natural range <>) of neural_bvector;
88  type neural_nmatrix   is array (natural range <>) of neural_nvector;
89  alias neural_wmatrix  is neural_matrix;
90  type neural_dmatrix   is array (natural range <>) of neural_dvector;
91  type neural_qmatrix   is array (natural range <>) of neural_qvector;
92  type neural_omatrix   is array (natural range <>) of neural_ovector;
93
94  /*
95      Shorthand notations of the above-mentioned types
96  */
97
98  /* Scalars */
99  alias nrl_bit    is neural_bit;
100  alias nrl_nib    is neural_nibble;
101  alias nrl_wrd    is neural_word;
102  alias nrl_dwd    is neural_dword;
103  alias nrl_qwd    is neural_qword;
104  alias nrl_owd    is neural_oword;
105  /* Vectors */
106  alias nrl_vec    is neural_vector;
107  alias nrl_bvec   is neural_bvector;
108  alias nrl_nvec   is neural_nvector;
109  alias nrl_wvec   is neural_wvector;
110  alias nrl_dvec   is neural_dvector;
111  alias nrl_qvec   is neural_qvector;
112  alias nrl_ovec   is neural_ovector;
113  /* Matrices */
114  alias nrl_mat    is neural_matrix;
115  alias nrl_bmat   is neural_bmatrix;
116  alias nrl_nmat   is neural_nmatrix;
117  alias nrl_wmat   is neural_wmatrix;
118  alias nrl_dmat   is neural_dmatrix;
119  alias nrl_qmat   is neural_qmatrix;
120  alias nrl_omat   is neural_omatrix;
121
122  -- Custom helper types for use in finding files' dimensions
123  -- (i.e., number of neurons and weights within each neuron)

```

```

124     type dimension is record
125         rows : natural;
126         cols : natural;
127     end record;
128     type dimensions_array is
129         array (natural range <>) of dimension;
130
131     -- An enumeration/structure of a layer's neuron(s)' weights (array
132     -- of array) AND each neuron's associated bias (array).
133     type layer_parameters is record
134         dims      : dimension; -- Array limitations' workaround
135         weights   : neural_matrix;
136         biases    : neural_vector;
137     end record;
138
139     -- NOTE: This is more limited than it seems, as each record would
140     -- have to be of exactly the same length in a VHDL array of records
141     -- On the other hand, dynamically allocated arrays, which are
142     -- able to hold variable-length elements, are NOT allowed in
143     -- synthesizable VHDL at all.
144     -- See: https://stackoverflow.com/a/61031840
145     --
146     -- Array of records to hold the parameters (i.e., weights and
147     -- biases) of each and every layer, amounting to the entire network
148     type network_layers is
149         array (natural range <>) of layer_parameters;
150
151
152     -- -----
153     -- rtl_synthesis off
154     -- pragma translate_off
155     -- -----
156     -- Simulation-Time Pointer types: You may access these types'
157     -- dereferenced variables' values using ".all" suffixes after them.
158     -- Also, you may have "dynamically sized" arrays by concatenating
159     -- new "new" definitions with previous ones. Lastly, you may re-
160     -- use declared pointer types by using "deallocate()" on them.
161     type neural_word_ptr    is access neural_word;
162     type neural_vector_ptr  is access neural_vector;
163     type neural_matrix_ptr  is access neural_matrix;
164     type network_layers_ptr is access network_layers;
165
166     -- NOTE: VHDL doesn't permit files of multidimensionals (e.g.,
167     -- matrices) directly, so some improvisation with pointer types
168     -- is required in order to retrieve them as arrays of arrays.
169     --
170     -- This could also be a file of fixed-points as fixed_pkg provides
171     -- textio.read, although ONLY in binary, octal, or hex formats.
172     type neural_file is file of neural_word;
173
174     -- -----
175     -- pragma translate_on
176     -- rtl_synthesis on
177     -- -----
178
179     -- NOTE: These "partially" resize their inputs; they are used as
180     -- tricks to force synthesizers tools to let us use arrays of
181     -- differing element sizes. In practice, these 'resize's only
182     -- "expand" their inputs and don't fully scale the data; they
183     -- _partially_ fill-in the expanded result.
184     function resize(
185         input_arr : in neural_vector;

```

```

186         target_dim : in dimension
187     ) return neural_vector;
188
189     function resize(
190         input_mat : in neural_matrix;
191         target_dim : in dimension
192     ) return neural_matrix;
193
194
195     -- NOTE: These refuse to work with unconstrained generic types;
196     -- so, you will have to define the manually for each constant
197     -- decleration, if you want to use them.
198     /*
199     attribute num_inputs : natural;
200     attribute num_outputs : natural;
201
202     attribute num_inputs of neural_vector : type is
203         neural_vector'length;
204     attribute num_outputs of neural_vector : type is
205         neural_vector'length;
206
207     attribute num_inputs of neural_matrix : type is
208         neural_matrix(0)'length; -- or 'length(dimension) for 2-D array
209     attribute num_outputs of neural_matrix : type is
210         neural_matrix'length;
211
212     -- NOTE: Records can't have methods; there's no protected record.
213     -- So, we can use the rather obscure 'group' keyword to accomplish
214     -- a similar objective.
215     group layer_group_type is ( constant <> );
216     attribute num_neurons : natural;
217
218     group layer_group : layer_group_type (layer_parameters);
219     attribute num_neurons of layer_group : group is
220         layer_parameters.biases'length; -- or weights'length
221     */
222
223     -- Another way of realizing the above idea is by using functions.
224     /*
225     function len_inputs(
226         input_arr : in neural_vector
227     ) return natural;
228     function len_outputs(
229         input_arr : in neural_vector
230     ) return natural;
231
232     function len_inputs(
233         input_mat : in neural_matrix
234     ) return natural;
235     function len_outputs(
236         input_mat : in neural_matrix
237     ) return natural;
238
239     function len_neurons(
240         input_arr : in neural_vector
241     ) return natural;
242     function len_neurons(
243         input_arr : in neural_vector
244     ) return natural; -- array of naturals?
245     */
246 end package fixed_neural_pkg;
247

```



```

248 -- TODO: Whenever Synthesis tools support protected types inside
249 -- synthesis (for constant initializations), as opposed to just
250 -- simulation, convert the methods inside this package body to
251 -- be inside a 'protected' type of 'neural_vector' and '_matrix'.
252 package body fixed_neural_pkg is
253
254     function resize(
255         input_arr  : in neural_vector;
256         target_dim : in dimension
257     ) return neural_vector is
258         variable resized_arr : neural_vector (0 to target_dim.rows-1);
259     begin
260         resized_arr(input_arr'range) := input_arr; -- Partially fill it
261
262         return resized_arr;
263     end function;
264
265     function resize(
266         input_mat  : in neural_matrix;
267         target_dim : in dimension
268     ) return neural_matrix is
269         variable resized_mat : neural_matrix
270             (0 to target_dim.rows-1) (0 to target_dim.cols-1);
271     begin
272         for row in input_mat'range loop
273             --resized_mat(row) := resize(input_mat(row), target_dim);
274             resized_mat(row)(input_mat(row)'range) := input_mat(row);
275         end loop;
276
277         return resized_mat;
278     end function;
279
280 end package body fixed_neural_pkg;
281
282
283
284 -- "'guard_bits' defaults to 'fixed_guard_bits,' which defaults
285 -- to 3. Guard bits are used in the rounding routines. If guard
286 -- is set to 0, the rounding is automatically turned off.
287 -- These extra bits are added to the end of the numbers in the
288 -- division and "to_real" functions to make the numbers more
289 -- accurate." (Fixed point package user's guide by Mr. David Bishop)
290 library ieee;
291 library config;
292 use config.constants.all;
293 package fixed_pkg_for_neural is new ieee.fixed_generic_pkg
294     generic map ( -- NOTE: ieee_proposed pre VHDL-08
295         fixed_round_style      => c_FIXED_ROUND,
296         fixed_overflow_style   => c_FIXED_OFLOW,
297         fixed_guard_bits       => c_GUARD_BITS,
298         no_warning             => false
299     );
300
301 library neural;
302 library config;
303 use config.constants.all;
304 package neural_typedefs is new neural.fixed_neural_pkg
305     generic map (
306         --INTEGRAL_BITS      => c_WORD_INTG, -- NOTE: Signed
307         --FRACTIONAL_BITS   => c_WORD_FRAC,
308         g_FIXED_PKG_INSTANCE => work.fixed_pkg_for_neural
309     );

```

```

310 -- After instantiation, you may use the Packages above as follows:
311 -- use work.fixed_pkg_for_neural.all, work.neural_typedefs.all;
312 --
313 -- Alternatively, you may use their bundle as a context
314 -- (defined near the end of this file).
315
316 -- NOTE: IEEE's official fixed-point package has a bug (in
317 -- IEEE 1076-2008, a.k.a. VHDL-08), where to_ufixed(sfixed) isn't
318 -- defined in the Package's header file, while its parallel
319 -- to_sfixed(ufixed) is. So, the solution is to copy-paste the
320 -- declaration of said function from fixed_generic_pkg-body.vhdl
321 -- into our local library.
322 -- SEE: https://gitlab.com/IEEE-P1076/VHDL-Issues/-/issues/269
323 library ieee;
324     use ieee.std_logic_1164.all; -- Added to fix ModelSim's errors
325
326 library work;
327     use work.fixed_pkg_for_neural.all, work.neural_typedefs.all;
328
329 -- TODO: Include arrays of std_(u)logic_vector types?
330 --type t_slv_arr is array (natural range <>) of std_logic_vector;
331 --type t_suv_arr is array (natural range <>) of std_ulogic_vector;
332 package fixed_generic_pkg_bugfix is
333     function to_ufixed(
334         arg : UNRESOLVED_sfixed
335     ) return UNRESOLVED_ufixed;
336 end package fixed_generic_pkg_bugfix;
337
338 package body fixed_generic_pkg_bugfix is
339     -- null array constants
340     constant NAUF : UNRESOLVED_ufixed (0 downto 1) :=
341         (others => '0');
342
343     -- Special version of "minimum" to do some
344     -- boundary checking with errors
345     function mine(l, r : INTEGER) return INTEGER is
346     begin
347         if (L = INTEGER'low or R = INTEGER'low) then
348             report fixed_generic_pkg_bugfix'instance_name
349                 & " Unbounded number passed, was a literal used?"
350                 severity error;
351             return 0;
352         end if;
353
354         return minimum (L, R);
355     end function mine;
356
357     -- converts an sfixed into a ufixed. The output is the same
358     -- length as the input, because abs("1000") = "1000" = 8.
359     function to_ufixed(
360         arg : UNRESOLVED_sfixed
361     ) return UNRESOLVED_ufixed is
362         constant left_index  : INTEGER := arg'high;
363         constant right_index : INTEGER := mine(arg'low, arg'low);
364         variable xarg        :
365             UNRESOLVED_sfixed (left_index+1 downto right_index);
366         variable result      :
367             UNRESOLVED_ufixed (left_index downto right_index);
368     begin
369         if arg'length < 1 then
370             return NAUF;
371         end if;

```

```

372
373     xarg    := abs(arg);
374     result := UNRESOLVED_ufixed (
375         xarg (left_index downto right_index)
376     );
377
378     return result;
379
380 end function to_ufixed;
381 end package body fixed_generic_pkg_bugfix;
382
383 -- After instantiation, you may use the Packages above as follows:
384 -- use work.fixed_pkg_for_neural.all, work.neural_typedefs.all;
385 --
386 -- Alternatively, you may use their bundle as a context:
387 context neural_context is -- VHDL-2008 feature
388     library neural;
389     use neural.fixed_pkg_for_neural.all,
390         neural.neural_typedefs.all;
391     use neural.fixed_generic_pkg_bugfix.all; -- REQUIRED!
392 end context neural_context;
393
394 -- NOTE: Unfortunately, even if I play by Vivado Simulator's rules by
395 -- placing the VHDL-93 compatibility version of IEEE's fixed_pkg into
396 -- a local directory AND commenting-out all homographes of std_logic_
397 -- vectors AND removing all references to 'line' datatypes, it still
398 -- finds a way to crash abruptly in other areas (e.g., File I/O),
399 -- without any log whatsoever, in simulation; working with Vivado's
400 -- Simulator is pointless, as even a ModelSim version from 8 years ago
401 -- (as of 2024) far outperforms it.
402
403 -- -----
404 -- rtl_synthesis off
405 -- pragma translate_off
406 -- -----
407
408 -- NOTE: Here, in simulation, we OVERWRITE the previous declaration;
409 -- this is done because there is no other way to detect whether the
410 -- code is being simulated and skip synth's code, in the latter case we
411 -- also need to use a custom version of IEEE's fixed_pkg due to Vivado.
412 /*
413 context neural_context is
414     library ieee_proposed;
415     use ieee_proposed.fixed_pkg.all;
416     -- NOTE: Xilinx Vivado does not support fixed- or floating-point
417     -- packages for use within its simulator, even though it can
418     -- synthesize them just fine; as a workaround, use
419     -- a LOCAL 'ieee_proposed' instead of 'ieee'.
420     --
421     -- SEE:
422     --     docs.xilinx.com/r/en-US/ug900-vivado-logic-simulation/
423     --     Fixed-and-Floating-Point-Packages
424     --
425     --     insights.sigasi.com/tech/list-known-vhdl-metacomment-pragmas
426     library config;
427     use config.neural_typedefs.all;
428     use config.fixed_generic_pkg_bugfix.all; -- REQUIRED!
429 end context neural_context;
430 */
431
432 -- -----
433 -- pragma translate_on

```

```
434 -- rtl_synthesis on
435 -- -----
436
437 -- -----
438 -- END OF FILE: neural_typedefs.vhd
439 -- -----
```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- activation.vhd is a part of Innervator.
4  -- -----
5
6  library ieee;
7      use ieee.std_logic_1164.all;
8      --use ieee.math_real.MATH_E; -- NOT synthetizable; use as generics!
9
10 library work;
11     context work.neural_context;
12
13
14 -- TODO: Implement more activation functions (e.g., ReLU, etc.)
15 package activation is
16
17     alias in_type is neural_dword;
18     alias out_type is neural_bit;
19
20     function sigmoid(
21         x : in_type
22     ) return out_type;
23
24 end package activation;
25
26
27 package body activation is
28
29     function sigmoid(
30         x : in_type
31     ) return out_type is
32         -- TODO: Automate the generation of these constant look-up
33         -- parameters to be done at compile-time, based on their sizes.
34
35         -- Boundaries beyond which linear approximation begins erring.
36         constant UPPER_BOUND : in_type :=
37             to_sfixed(+2.0625, in_type'high, in_type'low);
38         constant LOWER_BOUND : in_type :=
39             to_sfixed(-2.0625, in_type'high, in_type'low);
40
41         constant LINEAR_M : in_type := -- M (the coefficient)
42             to_sfixed(0.1875, in_type'high, in_type'low);
43         constant LINEAR_C : in_type := -- C (the displacement)
44             to_sfixed(0.5, in_type'high, in_type'low);
45
46         -- Look-Up constants for when the input value falls beyond
47         -- the Linear Approximation's acceptable error range.
48         --
49         -- NOTE: ZERO and ONE do not correspond to exactly 0.0 and 1.0;
50         -- instead, they are a bit "leaky," similar to Sigmoid's output
51         --
52         -- Actually, these are not really required and only add to the
53         -- problem by introducing additional logic timing delay.
54         /*
55         constant ZERO : out_type :=
56             to_ufixed(0.0625, out_type'high, out_type'low);
57         constant ONE : out_type :=
58             to_ufixed(0.9375, out_type'high, out_type'low);
59         */
60
61         -- A temporary variable is required, because only VHDL-2019

```

```
62      -- allows conditional assignment (i.e., when...else) in
63      -- front of return statements, and we are using VHDL-2008.
64      variable result : out_type;
65
66      --attribute use_dsp : string;
67      --attribute use_dsp of result : variable is "no";
68      begin
69
70      -- It is rather hard to implement exponents and logarithms
71      -- in an FPGA; it would use too many logic blocks. Even
72      -- approximating Sigmoid using the absolute value would
73      -- still require 'x' to be divided, very slowly, by (1 + |x|).
74      -- So, a linear approximation is used instead.
75      result := --ZERO when x < LOWER_BOUND else
76                --ONE  when x > UPPER_BOUND else
77                resize(to_ufixed((LINEAR_M * x) + LINEAR_C), result);
78
79      return result;
80
81      end function sigmoid;
82
83  end package body activation;
84
85  -- -----
86  -- END OF FILE: activation.vhd
87  -- -----
```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- file_parser.vhd is a part of Innervator.
4  -- -----
5
6
7  -- TODO: Because the std.textio is so limited in VHDL, maybe
8  -- try and investigate VHDL's VHPI to call C/C++ functions?
9
10 library std;
11     use std.textio.all;
12
13 library ieee;
14     use ieee.std_logic_1164.all;
15     --use ieee.std_logic_textio.all;
16
17 -- TODO: Make these subtypes given through generics, in VHDL-2019.
18 library neural;
19     context neural.neural_context;
20
21 -- NOTE: The reason this package has been divided into two, including
22 -- an auxiliary one, is because I could not use the same functions
23 -- defined in a single package's body inside constant initializations
24 -- of its header (since the functions were not 'elaborated' on yet);
25 -- while you might think that 'deferred constants' were a solution to
26 -- this, they actually were not.
27 --     Because I also needed to use said deferred constants inside
28 -- the subtypes/types defined in the actual package's header, I
29 -- would trigger simulator errors related to "illegal use of deferred
30 -- constants." So, this was the only solution to access both.
31 --     NOTE: Actually, the above-mentioned note no longer applies,
32 -- because I was able to find a way to constrain array types as
33 -- subtypes which could then be connected to external, unconstrained
34 -- types in constant evaluations, but it is nonetheless a good
35 -- practice to separate auxiliary functionality here.
36 package file_parser_aux is
37
38     -- NOTE: Because the fixed-point package only provides procedures
39     -- to read its types from files in the textual format of binary,
40     -- octal, and hexadecimal (i.e., NOT decimal), I initially chose
41     -- to read the inputs as floating-point 'real' types which would
42     -- then be converted into fixed-point internally. HOWEVER, this
43     -- brought me to the second issue: Xilinx Vivado has a bug in its
44     -- implementation of the 'read()' procedure for 'real' datatypes,
45     -- which prevents you from reading more than 7 elements from a
46     -- single row. THUS, I had to resort to the lesser-known 'sread()'
47     -- procedure, which takes in a (forcibly constrained) string,
48     -- strips the given line of any and all types of whitespace
49     -- characters and, lastly, tokenizes each whitespace-separated
50     -- collection of characters into the input string; this way,
51     -- I could get past Vivado's bug and convert the string
52     -- representation of floating-points into a 'real' and finally
53     -- an actual fixed-point... or, could I?
54     --     Vivado, unsurprisingly, did NOT care to implement 'sread()'
55     -- AT ALL; while it appears to synthesize correctly, it hard-
56     -- crashes the simulator with NO log and also does not even work
57     -- in synthesis itself (tested using 'assert' statements).
58     -- To make matters even worse, Vivado implemented the fixed-point
59     -- package---an official IEEE package---in an incredibly poor and
60     -- "hacky" way; not only does the simulator outright REFUSE to work
61     -- with the ieee.fixed_generic_pkg, but it also suggests, even as of

```

```

62  -- the 2023.2 version, that you use the now-OUTDATED workaround of
63  -- replacing ieee with ieee_proposed as to use the VHDL-93 version
64  -- of the Package. The problem here is that, sometime around 2021,
65  -- they removed ieee_proposed from the package list, making their
66  -- OWN workaround OBSOLETE. Not only that but, because they again
67  -- failed to implement std.textio's 'line' type properly, you also
68  -- cannot place the raw files of fixed-point package into your own
69  -- project's directory and resolve the issue; Vivado does not even
70  -- let you define or use a procedure that takes in a 'line' type
71  -- parameter. This also brings us to the next issue: I could NOT
72  -- use fixed-point package's implementation of sfixed/ufixed parser
73  -- [i.e., read()]. SEE:
74  -- https://support.xilinx.com/s/question/0D52E000061LghFSAS
75  -- ULTIMATELY, I had to resort back to using Vivado's broken
76  -- implementation of read() (which is the ONLY choice) and
77  -- flatten my input files to not contain more than a single
78  -- element per line.
79  -- But then I encountered yet another problem: Vivado's
80  -- 'readline()' behaved very strangely, reporting the 'length
81  -- attribute to be >0 even for otherwise-blank lines. The issue
82  -- as I later found out, was that Vivado, for some reason, treats
83  -- files with a '.txt' extension in a special manner and readline()
84  -- completely breaks on those; so, I had to rename it to something
85  -- else (e.g., '.dat'). This should not happen since 'text' is
86  -- defined as 'file of string' in std.textio, and, if I wanted to
87  -- have my file read as binary or any other type, I would define
88  -- my own 'type binary is file of bit', yet Vivado makes this
89  -- (undocumented) assumption on its own... SEE:
90  -- reddit.com/r/FPGA/comments/16th9ok
91  -- /textio_not_reading_negative_integers_from_file_in
92  -- That seemed to solve the Issue only in simulation, although
93  -- read() would still not read more than 7 elements per line; in
94  -- synthesis, read() was STILL broken beyond repair. Eventually,
95  -- I found that Vivado actually implemented the read() procedure
96  -- correctly ONLY in case of 'bit_vector' and 'std_logic_vector';
97  -- NOTHING else---not even simple positive integers---works.
98  -- Also, I found out that another limitation is that you CANNOT
99  -- detect "blank" lines; for some reason, Vivado will always report
100 -- the current line's 'length attribute to be the same throughout
101 -- an entire file. Fortunately, this could be worked-around by
102 -- merely using rows full of 'X', 'Z', or 'U' as "delimiters."
103 -- inside COMMENTS, and its read() was completely dysfunctional).
104 -- Never have I had to wrestle so much with a toolchain just to
105 -- get something so incredibly trivial to work, but this seems to
106 -- be commonplace within EDA tools: I also looked into Synopsys'
107 -- Synplify to see if they implemented file I/O properly there,
108 -- but Synplify seemed even buggier than Vivado (e.g., it would
109 -- break when you used ASCII grave accents and specific keywords
110 -- such as protected).
111 --
112 -- NOTE: Do NOT use 'ulogic' as Vivado's read() is broken for it.
113 subtype read_t is std_logic_vector(neural_word'length-1 downto 0);
114 --
115 -- NOTE: Also, since you cannot leave out the 'out' parameters of
116 -- procedures as 'open' in VHDL (unlike entities and components),
117 -- I NEED to use this "dummy" placeholder even if we don't actually
118 -- care about the read values and only want to count their numbers.
119 alias dummy_t is read_t;
120 -- Be aware that you cannot combine 'sread()' (which is broken
121 -- anyway) with a non-static while-loop, or else Vivado will
122 -- still complain about using an unsynthesizable procedure.
123 --constant FAKE_LIMIT : natural := 2**16 - 1;

```



```

124
125 constant ROW_DELIMITER : read_t := (others => 'X');
126 --constant NULL_SLV      : std_logic_vector (0 downto 1) :=
127 --      (others => '0');
128
129
130 -- Helper macros; TODO: Make this package a generic based on these
131 function get_weights_file(
132     layer_path : in string;
133     layer_idx   : in natural
134 ) return string;
135
136 function get_biases_file(
137     layer_path : in string;
138     layer_idx   : in natural
139 ) return string;
140
141 impure function get_num_layers(
142     network_path : in string
143 ) return natural;
144
145 impure function get_network_dimensions(
146     network_dir : in string;
147     num_layers   : in natural
148 ) return dimensions_array;
149
150 function max(x : dimensions_array) return dimension;
151
152 end package file_parser_aux;
153
154 package body file_parser_aux is
155
156     function get_weights_file(
157         layer_path : in string;
158         layer_idx   : in natural
159     ) return string is
160     begin
161         return layer_path & "/weights_" &
162             natural'image(layer_idx) & ".dat";
163     end function get_weights_file;
164
165     function get_biases_file(
166         layer_path : in string;
167         layer_idx   : in natural
168     ) return string is
169     begin
170         return layer_path & "/biases_" &
171             natural'image(layer_idx) & ".dat";
172     end function get_biases_file;
173
174     -- Returns the number of layers in a network, depending on how
175     -- many 'weights' files were present in a given directory.
176     -- Unfortunately, as a side effect, it has to create an extra
177     -- file named "weights_{N+1}.dat" to bypass Vivado's limitations.
178     --
179     -- NOTE: Due to Vivado's nonstandard implementation of file_open(),
180     -- we cannot depend on "open_status = name_error," because Vivado
181     -- just quits whenever it encounters a file that cannot be opened.
182     -- SEE: support.xilinx.com/s/question/0D54U00008CO8pTSAT/
183     --      bug-fileopen-is-not-consistent-with-ieee-standards
184     impure function get_num_layers(
185         network_path : in string

```

```

186 ) return natural is
187     file      test_handle  : text;
188     variable open_status   : file_open_status;
189     variable file_no       : natural := 0;
190 begin
191     -- Exit when there are no more layers/files to process.
192     try_file : loop
193         -- NOTE: In VHDL-2008, there is no 'read_write_mode' as
194         -- in VHDL-2009; this is a rather lengthy workaround.
195         --
196         -- NOTE: Do NOT use 'write_mode'; it wipes the file out.
197         -- use 'append_mode' instead.
198         file_open(
199             open_status,
200             test_handle,
201             get_weights_file(network_path, file_no),
202             append_mode -- Create it if it doesn't exist
203         );
204         file_close(test_handle);
205
206         -- 'open_status' doesn't work in Vivado.
207         exit when (open_status /= open_ok);
208
209         file_open(
210             open_status,
211             test_handle,
212             get_weights_file(network_path, file_no),
213             read_mode -- Re-open it in readmode
214         );
215
216         was_empty : if endfile(test_handle) then
217             file_close(test_handle);
218             exit; -- Exit entirely
219         else
220             -- The file opened was valid & there exists a layer.
221             file_no := file_no + 1;
222         end if was_empty;
223
224         file_close(test_handle);
225     end loop try_file;
226
227     return file_no;
228 end function get_num_layers;
229
230 -- This function returns the number of rows and columns in a file.
231 impure function get_layer_dimension(
232     layer_dir : in string;
233     layer_idx : in natural
234 ) return dimension is
235     constant sample_file      : string :=
236         get_weights_file(layer_dir, layer_idx);
237     file      file_to_test      : text open read_mode is sample_file;
238     variable current_line      : line;
239     variable dummy_element     : dummy_t;
240     variable read_succeeded    : boolean;
241     variable n_rows            : natural := 1; -- No ending delimiter
242     variable n_cols            : natural := 0;
243     variable layer_dimension   : dimension;
244 begin
245
246     count_rows : while not endfile(file_to_test) loop
247         --count_cols : while not endfile(file_to_test) loop

```

```

248
249     readline(file_to_test, current_line);
250     read(current_line, dummy_element, read_succeeded);
251
252     if dummy_element = ROW_DELIMITER then
253         n_rows := n_rows + 1;
254     end if;
255
256     if n_rows = 1 then -- Avoid re-counting multiple times
257         n_cols := n_cols + 1;
258     end if;
259
260     --end loop count_cols;
261 end loop count_rows;
262
263 assert n_rows > 0
264     report "No rows existed in layer no. "
265         & natural'image(layer_idx) & "."
266         severity failure;
267 assert n_cols > 0
268     report "No columns existed in layer no. "
269         & natural'image(layer_idx) & "."
270         severity failure;
271
272     layer_dimension.rows := n_rows;
273     layer_dimension.cols := n_cols;
274
275     return layer_dimension;
276 end function get_layer_dimension;
277
278 -- This returns an array of rows and cols in a network's layers
279 impure function get_network_dimensions(
280     network_dir : in string;
281     num_layers  : in natural
282 ) return dimensions_array is
283     variable network_dimensions :
284         dimensions_array (0 to num_layers-1);
285 begin
286
287     per_layer : for idx in network_dimensions'range loop
288         network_dimensions(idx) :=
289             get_layer_dimension(network_dir, idx);
290     end loop per_layer;
291
292     return network_dimensions;
293 end function get_network_dimensions;
294
295
296 -- Returns the maximum of either columns# and rows# together
297 -- as a single dimension; takes an array of dimensions.
298 function max(x : dimensions_array) return dimension is
299     variable x_max : dimension := x(x'low); -- Start with the first
300 begin
301     -- Starting from 'low+1 because 'low was assigned beforehand.
302     comparator : for i in x'low+1 to x'high loop
303         x_max.rows := x(i).rows when (x(i).rows > x_max.rows);
304         x_max.cols := x(i).cols when (x(i).cols > x_max.cols);
305     end loop comparator;
306
307     return x_max;
308 end function max;
309

```

```

310 end package body file_parser_aux;
311
312
313 library std;
314     use std.textio.all;
315
316 library ieee;
317     use ieee.std_logic_1164.all;
318
319 library work;
320     use work.file_parser_aux.all;
321
322 library neural;
323     context neural.neural_context;
324
325 package file_parser is
326
327     impure function parse_network_from_dir(
328         network_dir : string
329     ) return network_layers;
330
331 end package file_parser;
332
333 -- NOTE: In VHDL, you are not able to pass or return ranges (e.g.,
334 -- '7 downto 0') to or from functions/procedures; so, the start/end
335 -- both have to be given SEPARATELY.
336 --
337 -- NOTE: In VHDL, you also cannot pass 'file' handles/objects; so,
338 -- the solution is to redundantly open/close the files EACH time.
339 package body file_parser is
340     -- NOTE Due to VHDL limitations, and to avoid allocators (i.e.,
341     -- 'access', 'new', and dynamic concatenation), which are
342     -- completely disallowed in synthesis (EVEN to just initialize
343     -- other constants) we have to know the number of rows/cols in
344     -- each file as a constant BEFORE calling these functions. In
345     -- other words, the '_dim' parameters are the actual dimensions
346     -- and not the ones we'd subsequently use for the Resizing Trick.
347     -- (Thanks to Mr. Brian Padalino for his dimension-measuring idea.)
348
349     -- NOTE: Make sure to review the following link:
350     -- support.xilinx.com/s/question/0D54U00008ADvMRSA1
351     -- TL;DR: Vivado (as of 2023.2) has a bug within its
352     -- 'read()' procedure of 'real' datatypes that prevents
353     -- reading out more than 7 times from a single line.
354     --
355     -- This means that "rows" hereinafter refer to a _flattened_ list
356     -- of elements on multiple lines, NOT multiple elements per line.
357
358
359     -- Variant for arrays (e.g., biases)
360     impure function parse_elements(
361         file_path : in string;
362         file_dim  : in dimension
363     ) return neural_vector is
364         file      file_handle : text open read_mode is file_path;
365         variable file_row     : line;
366         variable row_elem    : read_t; -- Intermediary for conversion
367         variable result_arr  : neural_vector (0 to file_dim.rows-1);
368     begin
369         -- No need for an outer loop; there is only ONE "row" of biases
370
371         parse_row : for col in result_arr'range loop

```

```

372
373     readline(file_handle, file_row);
374     read(file_row, row_elem);
375
376     -- Convert to the internally used fixed-point type.
377     result_arr(col) := to_sfixed(
378         row_elem, neural_word'high, neural_word'low
379     );
380
381     end loop parse_row;
382
383     return result_arr;
384 end function parse_elements;
385
386 -- Variant for nested arrays/matrices (e.g., weights)
387 impure function parse_elements(
388     file_path : in string;
389     file_dim  : in dimension
390 ) return neural_matrix is
391     file      file_handle : text open read_mode is file_path;
392     variable file_row     : line;
393     variable row_elem     : read_t; -- Intermediary for conversion
394     variable dummy_length : natural;
395     variable result_mat   : neural_matrix
396         (0 to file_dim.rows-1) (0 to file_dim.cols-1);
397 begin
398
399     parse_rows : for row in result_mat'range loop
400         parse_cols : for col in result_mat(0)'range loop
401
402             readline(file_handle, file_row);
403             read(file_row, row_elem);
404
405             -- If the row hasn't ended, convert and store the item.
406             result_mat(row)(col) := to_sfixed(
407                 row_elem, neural_word'high, neural_word'low
408             );
409
410         end loop parse_cols;
411
412         -- Skip the delimiters (except at EOF)
413         if not endfile(file_handle) then
414             readline(file_handle, file_row);
415         end if;
416     end loop parse_rows;
417
418     return result_mat;
419 end function parse_elements;
420
421
422 -- NOTE: A very frustrating limitation was that anything earlier
423 -- than VHDL-2019 cannot use 'block' in subprograms (e.g., function
424 -- bodies), and I really needed to use 'block's because they allow
425 -- you to have a varying number of locally scoped constant types.
426 -- Because I have an unknown (though not at compile-time) number of
427 -- layers to parse, I cannot simply hard-code 10 or 100 constant
428 -- types. So, one way was to constrain the entire holder (record)
429 -- inside the function's declaration region, but this proved a
430 -- challenge: I now could not return my value from the function in
431 -- any way.
432 --     Constraining a record inside a function's LOCAL declaration
433 -- area meant that I could not use it to define the function's own

```

```

434 -- return type, much less use it outside the function for the
435 -- caller's type. Other than using VHDL-2019 (which, as of 2024,
436 -- is not supported anywhere), a solution to this was using a
437 -- for-generate loop (which DOES have a locally scoped declaration
438 -- region) inside an entity. However, this meant having to
439 -- instantiate an entity each time you want to use the parser;
440 -- the syntax would look rather clunky that way.
441 -- Ultimately, I chose to take advantage of VHDL-2008's generic
442 -- packages and perform some of the processing (i.e., finding the
443 -- number of layers or their dimensions), which will be used by the
444 -- custom constrained record types, earlier in the package's header
445 impure function parse_network_from_dir(
446     network_dir : string
447 ) return network_layers is
448
449     constant NUM_LAYERS : natural :=
450         get_num_layers(network_dir);
451     constant LAYER_DIMS : dimensions_array :=
452         get_network_dimensions(network_dir, NUM_LAYERS);
453     constant MAX_DIM : dimension :=
454         max(LAYER_DIMS);
455
456     subtype constr_params_arr_t is network_layers
457         (0 to NUM_LAYERS-1)
458         (
459             weights (0 to MAX_DIM.rows-1) (0 to MAX_DIM.cols-1),
460             biases  (0 to MAX_DIM.rows-1)
461         );
462     -- SEE: https://gitlab.com/IEEE-P1076/VHDL-Issues/-/issues/312
463     subtype constr_params_t is constr_params_arr_t'element;
464
465     variable current_layer : constr_params_t;
466     variable result_arr : constr_params_arr_t;
467 begin
468
469     per_layer : for i in 0 to NUM_LAYERS-1 loop
470
471         current_layer.dims := LAYER_DIMS(i);
472         current_layer.biases := resize(parse_elements(
473             get_biases_file(network_dir, i), LAYER_DIMS(i)
474         ), MAX_DIM);
475         current_layer.weights := resize(parse_elements(
476             get_weights_file(network_dir, i), LAYER_DIMS(i)
477         ), MAX_DIM);
478
479         result_arr(i) := current_layer;
480     end loop per_layer;
481
482     return result_arr;
483 end function parse_network_from_dir;
484
485
486 end package body file_parser;
487
488
489 -- -----
490 -- END OF FILE: file_parser.vhd
491 -- -----

```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- pipeliner.vhd is a part of Innervator.
4  -- -----
5
6
7  -- For background, sometimes operations/logic inside a process might
8  -- take too long, due to the limitations of the physical world, and
9  -- therefore not be ready within a single clock cycle; this is referred
10 -- to as "timing violation" in FPGA design, and it is especially
11 -- prevalent when using too many combination logic and/or having too
12 -- fast of a clock (e.g., 500MHz -> 2 ns). So, a workaround is to
13 -- "delay" (i.e., spread) the processing of data over multiple clock
14 -- cycles; this is done by assigning inputs/outputs to one or more
15 -- "register" signals, which makes said ins/outs reach the outside
16 -- entities by multiple clock cycles (e.g., 50 ns instead of 10 ns)
17 -- as to give the FPGA fabric more time to complete the operations
18 -- we ask it to. Also, certain built-in blocks, such as DSPs, will
19 -- "pull in" our externally registered (pipelined) signals and have
20 -- the same delaying effect.
21 --     Another type of delay is "routing" delay, which is actually
22 -- more common nowadays and in newer FPGAs. Routing delay refers
23 -- to the actual, physical delay caused by electricity moving too
24 -- "slowly" in an internal FPGA route (i.e., wire), and it is often
25 -- caused when you try to access specific hard-macro components
26 -- (e.g., a DSP or BRAM) that only exists in a specific location
27 -- of your FPGA chip, while you have too much logic depend on that,
28 -- which, in turn, makes it physically impossible for the router to
29 -- place them all next to said hard-macro instantiation. The solution
30 -- is exactly the same: convert the routed path to a "multi-cycle"
31 -- clock using flip-flop chains.
32 --     Note that delaying does imply that external users have to
33 -- "wait" for the processing to be finished before they can supply
34 -- new data to us; while the old data is going through flip-flop
35 -- chains, new data could follow after just a single clock cycle.
36 -- If the data D1 is given at nanosecond 0 and D2 at 10, then
37 -- (in a double-registering pipeline), output O1 would be returned
38 -- at nanosecond 20 AND O2 would STILL appear at 30. In other words,
39 -- we didn't have to wait for O1 to be fully done before passing D2,
40 -- and you could think of it as a car assembly line.
41 --     Be aware that timing violations do not 100% mean that your
42 -- design won't work in practice. Synthesis tools account for all
43 -- extremities when calculating timing violations; they consider
44 -- worst-case scenarios and ascertain that your described logic will
45 -- finish under those circumstances within the allocated clock cycle.
46 -- If you proceed with having tiny timing errors, your design might
47 -- work just fine, but the moment your FPGA package gets too heated
48 -- or too cold (for example), it'll have a higher tendency to fail.
49 --     Also, this closely relates to a "synchronizer," which also
50 -- uses delay lines to solve a different problem (metastability).
51 --     Lastly, you might have to enable "retiming" optimizations in
52 -- your synthesis tools so that they may take advantage of pipelining.
53
54 library ieee;
55     use ieee.std_logic_1164.all;
56
57 -- NOTE: Currently, there is a language inconsistency (reported by me)
58 -- in VHDL that disallows output ports from being able to use
59 -- "aggregated" expressions (i.e., tuples) such as (1, 2, 3) or
60 -- (A, B, C):
61 --     https://gitlab.com/IEEE-P1076/VHDL-Issues/-/issues/311

```



```

62 --
63 -- had they been allowed, the entity could have been used like so:
64 --
65 --     -- Here, we expect that the entity constructs a "trit vector"
66 --     -- type internally, using it to denote its input/output signals.
67 --     test_instance : entity work.pipeliner
68 --         generic map (2, trit)
69 --         port map (
70 --             (alpha, beta, gamma), -- This works, as of VHDL-08.
71 --             (alpha_reg, beta_reg, gamma_reg) -- EXPRESSION ERROR!
72 --         );
73 --
74 -- However, because dis-aggregation is not done on 'out' ports, and
75 -- because arrays cannot be constructed based on a base type in VHDL
76 -- versions earlier than 2019, the next interface could have been this:
77 --
78 --     -- trit_vector is defined by us earlier in the Architecture.
79 --     test_instance : entity work.pipeliner
80 --         generic map (2, trit_vector, 3)
81 --         port map (
82 --             (alpha, beta, gamma), -- This works, as of VHDL-08.
83 --             o_signals(0) => alpha_reg,
84 --             o_signals(1) => beta_reg,
85 --             o_signals(2) => gamma_reg
86 --         );
87 --
88 -- Sadly, you cannot "convince" VHDL's strong typing system that the
89 -- generically provided type is, in fact, an array that you can index
90 -- over; this means that the next alternative was either "dividing" the
91 -- interface to take/return one signal at a time, or continue with
92 -- tuples but provide a duplicate overload for every single array type.
93 --
94 -- NOTE: The aforementioned topic also applies equally to generic
95 -- subprograms (i.e., generically typed procedures) with the difference
96 -- being that they cannot have indexed/composite port assignments
97 -- whatsoever; the compiler complains about non-static (?) ports.
98 --
99 -- NOTE: Subprograms (like procedures) can take generic arguments,
100 -- similar to entities. However, a nuanced difference is that an
101 -- entity's 'generic(...);' and 'port(...);' clauses are statements
102 -- (hence the semicolon; after them), while 'generic(...)' and
103 -- 'parameter(...)' are NOT statements and do not have semicolons.
104 -- This also means that we cannot declare constants in-between the
105 -- 'procedure' and 'is' keywords, unlike entities.
106 package pipeliner is
107     procedure registrar
108         generic (
109             constant g_NUM_STAGES :
110                 in natural range 2 to natural'high := 3;
111             type t_ARG_TYPE
112         )
113     parameter (
114         signal i_clk      : in std_ulogic;
115         signal i_signal    : in t_ARG_TYPE;
116         signal o_signal    : out t_ARG_TYPE
117     );
118
119 end package pipeliner;
120
121 -- Instantiate in the following format:
122 --
123 --     procedure register_signal is new core.pipeliner.registrar

```



```

124 --      generic map (3, std_logic_vector);
125 --
126 --      register_signal(i_clk, sig_unreg, sig_reg);
127 --
128 package body pipeliner is
129
130     -- Works on a single signal at a time
131     procedure registrar -- Singular
132         generic (
133             constant g_NUM_STAGES :
134                 in natural range 2 to natural'high := 3;
135             type t_ARG_TYPE
136         )
137         parameter (
138             signal i_clk      : in std_ulogic;
139             signal i_signal   : in t_ARG_TYPE;
140             signal o_signal   : out t_ARG_TYPE
141         )
142     is
143         type t_arg_arr is array (natural range <>) of t_ARG_TYPE;
144         -- NOTE: Even though this is a variable, and variables are
145         -- purported to update their values immediately, these are,
146         -- in reality, no different from normal signals in this case.
147         -- Unlike simulation, hardware is not 100% perfect and a
148         -- daisy-chain of variables mean that the following variables
149         -- would be longer "wires" connecting them to preceding ones,
150         -- meaning that there could be a delay (albeit very tiny)
151         -- between the time it takes for the preceding variables to
152         -- update their own values and the time it takes for variables
153         -- later in the chain to take effect of said updates.
154         -- I find a lot of tutorials and beginners' guides very
155         -- misleading as a result of this; had such pitfalls been
156         -- mentioned, variables would not be treated so arcanelly.
157         -- Also, since variables could effectively be used to
158         -- replicate signals' delayed assignment behavior in synthesis,
159         -- one cannot help but wonder why VHDL just doesn't let us
160         -- declare signals directly in the first place; the reason
161         -- is rather arbitrary and lies in (now-ancient) design
162         -- choices of VHDL and Ada.
163         variable pipeline_regs : t_arg_arr (g_NUM_STAGES-2 downto 0);
164
165     begin
166         pipeline_regs(0) := i_signal when rising_edge(i_clk);
167
168         for i in 1 to g_NUM_STAGES-2 loop
169             pipeline_regs(i) :=
170                 pipeline_regs(i - 1) when rising_edge(i_clk);
171         end loop;
172
173         o_signal <= pipeline_regs(pipeline_regs'high);
174     end procedure registrar;
175
176
177     -- NOTE: The rationale behind using array types was to
178     -- replicate the same "variadic" number of parameters
179     -- that we have in C.
180     --
181     -- TODO: Whenever VHDL-2019 is widespread, make this take
182     -- in a base type and construct its array ITSELF; this
183     -- way, users won't have to define an otherwise-unused
184     -- array type just to pass it to this entity, and they can
185     -- use a single entity instantiation for multiple signals.

```

```

186  --
187  -- https://gitlab.com/IEEE-P1076/VHDL-Issues/-/issues/311
188  /*
189  procedure register_signals -- Plural
190      generic (
191          constant g_NUM_STAGES :
192              in natural range 2 to natural'high := 3;
193
194          type t_ELEMENT is private;
195          type t_ARRAY is array (natural) of t_ELEMENT
196      )
197      parameter (
198          signal i_clk      : in std_ulogic;
199          signal i_signals : in  t_ARRAY (0 downto 0);
200          signal o_signals : out t_ARRAY (0 downto 0)
201      )
202  is
203      -- NOTE: We subtract by 2 because the incoming (input) signal
204      -- itself also counts; when we talk about "double registering"
205      -- signals, we know that there already WAS a flip-flop register
206      -- (i.e., the 'in' signal) and we need to add just 1 more to it
207      constant NUM_ELEMS : positive := i_signals'length;
208      -- Constrained subtypes of the array's elements.
209      subtype t_CONSTRAINED is t_ARG_TYPE (i_signals'element'range);
210
211      procedure register_custom_typed is new registrar
212          generic map (g_NUM_STAGES, t_CONSTRAINED);
213  begin
214      -- Compile-time assertions will go here to ascertain
215      -- that the 'input' and 'output' are of equal lengths.
216      --assert i_signals'length = o_signals'length
217          --report "Aggregated input/output tuples " &
218          --      "are not equal in size."
219          --severity failure;
220
221      per_signal : for i in 0 to NUM_ELEMS-1 loop
222          register_custom_typed(i_clk, i_signals(i), o_signals(i));
223      end loop per_signal;
224  end procedure register_signals;
225  */
226
227 end package body pipeliner;
228
229
230 library ieee;
231 use ieee.std_logic_1164.all;
232
233 entity pipeliner_single is
234     generic (
235         g_NUM_STAGES : positive := 2;
236         type arg_type
237     );
238     port (
239         i_clk      : in std_ulogic;
240         i_signal : in  arg_type;
241         o_signal : out arg_type
242     );
243
244     type arg_arr_type is array (natural range <>) of arg_type;
245 end entity pipeliner_single;
246
247 architecture behavioral of pipeliner_single is

```

```

248 -- Synthesis tools will often replace a series of flip-flops
249 -- with better primitives, like shift-registers, that "achieve"
250 -- the same delaying effect. However, we might sometimes WANT
251 -- to use flip-flops specifically, so we can turn off that
252 -- optimization by using vendor-specific attribute definitions.
253 --
254 -- Also, note that avoiding the usage of explicit reset signals
255 -- may also result in the same shift-register (SRL) conversion.
256 -- ednasia.com/coding-consideration-for-pipeline-flip-flops
257
258 /* Xilinx ISE */
259 attribute register_balancing : string;
260 /* Xilinx Vivado */
261 attribute shreg_extract : string; -- No shift-reg. conversion
262 /* Altera Quartus */
263 attribute syn_allow_retiming : boolean;
264
265 -- TODO: Do we also apply these to input/output or just cascades?
266 attribute register_balancing of
267     i_signal : signal is "backward";
268 attribute shreg_extract of
269     i_signal : signal is "no";
270 attribute syn_allow_retiming of
271     i_signal : signal is true;
272 -- Output
273 attribute register_balancing of
274     o_signal : signal is "backward";
275 attribute shreg_extract of
276     o_signal : signal is "no";
277 attribute syn_allow_retiming of
278     o_signal : signal is true;
279 begin
280
281 -- NOTE: Concurrent assignments mean that there is no
282 -- delay involved; both wires "connect" and act as one.
283 --
284 -- "Delaying" by 1 stage means that we will not have
285 -- to place any additional flip-flops in the middle.
286 -- The reason is that the unregistered i_signal
287 -- would also have a propagation delay of 1 whenever
288 -- something is assigned to it.
289 single_pipeline : if g_NUM_STAGES = 1 generate
290     o_signal <= i_signal; -- CONCURRENT assignment
291 -- Otherwise, implement actual delay with flip-flops.
292 else generate
293     constant STAGES_HIGH : natural := g_NUM_STAGES - 2;
294     signal pipeline_regs : arg_arr_type
295         (0 to STAGES_HIGH);
296
297     attribute register_balancing of
298         pipeline_regs : signal is "backward";
299     attribute shreg_extract of
300         pipeline_regs : signal is "no";
301     attribute syn_allow_retiming of
302         pipeline_regs : signal is true;
303 begin
304     pipeline_chain : for j in 0 to STAGES_HIGH generate
305         pipeline_register : process (i_clk) begin
306             if rising_edge(i_clk) then
307                 pipeline_regs(j) <= i_signal when j = 0
308                     else pipeline_regs(j-1);
309             end if;

```

```
310         end process pipeline_register;
311     end generate pipeline_chain;
312
313     o_signal <= pipeline_regs(STAGES_HIGH); -- CONCURRENT assignment
314 end generate;
315
316 end architecture behavioral;
317
318
319
320 -- -----
321 -- END OF FILE: pipeliner.vhd
322 -- -----
```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- neuron.vhd is a part of Innervator.
4  -- -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 library work;
11     context work.neural_context;
12
13 library config;
14     use config.constants.all;
15
16 library core;
17
18 -- TODO: Generically 'type' these, at least in VHDL-2019.
19 entity neuron is
20     generic (
21         g_NEURON_WEIGHTS : neural_wvector;
22         g_NEURON_BIAS     : neural_word;
23         /* Sequential (pipeline) controllers */
24         -- Number of inputs to be processed at a time (default = all)
25         g_BATCH_SIZE      : positive := g_NEURON_WEIGHTS'length;
26         -- Number of pipeline stages/cycle delays (default = none)
27         g_PIPELINE_STAGES : natural := 0
28     );
29     -- NOTE: There are also other types such as 'buffer' and the
30     -- lesser-known 'linkage' but they are very situation-specific.
31     -- NOTE: Apparently, it is better to use Active-High (like o_done
32     -- when '1' instead of o_busy '1') internal signals in most FPGAs.
33     port (
34         -- NOTE: Do NOT name these as mere 'input' or 'output' because
35         -- std.textio also defines those, and they can conflict.
36         i_inputs : in  neural_bvector (0 to g_NEURON_WEIGHTS'length-1);
37         o_output : out neural_bit; -- The "Action Potential"
38         /* Sequential (pipeline) controllers */
39         i_clk    : in  std_ulogic; -- Clock
40         i_rst    : in  std_ulogic; -- Reset
41         i_fire   : in  std_ulogic; -- Start/fire up the neuron
42         o_done   : out std_ulogic -- Are we done processing the batch?
43     );
44
45     constant NUM_INPUTS : positive := i_inputs'length;
46 begin
47
48     -- TODO: See if the 'instance_name or _path attributes are
49     -- supported in Vivado synthesis; if so, use them here.
50     assert g_NEURON_WEIGHTS'length mod g_BATCH_SIZE = 0
51         report "Size of input data is not evenly divisble " &
52             "by the given batch size."
53             severity failure;
54
55 end entity neuron;
56
57
58 -- TODO: While anything with registers and flip-flops can be called
59 -- a pipeline, it might still not be very appropriate to call this
60 -- a pipelined neuron, because the actual registers are used to
61 -- resolve routing delays and, essentially, function as multi-cycle

```

```

62 -- paths. (Revamp this to actually function like a pipeline.)
63 architecture pipelined of neuron is
64     -- TODO: For UNJUSTIFIED reasons, you cannot declare a signal
65     -- within a process' declaratory section, and yet variables can
66     -- be FUNCTIONALLY the same, in case of counters. While you
67     -- could use 'block' clauses to wrap the process and have
68     -- "locally scoped" signals that way, it is still going to add
69     -- an extra indention nest, and it is not very ideal.
70     -- TODO: Decide if we want to move these otherwise-local signals
71     -- to become variables in their respective processes. (One dis-
72     -- advantage is that simulators might not show them in waveforms.)
73
74     -- This is the pipeline's propagation delay counter. For example,
75     -- in a pipeline with 3 stages, the "external" input will take
76     -- 3 clock cycles to arrive "inside" the pipelined processor,
77     -- and said processor's output (back to the external source)
78     -- would also take 3 clock cycles to reach.
79     -- Hence, we need to wait (i.e., count each cycle) until the
80     -- data is fully loaded into the pipeline, in both directions.
81     -- This might need a re-design to be more clear; the reason
82     -- it starts at 1 is that the first batch in the pipeline will
83     -- have already gotten filled in the "pre-processing" stage.
84     signal pipe_delay : natural range 0 to g_PIPELINE_STAGES-1 := 1;
85     -- Current iteration number (total number of
86     -- iterations = number of data / size of batches).
87     --
88     -- Because the pipeline is always "ahead" of the processing
89     -- multiplier by the number of stages (in clock cycles),
90     -- two separate indices are used to keep track.
91     signal pipe_iter_idx : natural range 0 to NUM_INPUTS := 0;
92     signal proc_iter_idx : natural range 0 to NUM_INPUTS := 0;
93
94     -- The Neuron's Finite-State Machine (FSM)
95     type neuron_state_t is (
96         idle, initializing, processing, finalizing, activating, done
97     );
98     signal neuron_state : neuron_state_t := idle;
99
100     -- This is the "localized" version of the input signal; given that
101     -- our given input itself might reset or become cleared right after
102     -- 'i_fire' is set to high, we need to sample and locally store the
103     -- input at that time for later use within the processing stages.
104     signal inputs_local : i_inputs'subtype;
105     -- These are the unregistered input and registered output signals.
106     -- TODO: Somehow use 'subtype and 'element to derive these
107     -- Input data (DSP input 1)
108     signal inputs_unreg : neural_bvector (0 to g_BATCH_SIZE-1);
109     signal inputs_reg : neural_bvector (0 to g_BATCH_SIZE-1);
110     -- Weights (DSP input 2)
111     signal weights_unreg : neural_wvector (0 to g_BATCH_SIZE-1);
112     signal weights_reg : neural_wvector (0 to g_BATCH_SIZE-1);
113     -- Internal DSP multiplier (product) pipeline
114     --
115     -- NOTE: If you are using a very small (i.e., < 4) number of
116     -- bits for either the integral or fractional part, you may
117     -- consider using a slightly larger multiple of the _word type
118     -- (e.g., neural_word or neural_qword) here for the accumulator
119     -- to accomodate for the many additions that occur within
120     -- the inner for-loop and would otherwise overflow. After the
121     -- activation function/clamping takes place, and the variable
122     -- gets its range restricited within [0, 1), we can safely
123     -- resize it back to a smaller bit width.

```

```

124  -- Also, this might result in the synthesizer using
125  -- available DSP (dedicated multiplier) blocks on your
126  -- FPGA, which would conserve other logic resources.
127  signal products_unreg    : neural_dvector (0 to g_BATCH_SIZE-1);
128  signal products_reg      : neural_dvector (0 to g_BATCH_SIZE-1);
129  -- Multiplied-Accumulated weighted sum (DSP output)
130  signal outputs_unreg     : neural_dvector (0 to g_BATCH_SIZE-1);
131  signal outputs_reg       : neural_dvector (0 to g_BATCH_SIZE-1);
132  -- The activation function (not batched)
133  signal activation_unreg   : neural_bit;
134  signal activation_reg     : neural_bit;
135
136  function activation_function(
137      x : neural_dword
138  ) return neural_bit is
139  begin
140      -- TODO: Automatically select between activations, as needed.
141      return work.activation.sigmoid(x);
142  end function activation_function;
143
144
145  -- Yet another VHDL annoyance:
146  -- stackoverflow.com/questions/31044965/
147  -- procedure-call-in-loop-with-non-static-signal-name
148  -- In short, procedures cannot take elements from array
149  -- signals, such as test_sig(i), even if the index is
150  -- static. For that reason, we have to take indices
151  -- rather than pre-indexed array elements.
152  --
153  -- Multiplier-Accumulator ("MAC")
154  procedure multiply_accumulate(
155      constant idx          : in natural;
156      signal mul_a           : in neural_wvector;
157      signal mul_b           : in neural_bvector;
158      signal acc_in          : in neural_dvector;
159      signal acc_out         : out neural_dvector;
160      signal prod_unreg      : out neural_dvector;
161      signal prod_reg        : in neural_dvector
162  ) is
163      variable products     : neural_dword;
164      variable summation    : neural_dword;
165
166      -- NOTE: Unfortunately, procedures' 'out' parameters
167      -- cannot be assigned to 'open', unlike actual entities/
168      -- components; use dummies or placeholders as a workaround.
169      variable dummy_carry  : std_ulogic;
170  begin
171      products := resize(
172          mul_a(idx)
173          * -- Multiply
174          resize( -- Resize, if needed
175              to_sfixed(mul_b(idx)),
176              mul_a(idx)),
177          acc_in(idx));
178
179      prod_unreg(idx) <= products;
180
181      add_carry(
182          L      => acc_in(idx),
183          R      => prod_reg(idx),
184          c_in   => '0',
185          result => summation,

```

```

186         c_out => dummy_carry -- IGNORED!
187     );
188
189     acc_out(idx) <= summation;
190 end procedure multiply_accumulate;
191
192 -- These cause simulation/synthesis mismatch
193 /*
194 procedure register_inputs is new core.pipeliner.registrar
195     generic map (2, inputs'element);
196
197 procedure register_weights is new core.pipeliner.registrar
198     generic map (2, g_NEURON_WEIGHTS'element);
199 */
200 begin
201     -- NOTE: This form of pipelining would only fix timing issues
202     -- related to physical routing, not resource/logic consumption.
203     --
204     -- When g_PIPELINE_STAGES is 0, the _reg output and _unreg
205     -- input get concurrently connected (with no delay).
206     no_pipeline : if g_PIPELINE_STAGES = 0 generate
207         create_pipeline : for i in 0 to g_BATCH_SIZE-1 generate
208             inputs_reg(i)    <= inputs_unreg(i);
209             weights_reg(i)    <= weights_unreg(i);
210             products_reg(i)   <= products_unreg(i);
211             outputs_reg(i)    <= outputs_unreg(i);
212         end generate create_pipeline;
213
214         activation_reg        <= activation_unreg;
215
216     else generate
217         create_pipeline : for i in 0 to g_BATCH_SIZE-1 generate
218             register_inputs    : entity core.pipeliner_single
219                 generic map (g_PIPELINE_STAGES, neural_bit)
220                 port map (i_clk, inputs_unreg(i), inputs_reg(i));
221             register_weights   : entity core.pipeliner_single
222                 generic map (g_PIPELINE_STAGES, neural_word)
223                 port map (i_clk, weights_unreg(i), weights_reg(i));
224             register_products  : entity core.pipeliner_single
225                 generic map (g_PIPELINE_STAGES, neural_dword)
226                 port map (i_clk, products_unreg(i), products_reg(i));
227             register_outputs   : entity core.pipeliner_single
228                 generic map (g_PIPELINE_STAGES, neural_dword)
229                 port map (i_clk, outputs_unreg(i), outputs_reg(i));
230         end generate create_pipeline;
231
232         register_activation    : entity core.pipeliner_single
233             generic map (g_PIPELINE_STAGES, neural_bit)
234             port map (i_clk, activation_unreg, activation_reg);
235
236     end generate;
237
238     -- NOTE: Combinational (i.e., un-clocked and stateless) logic,
239     -- sensitive only to changes in inputs, will be much "faster" and
240     -- perform everything in a SINGLE clock cycle. However, it will
241     -- also use a much, much higher number of logic blocks in the FPGA,
242     -- meaning that a single neuron with 64 inputs could potentially
243     -- take up 10% of a small FPGA's (e.g., Artix-7) LUTs.
244     --
245     -- A workaround is to convert the combination logic to a
246     -- sequential (i.e., clocked and stateful) one, where weighted sums
247     -- are calculated in small "batches" in each clock cycle; this
248     -- does have the disadvantage of requiring MULTIPLE clock cycles

```



```

248  -- for the entire calculation to be done (e.g., for 64 inputs and
249  -- a 100MHz clock, the combinational approach would take 10ns while
250  -- the sequential one, with segments of 2, might take 320+20ns).
251  --     Additionally, if you go with the combination approach while
252  -- keeping track of the previous states, you can introduce latches.
253  -- Lastly, if you go with the sequential approach, you also need to
254  -- have additional communication mechanism with the external logic
255  -- to let them know whenever this neuron is done processing its
256  -- batch or whenever it should begin processing the given batch;
257  -- otherwise, since your sequential process is already clocked and
258  -- you can't use the 'input' signal's event in its sensitivity list
259  -- you would have to maintain its previous states and compare them.
260  --
261  -- TODO: Improve the pipelining to actually overlap and be
262  -- continuously fed from the input data.
263  neuron_loop : process (i_clk, i_rst) is
264      -- NOTE: Use 'variable' as opposed to a 'signal' because these
265      -- for-loops are supposed to unroll inside a _single_ tick of
266      -- the Process, meaning that any subsequent assignments to
267      -- a 'singal' accumulator would be DISCARDED; by using
268      -- variables, we can resolve this issue.
269      --
270      -- NOTE: Somehow, using an initial value here adds a huge
271      -- ~3 ns setup timing slack; this should NOT be happening!
272      variable weighted_sum : neural_dword; --:=
273      --resize(g_NEURON_BIAS, neural_dword'high, neural_dword'low);
274
275      -- Number of iterations (if batch processing is enabled)
276      constant ITER_HIGH : natural := NUM_INPUTS - g_BATCH_SIZE;
277
278      procedure perform_reset is
279      begin
280          neuron_state <= idle;
281      end procedure perform_reset;
282  begin
283
284      if not c_RST_SYNC and i_rst = c_RST_POLE then perform_reset;
285      elsif rising_edge(i_clk) then
286          if c_RST_SYNC and i_rst = c_RST_POLE then perform_reset;
287          else
288
289              case neuron_state is
290              when idle =>
291                  neuron_state <= idle;
292
293                  -- Reset back to default values
294                  proc_iter_idx <= 0;
295                  pipe_iter_idx <= g_BATCH_SIZE; -- 0 in the loop
296                  pipe_delay <= 1; -- 1 delay's accounted here
297
298                  o_output <= to_ufixed(0, o_output);
299                  o_done <= '0';
300
301                  weighted_sum := -- Start with the Bias
302                  resize(g_NEURON_BIAS, weighted_sum);
303
304                  inputs_local <= (others =>
305                      to_ufixed(0, inputs_local'element'high,
306                      inputs_local'element'low)
307                  );
308                  -- Because we will be "adding" these before
309                  -- they are truly filled with calculated

```

```

310      -- values, we initialize them to a known
311      -- (i.e., 0) value at first.
312      products_unreg <= (others =>
313          to_sfixed(0, products_unreg'element'high,
314              products_unreg'element'low)
315      );
316      outputs_unreg <= (others =>
317          to_sfixed(0, outputs_unreg'element'high,
318              outputs_unreg'element'low)
319      );
320      /*
321      inputs_unreg <= (others =>
322          to_ufixed(0, inputs_unreg'element));
323      weights_unreg <= (others =>
324          to_sfixed(0, weights_unreg'element));
325      */
326
327      if i_fire = '1' then
328          -- Save the external inputs (which
329          -- can change after i_fire).
330          inputs_local <= i_inputs;
331
332          for i in 0 to g_BATCH_SIZE-1 loop
333              inputs_unreg(i) <= i_inputs(i);
334              weights_unreg(i) <= g_NEURON_WEIGHTS(i);
335          end loop;
336
337          -- NOTE: when pipeline stages == 1,
338          -- skip initializing and do the processing.
339          -- Switching cases counts as 1 delay itself
340          if g_PIPELINE_STAGES = 1 then
341              neuron_state <= processing;
342          else
343              neuron_state <= initializing;
344          end if;
345
346      end if;
347      -----
348
349      -- Here, we "wait" (for a number of clock cycles
350      -- equal to pipeline stages) so that the first data
351      -- arrives through the pipeline; otherwise, the
352      -- weighted sum would have uninitialized values.
353      when initializing =>
354          neuron_state <= initializing;
355
356          -- Even though we are "waiting," we should still
357          -- continue to fill the upcoming pipeline stages
358          --
359          -- TODO: Account for the scenario where the
360          -- pipeline stages exceed the number of inputs
361          -- (stop filling the pipeline at that point.)
362          for i in 0 to g_BATCH_SIZE-1 loop
363              inputs_unreg(i) <=
364                  inputs_local(pipe_iter_idx+i);
365              weights_unreg(i) <=
366                  g_NEURON_WEIGHTS(pipe_iter_idx+i);
367          end loop;
368          pipe_iter_idx <= pipe_iter_idx + g_BATCH_SIZE;
369
370          if (pipe_delay < g_PIPELINE_STAGES-1) then
371              pipe_delay <= pipe_delay + 1;

```

# Innervator

```

372         else
373             pipe_delay    <= 0; -- Reuse for 'finalizing'
374             neuron_state <= processing;
375         end if;
376     -- -----
377
378     when processing =>
379         neuron_state <= processing;
380
381         -- TODO: Have a generic switch to toggle
382         -- between computing the activation function
383         -- DURING the last batch iteration (1 clock
384         -- cycle less latency), or AFTER a clock cycle
385         -- passes, like now (better logic timing).
386
387         -- NOTE: The pipeline still needs to be filled
388         -- at g_NUM_STAGES ahead-of-time, but it will
389         -- also need to stop sooner; this is why we
390         -- keep track and stop it separately.
391         -- TODO: Account for when batch processing's off
392         pipeline_unsaturated : if
393             pipe_iter_idx < NUM_INPUTS
394         then
395             for i in 0 to g_BATCH_SIZE-1 loop
396                 -- Continue filling the pipeline, which
397                 -- will eventually reach the multiplier
398                 -- after a number of clock cycles (i.e.,
399                 -- pipeline stages).
400                 inputs_unreg(i) <=
401                     inputs_local(pipe_iter_idx+i);
402                 weights_unreg(i) <=
403                     g_NEURON_WEIGHTS(pipe_iter_idx+i);
404             end loop;
405             pipe_iter_idx <=
406                 pipe_iter_idx + g_BATCH_SIZE;
407         end if pipeline_unsaturated;
408
409
410         -- NOTE: This loop is unrolled into actual
411         -- hardware; this is why we don't multiply
412         -- an entire matrix all in one pass (it
413         -- would be far too much in one clock cycle)
414         for i in 0 to g_BATCH_SIZE-1 loop
415             -- This is a running accumulator; the
416             -- result of the multiplication of each
417             -- weight by its associated input is
418             -- resized (IF NEEDED) to the size of
419             -- the Accumulator and then added to it
420             multiply_accumulate(
421                 -- VHDL limitation workaround
422                 idx      => i,
423                 -- Numbers to multiply
424                 mul_a     => weights_reg,
425                 mul_b     => inputs_reg,
426                 -- Accumulator
427                 acc_in    => outputs_reg,
428                 acc_out   => outputs_unreg,
429                 -- Internal multiplier pipeline
430                 prod_unreg => products_unreg,
431                 prod_reg  => products_reg
432             );
433         end loop;

```

# Innervator

```

434     proc_iter_idx <= proc_iter_idx + g_BATCH_SIZE;
435
436     -- NOTE: Short-circuited to one at compile-time
437     sum_calculated : if
438         -- Not processing in batches (iterate once)
439         (ITER_HIGH = 0 and
440          proc_iter_idx /= 0) or
441         -- OR: Processing in batches
442         (ITER_HIGH /= 0 and
443          proc_iter_idx >= ITER_HIGH)
444     then
445         neuron_state <= finalizing;
446     end if sum_calculated;
447 -----
448
449     -- Here, similar to initializing, we will wait for
450     -- the pipelined OUTPUT to "catch up" and finish.
451     when finalizing =>
452         neuron_state <= finalizing;
453
454         -- TODO: Do we also pipeline this?
455         for i in 0 to g_BATCH_SIZE-1 loop
456             weighted_sum := resize(
457                 weighted_sum
458                 + outputs_reg(i)
459                 + products_reg(i),
460                 weighted_sum);
461         end loop;
462
463         if (pipe_delay < g_PIPELINE_STAGES-1) then
464             pipe_delay <= pipe_delay + 1;
465         else
466             pipe_delay <= 0;
467             if g_PIPELINE_STAGES = 1 then
468                 neuron_state <= done;
469             else
470                 neuron_state <= activating;
471             end if;
472         end if;
473 -----
474
475     -- Wait for activation on the weighted sum.
476     -- (It is done here to avoid logic/gate delay
477     -- that'd occur in the previous case's branch,
478     -- because it also has to be pipelined.)
479     when activating =>
480         -- Activate only once
481         if pipe_delay = 0 then
482             activation_unreg <=
483                 activation_function(weighted_sum);
484         end if;
485
486         if (pipe_delay < g_PIPELINE_STAGES-1) then
487             pipe_delay <= pipe_delay + 1;
488         else
489             neuron_state <= done;
490         end if;
491 -----
492
493     when done =>
494         o_output <= activation_reg;
495

```

## Innervator

```
496         -- Signal that we're no longer busy,  
497         -- for a single clock cycle  
498         o_done      <= '1';  
499  
500         -- Back to idle state; await new data  
501         neuron_state <= idle;  
502         -----  
503  
504         -- Hardening in case of "unknown" states  
505         when others => -- 1-clock-long cleanup phase  
506             perform_reset;  
507         -----  
508     end case;  
509  
510     end if;  
511 end if;  
512 end process neuron_loop;  
513  
514  
515  
516 end architecture pipelined;  
517  
518  
519 -----  
520 -- END OF FILE: neuron.vhd  
521 -----
```

```

1  -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- layer.vhd is a part of Innervator.
4  -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 library work;
11     context work.neural_context;
12
13 library config;
14     use config.constants.all;
15
16
17 -- TODO: Generically 'type' these, at least in VHDL-2019.
18 entity layer is
19     generic (
20         g_LAYER_WEIGHTS      : neural_matrix;
21         g_LAYER_BIASES       : neural_vector;
22         /* Sequential (pipeline) controllers */
23         -- Number of inputs to be processed at a time (default = all)
24         g_BATCH_SIZE          : positive := g_LAYER_WEIGHTS'element'length;
25         -- Number of pipeline stages/cycle delays (default = none)
26         g_PIPELINE_STAGES    : natural := 0
27     );
28     port (
29         i_inputs  : in  neural_bvector
30             (0 to g_LAYER_WEIGHTS'element'length-1);
31         o_outputs : out neural_bvector
32             (0 to g_LAYER_WEIGHTS'length-1);
33         /* Sequential (pipeline) controllers */
34         i_clk      : in  std_ulogic; -- Clock
35         i_rst      : in  std_ulogic; -- Reset
36         i_fire     : in  std_ulogic; -- Start/fire up all the neurons
37         o_done     : out std_ulogic  -- Is the layer processing done?
38     );
39
40     -- NOTE: This assumes that the upper hierarchy (i.e., network)
41     -- supplies the sanitized/actual slices of the layer's parameters.
42     constant NUM_INPUTS   : positive := i_inputs'length;
43     -- NOTE: This one is also the number of "neurons" in this layer.
44     constant NUM_OUTPUTS  : positive := o_outputs'length;
45 end entity layer;
46
47
48 architecture dense of layer is -- [Structural arch.]
49     signal neurons_done :
50         std_ulogic_vector (0 to NUM_OUTPUTS-1);
51 begin
52
53     neural_layer : for i in 0 to NUM_OUTPUTS-1 generate
54         neuron_instance : entity work.neuron
55             generic map (
56                 g_NEURON_WEIGHTS => g_LAYER_WEIGHTS(i),
57                 g_NEURON_BIAS    => g_LAYER_BIASES(i),
58                 g_BATCH_SIZE     => g_BATCH_SIZE,
59                 g_PIPELINE_STAGES => g_PIPELINE_STAGES
60             )
61         port map (

```

```
62         i_inputs => i_inputs,
63         o_output => o_outputs(i),
64         i_clk     => i_clk,
65         i_rst     => i_rst,
66         i_fire    => i_fire,
67         o_done    => neurons_done(i)
68     );
69 end generate neural_layer;
70
71 -- TODO: Decide if the Layer's busy signal should be based
72 -- on ALL neurons' busy signals (or just one of them?)
73 --o_busy <= '0' when (neurons_done = (others => '0')) else
74 --      '1' when (neurons_done = (others => '1'));
75 o_done <= neurons_done(0);
76
77 end architecture dense;
78
79
80 -----
81 -- END OF FILE: layer.vhd
82 -----
```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- synchronizer.vhd is a part of Innervator.
4  -- -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 -- Metastability occurs in flip-flops when the input signal changes
11 -- too close to the clock edge, violating "setup and hold" times;
12 -- this leaves the flip-flop in an unresolved state where the
13 -- output can be unpredictable and could potentially cause errors
14 -- in connected logic. To mitigate this, "cascading" multiple
15 -- flip-flops at the input is a common solution; it provides
16 -- additional time for the metastable signal to settle into a
17 -- stable 0 or 1 before being utilized elsewhere in the circuit.
18 --     This is called "synchronizing" or "de-glitching."
19 --
20 -- NOTE: A good portion of this synchronization also overlaps with
21 -- "pipelining," because both use a series of clocked flip-flops.
22 -- However, a good reason to separate them is to (hopefully) have
23 -- the synthesization tool lay out pipelines or synchronizers
24 -- close to their own respective groups. Also, the synchronizer
25 -- might be extended later on (maybe to support multiple clocks)
26 -- and separating them early-on would be beneficial, in that case.
27 --
28 -- TODO: Have a variadic variant for std_(u)logic VECTORS.
29 entity synchronizer is
30     generic (
31         g_NUM_STAGES : positive := 2
32     );
33     port (
34         i_clk      : in  std_ulogic;
35         i_signal    : in  std_ulogic;
36         o_signal    : out std_ulogic
37     );
38 end entity synchronizer;
39
40 architecture behavioral of synchronizer is
41     -- Synthesis tools will often replace a series of flip-flops
42     -- with better primitives, like shift-registers, that "achieve"
43     -- the same delaying effect. However, we might sometimes WANT
44     -- to use flip-flops specifically, so we can turn off that
45     -- optimization by using vendor-specific attribute definitions.
46     --
47     -- Also, note that avoiding the usage of explicit reset signals
48     -- may also result in the same shift-register (SRL) conversion.
49     --     ednasia.com/coding-consideration-for-pipeline-flip-flops
50
51     /* Xilinx Vivado/XST */
52     -- Disable the conversion of flip-flops to shift-registers
53     attribute shreg_extract : string;
54     -- Specifies that registers receive async. data.
55     attribute async_reg      : boolean; -- Also implies DONT_TOUCH.
56
57     -- TODO: See if we should apply this to the input/output?
58     -- Input
59     attribute shreg_extract of i_signal : signal is "no";
60     attribute async_reg     of i_signal : signal is true;
61     -- Output

```



```

62     attribute shreg_extract of o_signal : signal is "no";
63     attribute async_reg    of o_signal : signal is true;
64 begin
65
66     -- NOTE: Concurrent assignments mean that there is no
67     -- delay involved; both wires "connect" and act as one.
68     --
69     -- "Delaying" by 1 stage means that we will not have
70     -- to place any additional flip-flops in the middle.
71     -- The reason is that the unregistered i_signal
72     -- would also have a propagation delay of 1 whenever
73     -- something is assigned to it.
74     single_pipeline : if g_NUM_STAGES = 1 generate
75         o_signal <= i_signal; -- CONCURRENT assignment
76     -- Otherwise, implement actual delay with flip-flops.
77     else generate
78         -- NOTE: We subtract by 2 because the incoming (input)
79         -- signal itself also counts; so, when we talk about "double
80         -- registering" something, we know that there already WAS
81         -- a flip-flop register (i.e., the 'in' signal) and
82         -- we need to add just 1 more to it.
83         constant STAGES_HIGH : natural := g_NUM_STAGES - 2;
84
85         -- NOTE: I initially wanted to place this vector as individual
86         -- values inside the for-generate loop, but you cannot refer
87         -- back to a previous iteration of for-generate to access (i-1);
88         -- the solution was to declare it as a vector here. SEE:
89         -- https://groups.google.com/g/comp.lang.vhdl/c/rm97yoJwcWc
90         signal sync_regs : std_ulogic_vector
91             (0 to STAGES_HIGH) := (others => '0');
92
93         attribute shreg_extract of sync_regs : signal is "no";
94         attribute async_reg    of sync_regs : signal is true;
95     begin
96         cascade_chain : for i in 0 to STAGES_HIGH generate
97             cascade_register : process (i_clk) begin
98                 if rising_edge(i_clk) then
99                     -- If it is the first instance, then we must use
100                     -- the actual input signal; no previously cascaded
101                     -- registers can exist before i=0.
102                     sync_regs(i) <=
103                         i_signal when i = 0 else sync_regs(i-1);
104                 end if;
105             end process cascade_register;
106         end generate cascade_chain;
107
108         o_signal <= sync_regs(STAGES_HIGH); -- CONCURRENT assignment
109     end generate;
110
111 end architecture behavioral;
112
113
114 -----
115 -- END OF FILE: synchronizer.vhd
116 -----

```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- uart.vhd is a part of Innervator.
4  -- -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 -- A hierarchical interface for a full- or half-Duplex asynchronous
11 -- receiver/transmitter in the 8-N-1 frame: eight (8) data bits,
12 -- no (N) parity bit, and one (1) stop bit, plus an implicit start
13 -- bit; in this case, only 80% of the throughput is used for the data.
14 entity uart is
15     -- In digital communications, "baud" (i.e., symbol rate) is equal
16     -- to the bitrate (bit-rate). However, when the communications is
17     -- modulated to analog, a baud_can_encode more than 1 bit.
18     generic (
19         -- TODO: Take data length as a generic.
20         g_CLK_FREQ : positive := 100e6;
21         g_BAUD      : positive range positive'low to g_CLK_FREQ := 9_600
22     );
23     -- NOTE: 'Buffer' data flows out of the entity, but the entity can
24     -- read the signal (allowing for internal feedback); however, the
25     -- signal cannot be driven from outside the entity, unlike inputs.
26     port (
27         -- UART-Rx/Tx (Receive/Transmit) Shared Ports
28         i_clk      : in  std_ulogic; -- Internal FPGA clock
29         i_rst      : in  std_ulogic; -- Reset
30         -- UART-Rx (Receive) Ports
31         i_rx_serial : in  std_logic; -- External connection (wire)
32         o_rx_done   : out std_ulogic; -- "Done Reading" signal
33         o_rx_byte   : out std_ulogic_vector (7 downto 0); -- LSB first
34         -- UART-Tx (Transmit) Ports
35         i_tx_send   : in  std_ulogic; -- "Start Sending" signal
36         i_tx_byte   : in  std_ulogic_vector (7 downto 0); -- LSB first
37         o_tx_active : out std_ulogic; -- Half-Duplex transmitters ONLY
38         o_tx_done   : out std_ulogic; -- "Done Transmitting" signal
39         o_tx_serial : out std_logic -- External connection (wire)
40     );
41
42     -- NOTE: constants here are applied to ALL architectures
43     constant TICKS_PER_BIT : positive :=
44         positive(g_CLK_FREQ / g_BAUD) - 1;
45     constant DATA_HIGH    : natural  := o_rx_byte'high;
46 begin
47 end entity uart;
48
49
50 library ieee;
51     use ieee.std_logic_1164.all;
52     use ieee.numeric_std.all;
53
54 -- UART's standalone sub-components (i.e., UART-RCVR and -XMTR)
55 package uart_pkg is
56     component uart_rcvr
57         generic (
58             g_CLK_FREQ : positive;
59             g_BAUD      : positive
60         );
61         port (

```

```
62         i_clk      : in  std_ulogic;
63         i_rst      : in  std_ulogic;
64
65         i_serial    : in  std_logic;
66         o_done      : out std_ulogic;
67         o_byte      : out std_ulogic_vector (7 downto 0)
68     );
69 end component uart_rcvr;
70
71 component uart_xmtr
72     generic (
73         g_CLK_FREQ : positive;
74         g_BAUD      : positive
75     );
76     port (
77         i_clk      : in  std_ulogic;
78         i_rst      : in  std_ulogic;
79
80         i_send      : in  std_ulogic;
81         i_byte      : in  std_ulogic_vector (7 downto 0);
82         o_active     : out std_ulogic;
83         o_done       : out std_ulogic;
84         o_serial     : out std_logic
85     );
86 end component uart_xmtr;
87 end package uart_pkg;
88
89
90 -----
91 -- END OF FILE: uart.vhd
92 -----
```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- uart_xcvr.vhd is a part of Innervator.
4  -- -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 library work;
11     use work.uart_pkg.all;
12
13 library config;
14     use config.constants.all;
15
16 -- A half- or full-duplex async. receiver/transmitter (in 8-N-1 frame)
17 architecture transceiver of uart is -- [Structural arch.]
18 begin
19
20     receiver_component : component work.uart_pkg.uart_rcvr
21         generic map (
22             g_CLK_FREQ => g_CLK_FREQ,
23             g_BAUD      => g_BAUD
24         )
25         port map (
26             i_clk      => i_clk,
27             i_rst      => i_rst,
28
29             i_serial    => i_rx_serial,
30             o_done      => o_rx_done,
31             o_byte      => o_rx_byte
32         );
33
34     transmitter_component : component work.uart_pkg.uart_xmtr
35         generic map (
36             g_CLK_FREQ => g_CLK_FREQ,
37             g_BAUD      => g_BAUD
38         )
39         port map (
40             i_clk      => i_clk,
41             i_rst      => i_rst,
42
43             i_send      => i_tx_send,
44             i_byte      => i_tx_byte,
45             o_active    => o_tx_active,
46             o_done      => o_tx_done,
47             o_serial    => o_tx_serial
48         );
49
50     -- [Place for other concurrent statements]
51
52 end architecture transceiver;
53
54
55 library work;
56     use work.all;
57
58 -- NOTE: 'configuration' in VHDL is a barely documented and arcane
59 -- keyword, and the excessive repetition of component/entity ports
60 -- might also be defeating its purpose.
61 --     From what I have gathered, components are "idealized"

```

```

62 -- _placeholders_ for future "realized" entities. In an electronics
63 -- sense, they are chip sockets for upcoming chips; most of the time
64 -- they are not useful and merely add an unneeded layer of abstraction.
65 -- However, they could sometimes be useful for giving a hierarchical
66 -- organization to "sub-entities."
67 --     Additionally, components should be declared in an architecture's
68 -- header and correspond exactly (name- & port-wise) to their entities.
69 -- They will also require separate instantiations in said architecture's
70 -- body (i.e., structural VHDL), often resulting in lots of duplicated
71 -- port assignments.
72 --     It is also possible to use configurations to "bind" a specific
73 -- instance (or even 'all' instances) of a component to a specific
74 -- entity-architecture pair or other sub-configurations. Afterward,
75 -- you may even instantiate the configuration itself as you would do
76 -- so with an entity or component.
77 --
78 -- NOTE: Here, it is important to note that 'for' does NOT refer to
79 -- a for-loop and means a literal 'for' (i.e., FOR X, use Y).
80 --
81 -- NOTE: You cannot leave ports of type 'input' as 'open' but you can
82 -- assign 'U', 'X', 'Z', or '-' to them to achieve the same effect;
83 -- among these, '-' makes the most sense, semantically, but 'Z' seems
84 -- to be the one replicating the 'open' effect in schematics.
85 configuration uart_xcvr of uart is -- [config. name] of entity
86     for transceiver -- (i.e., the encapsulating architecture)
87
88         for all : work.uart_pkg.uart_rcvr -- (i.e., component instance)
89             use entity work.uart (receiver)
90                 generic map (
91                     g_CLK_FREQ => g_CLK_FREQ,
92                     g_BAUD     => g_BAUD
93                 )
94                 port map (
95                     i_clk      => i_clk,
96                     i_rst      => i_rst,
97
98                     i_rx_serial => i_serial,
99                     o_rx_done  => o_done,
100                    o_rx_byte  => o_byte,
101
102                    i_tx_send   => 'Z',
103                    i_tx_byte   => (others => 'Z'),
104                    o_tx_active => open,
105                    o_tx_done   => open,
106                    o_tx_serial => open
107                );
108         end for;
109
110         for all : work.uart_pkg.uart_xmtr -- (i.e., component instance)
111             use entity work.uart (transmitter)
112                 generic map (
113                     g_CLK_FREQ => g_CLK_FREQ,
114                     g_BAUD     => g_BAUD
115                 )
116                 port map (
117                     i_clk      => i_clk,
118                     i_rst      => i_rst,
119
120                     i_rx_serial => 'Z',
121                     o_rx_done  => open,
122                     o_rx_byte  => open,
123

```

```
124         i_tx_send    => i_send,  
125         i_tx_byte    => i_byte,  
126         o_tx_active  => o_active,  
127         o_tx_done    => o_done,  
128         o_tx_serial  => o_serial  
129     );  
130 end for;  
131  
132 end for;  
133 end configuration uart_xcvr;  
134  
135  
136 -- -----  
137 -- END OF FILE: uart_xcvr.vhd  
138 -- -----
```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- network.vhd is a part of Innervator.
4  -- -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 library work;
11     context work.neural_context;
12
13 library config;
14     use config.constants.all;
15
16 entity network is
17     generic (
18         g_NETWORK_PARAMS    : network_layers;
19         g_BATCH_SIZE        : positive;
20         g_PIPELINE_STAGES   : natural
21     );
22     port (
23         i_inputs  : in  neural_bvector
24             (0 to g_NETWORK_PARAMS(g_NETWORK_PARAMS'low).dims.cols-1);
25         o_outputs : out neural_bvector
26             (0 to g_NETWORK_PARAMS(g_NETWORK_PARAMS'high).dims.rows-1);
27         /* Sequential (pipeline) controllers */
28         i_clk     : in  std_ulogic; -- Clock
29         i_rst     : in  std_ulogic; -- Reset
30         i_fire    : in  std_ulogic; -- Start/fire up all the layers
31         o_done    : out std_ulogic  -- Is the network done processing?
32     );
33
34     constant NUM_LAYERS    : positive := g_NETWORK_PARAMS'length;
35     -- Number of neurons in the first (i.e., input) layer.
36     constant NUM_INPUTS   : positive := i_inputs'length;
37     -- Number of neurons in the last (i.e., output) layer.
38     constant NUM_OUTPUTS  : positive := o_outputs'length;
39 end entity network;
40
41
42 architecture neural of network is
43
44     -- Helper functions to "deflate" (or sanitize) max-sized arrays
45     -- introduced due to the workaround around VHDL's lack of
46     -- variable-sized elements in arrays.  You can find more details
47     -- on the specifics of this workaround (and how to reverse it)
48     -- in the file_parser.vhd file.
49     function deflate( -- Arrays
50         inflated_data : neural_vector;
51         size_key      : natural
52     ) return neural_vector is
53         constant deflated_array : neural_vector (0 to size_key-1) :=
54             inflated_data (0 to size_key-1);
55     begin
56         return deflated_array;
57     end function deflate;
58
59     function deflate( -- Matrices (i.e., Nested Arrays of Arrays)
60         inflated_data : neural_matrix;
61         size_key      : dimension

```

```

62 ) return neural_matrix is
63     variable deflated_matrix : neural_matrix
64         (0 to size_key.rows-1) (0 to size_key.cols-1);
65 begin
66     deflate_rows : for i in deflated_matrix'range loop
67         deflated_matrix(i) := -- Deflate individual sub-arrays.
68             deflate(inflated_data(i), size_key.cols);
69     end loop deflate_rows;
70
71     return deflated_matrix;
72 end function deflate;
73
74 -- Unfortunately, VHDL has a language-level limitation where it
75 -- does not allow you to refer back to a previous instance of a
76 -- for-generate's local signals, even though you could manually
77 -- "unroll" it into a single 'block' clause, containing mangled
78 -- names of signals that cannot collision. A solution is to
79 -- employ the same max-dimension workaround from file_parser.vhd.
80 -- Fortunately, the max-dimension has already been calculated by
81 -- the file parser, earlier; we get to re-use it here.
82 --     Unused/dummy elements would get discarded and optimized
83 -- by the synthesizer, but this can still be very redundant
84 -- and clunky, especially for nested arrays and very large ones.
85 signal layers_done      : std_ulogic_vector (0 to NUM_LAYERS-1);
86 -- This is a NESTED array of 'neural_bvector' (an array type).
87 --
88 -- NOTE: We cannot use the 'element attribute here; because
89 -- each element of the parameter array is a record on its own,
90 -- toolchains will break apart on such complicated expressions.
91 signal layers_outputs : neural_bmatrix (0 to NUM_LAYERS-1)
92     (0 to g_NETWORK_PARAMS(0).weights'length-1);
93
94 -- The entire network is done whenever the last layer of it is.
95 signal network_done      : std_ulogic;
96 -- The output of the entire network is its layer layer's output.
97 -- This could be arg-max'ed or used as-is, as needed.
98 signal network_outputs : neural_bvector (0 to NUM_OUTPUTS-1);
99 begin
100
101 -- NOTE: This is a "feed-forward" neural network, meaning that
102 -- each layer's output is connected to the proceeding layer's
103 -- input, in a chain-like formation.
104 neural_network : for i in 0 to NUM_LAYERS-1 generate
105     -- NOTE: These are intermediary signals to get around a VHDL
106     -- limitation where you cannot have "conditionally mapped"
107     -- port or generic maps (even in static for-generate clauses);
108     -- we first assign the condition to a signal and then the port
109     --     Because these are considered "concurrent" signal
110     -- connections, there's no assignment or clock cycle delay.
111     --     Also, in < VHDL-2019, you cannot define these as
112     -- constants, because you cannot use when...else in them.
113     -- TODO: See if we can somehow get the UNCONSTRAINED subtype
114     -- of 'inputs' here, rather than hard-coding it.
115     signal inputs_im : neural_bvector
116         (0 to g_NETWORK_PARAMS(i).dims.cols-1);
117     signal i_fire_im : i_fire'ssubtype;
118
119     -- NOTE: in VHDL, you cannot constrain port assignments
120     -- directly, because these would count as "locally non-static
121     -- ranges," even though they are constants at compile-time;
122     -- we have to constrain each generated instance's input here.
123     constant sanitized_weights : neural_matrix :=

```



```

124         deflate(
125             g_NETWORK_PARAMS(i).weights,
126             g_NETWORK_PARAMS(i).dims
127         );
128     constant sanitized_biases : neural_vector :=
129         deflate(
130             g_NETWORK_PARAMS(i).biases,
131             g_NETWORK_PARAMS(i).dims.rows
132         );
133 begin
134     -- NOTE: Somehow, using when...else instead of if..generate
135     -- seems to result in indices such as -1; it seems that the
136     -- condition after the 'else' part is "evaluated" even if
137     -- it is not supposed to be [i.e., when i=0, it tries to do
138     -- i-1 and index the array as (-1)].
139     input_layer_condition : if i = 0 generate
140         inputs_im <= i_inputs;
141         i_fire_im <= i_fire;
142     else generate
143         inputs_im <= layers_outputs(i-1)
144             (0 to g_NETWORK_PARAMS(i).dims.cols-1);
145         i_fire_im <= layers_done(i-1);
146     end generate input_layer_condition;
147
148     neural_layer : entity work.layer (dense)
149         generic map (
150             -- NOTE: These arrays are "sliced" due to a workaround,
151             -- which was used to bypass the lack of variable-sized
152             -- arrays in VHDL; said workaround (explained in the
153             -- file_parser.vhd file) would declare an array to be
154             -- the "maximum" possible size (the biggest out of all
155             -- its elements) and keep track of their "true" sizes
156             -- in a separate field called .dims. Here, we have
157             -- simply sliced the "inflated" array, discarding the
158             -- unused/dummy elements based on the true sizes.
159             g_LAYER_WEIGHTS => sanitized_weights,
160             g_LAYER_BIASES  => sanitized_biases,
161             g_BATCH_SIZE    => g_BATCH_SIZE,
162             g_PIPELINE_STAGES => g_PIPELINE_STAGES
163         )
164         port map (
165             i_inputs => inputs_im,
166             o_outputs => layers_outputs(i)
167                 (0 to g_NETWORK_PARAMS(i).dims.rows-1),
168             i_clk    => i_clk, -- Clock
169             i_rst    => i_rst, -- Reset
170             -- The first layer (i.e., the "input layer") will
171             -- activate whenever the network is told to.
172             -- Each subsequent layer will "fire" (i.e., activate)
173             -- whenever its previous layer is "done" processing.
174             i_fire   => i_fire_im,
175             o_done   => layers_done(i) -- Is it done processing?
176         );
177
178     end generate neural_network;
179
180     network_done    <= layers_done(layers_done'high);
181     network_outputs <= layers_outputs(layers_outputs'high)
182         (0 to NUM_OUTPUTS-1);
183
184     -- The Network's 'done' signal is different than the that of the
185     -- neurons in the sense that it "stays" done; that is to let us

```

# Innervator

```

186  -- know that that it has finished its processing and remains IDLE,
187  -- ready to accept another set of data.  Neurons' done signals
188  -- lasted for a single clock cycle and then switched back to 0;
189  -- because each layer's done signal was connected to the following
190  -- layers' 'fire' signal, continuing to keep said layer's done
191  -- at 1 would result in proceeding layers never ending firing.
192
193  toggle_out : process (i_clk, i_rst) is
194      procedure perform_reset is
195      begin
196          o_done <= '1';
197          o_outputs <= (others => (others => '0'));
198      end procedure perform_reset;
199  begin
200
201      if not c_RST_SYNC and i_rst = c_RST_POLE then perform_reset;
202      elsif rising_edge(i_clk) then
203          if c_RST_SYNC and i_rst = c_RST_POLE then perform_reset;
204          else
205
206              if network_done = '1' then
207                  o_done <= '1';
208                  o_outputs <= network_outputs;
209                  -- "elsif" so 'network_done' lasts at least one cycle.
210              elsif i_fire = '1' then
211                  o_done <= '0';
212                  o_outputs <= (others => (others => '0'));
213              end if;
214
215          end if;
216      end if;
217  end process toggle_out;
218
219  /*
220  toggle_out : process (all) is
221  begin
222
223      if network_done = '1' then
224          o_done <= '1';
225          outputs <= network_outputs;
226          -- "elsif" so 'network_done' lasts at least one cycle.
227      elsif i_fire = '1' then
228          o_done <= '0';
229          outputs <= (others => (others => '0'));
230      end if;
231
232  end process toggle_out;
233  */
234  end architecture neural;
235
236
237  -- -----
238  -- END OF FILE: network.vhd
239  -- -----

```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- math.vhd is a part of Innervator.
4  -- -----
5
6  library ieee;
7      use ieee.std_logic_1164.all;
8
9  library neural;
10     context neural.neural_context;
11
12
13  -- TODO: Implement more activation functions (e.g., ReLU, etc.)
14  package math is
15
16      -- Returns the index of the highest number in an unsigned array.
17      -- TODO: What if two numbers are equal?
18      -- TODO: Maybe provide a pipelined/multi-cycle variant.
19      function arg_max(
20          values : neural_bvector
21      ) return natural;
22
23  end package math;
24
25
26  package body math is
27
28      function arg_max(
29          values : neural_bvector
30      ) return natural is
31          variable max_index    : natural;
32          variable current_max : neural_bvector'element;
33      begin
34
35          current_max := (others => '0');
36
37          for i in values'range loop
38              if values(i) > current_max then
39                  max_index := i;
40                  current_max := values(i);
41              end if;
42          end loop;
43
44          return max_index;
45      end function arg_max;
46
47  end package body math;
48
49  -- -----
50  -- END OF FILE: math.vhd
51  -- -----

```

```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- uart_rcvr.vhd is a part of Innervator.
4  -- -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 library config;
11     use config.constants.all;
12
13 -- A simplex async. receiver (in 8-N-1 frame)
14 architecture receiver of uart is -- [RTL arch.]
15     -- Receiver's Finite-State Machine (FSM)
16     type uart_rx_state_t is (
17         idle, started, reading, done
18     );
19     signal uart_rx_state : uart_rx_state_t := idle;
20
21     -- Synchronized signal (from a metastable input signal).
22     signal synced_serial : std_ulogic;
23
24     signal tick_cnt    : natural range 0 to TICKS_PER_BIT := 0;
25     signal bit_index   : natural range 0 to DATA_HIGH := 0;
26 begin
27
28     -- NOTE: The input signal is assumed to have been synchronized/
29     -- deglitched beforehand, at the top module.
30     synced_serial <= i_rx_serial; -- CONCURRENT assignment
31
32
33     -- NOTE: As a little history, the reason why the 'start bit' is
34     -- checked to be 'active low' as opposed to 'active high' is
35     -- because physical cable connections could run far and were also
36     -- susceptible to damage along the path; by constantly driving an
37     -- 'active high' signal you could know for sure that (long) pauses
38     -- meant a disruption in the line, unlike what would otherwise be
39     -- interpreted as intentional "silence" in an 'active low' setup.
40     receive : process (i_clk, i_rst) is
41         procedure perform_reset is
42         begin
43             uart_rx_state <= idle;
44         end procedure perform_reset;
45     begin
46
47         if not c_RST_SYNC and i_rst = c_RST_POLE then perform_reset;
48         elsif rising_edge(i_clk) then
49             if c_RST_SYNC and i_rst = c_RST_POLE then perform_reset;
50             else
51                 -- NOTE: There is no need for an 'others' default case,
52                 -- because we are dealing with enumerated types and not
53                 -- std_logics; all enumerated cases are accounted for.
54                 -- SEE: sigasi.com/tech/
55                 --         vhdl-case-statements-can-do-without-others
56                 -- Despite this, the 'others' case can still be used to
57                 -- harden the state machine (against radiation, maybe)
58                 -- and make it safer by recovering from unknown values.
59                 --
60                 -- NOTE: Always assign a value to signals in a state
61                 -- machine; otherwise, latches (hard to synth, slow, &

```

```

62 -- prone to metastability) may be inferred. Also, since
63 -- we cannot really use default values to solve that
64 -- issue here, we instead re-assign the same state in
65 -- branches that do not result in a change of state.
66 case uart_rx_state is
67     -- Due to UART being asynchronous, a "start bit" is
68     -- utilized by the external transmitter to signal
69     -- that the actual data bits are forthcoming; the
70     -- start bit is simply a falling edge (from idle
71     -- high to a low pulse) immediately followed by
72     -- user data bits.
73     when idle =>
74         uart_rx_state <= idle;
75
76         -- Reset back to default values
77         o_rx_done <= '0';
78         o_rx_byte <= (others => '0');
79         bit_index <= 0;
80         tick_cnt <= 0;
81
82         if synced_serial = '0' then
83             uart_rx_state <= started;
84         end if;
85     -- -----
86
87     -- It is important to note that, even after the
88     -- start bit is detected, we still re-check the
89     -- start bit at its middle (one-half 'bit time')
90     -- as to make sure it was really valid. If not,
91     -- it is considered a spurious pulse (noise)
92     -- and is ignored.
93     -- Also, by waiting for the mid-point and
94     -- only then moving to the next state, we'd
95     -- also be shifting future data reads (sampling)
96     -- to their respective mid-points.
97     when started =>
98         uart_rx_state <= started;
99
100         if (tick_cnt < TICKS_PER_BIT / 2) then
101             tick_cnt <= tick_cnt + 1; -- Not middle
102         else -- Reached the middle
103             if (syncd_serial = '0') then -- Real start
104                 tick_cnt <= 0; -- Found mid; reset cnt
105                 uart_rx_state <= reading; -- Begin read
106             else -- Spurious pulse; ignore
107                 uart_rx_state <= idle;
108             end if;
109         end if;
110     -- -----
111
112     when reading =>
113         uart_rx_state <= reading;
114
115         if (tick_cnt < TICKS_PER_BIT) then
116             tick_cnt <= tick_cnt + 1;
117         else
118             tick_cnt <= 0;
119             -- NOTE: Data is sent with Least-Sig. Byte
120             -- first; holding type should be 'downto'
121             o_rx_byte(bit_index) <= synced_serial;
122
123             -- Check if more bits remain to be received

```

# Innervator

```

124         if (bit_index < DATA_HIGH) then
125             bit_index <= bit_index + 1;
126         else
127             bit_index <= 0;
128             uart_rx_state <= done;
129         end if;
130     end if;
131     -----
132
133     -- Once the data bits finish, a "stop bit" will
134     -- indicate so; the stop bit is normally a
135     -- transition back to the idle state or remaining
136     -- at the high state for an extra bit time.
137     -- A second (optional) stop bit can be configured,
138     -- usually to give the receiver time to get ready
139     -- for the next frame, but that is uncommon.
140     --     Considering how we already know the frame's
141     --     size beforehand, it might seem unnecessary to
142     --     even have a stop bit. However, the UART
143     --     does not depend on the start bit itself, for
144     --     synchronization, but the falling_edge_between
145     --     the previous stop bit AND the start bit.
146     --     There might not be such an edge without both
147     --     the start and stop bits.
148     when done =>
149         uart_rx_state <= done;
150
151         if (tick_cnt < TICKS_PER_BIT) then
152             tick_cnt <= tick_cnt + 1;
153         else
154             -- TODO: Error-check the stop bit as an
155             -- extra protection against de-syncing.
156             --if (synced_serial /= '1') then
157             --    perform_reset;
158             --end if;
159
160             o_rx_done <= '1';
161             tick_cnt <= 0;
162             uart_rx_state <= idle;
163         end if;
164     -----
165
166     -- Hardening in case of "unknown" states
167     when others => -- 1-clock-long cleanup phase
168         perform_reset;
169     -----
170 end case;
171
172     end if;
173 end if;
174 end process receive;
175
176 end architecture receiver;
177
178
179 -----
180 -- END OF FILE: uart_rcvr.vhd
181 -----

```

```
1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  -- uart_xmtr.vhd is a part of Innervator.
4  -- -----
5
6
7  library ieee;
8      use ieee.std_logic_1164.all;
9
10 library config;
11     use config.constants.all;
12
13 -- TODO (when needed)
14 -- A simplex async. transmitter (in 8-N-1 frame)
15 architecture transmitter of uart is -- [RTL arch.]
16     -- Placeholder
17 begin
18     -- Placeholder
19 end architecture transmitter;
20
21
22 -- -----
23 -- END OF FILE: uart_xmtr.vhd
24 -- -----
```

```
1  -- -----
2  --  SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  --  prescaler.vhd is a part of Innervator.
4  -- -----
5
6
7
8
9
10
11
12 -- -----
13 --  END OF FILE: prescaler.vhd
14 -- -----
```



```

1  -- -----
2  -- SPDX-License-Identifier: LGPL-3.0-or-later or CERN-OHL-W-2.0
3  --
4  -- Innervator: Hardware Acceleration for Neural Networks
5  --
6  -- Copyright (C) 2024 Fereydoun Memarzanjany
7  --
8  -- This hardware-descriptive model is free hardware design dual-
9  -- licensed under the GNU LGPL or CERN OHL v2 Weakly Reciprocal: you
10 -- can redistribute it and/or modify it under the terms of the...
11 --     * GNU Lesser General Public License as published by
12 --       the Free Software Foundation, either version 3 of the License,
13 --       or (at your option) any later version; OR
14 --     * CERN Open Hardware Licence Version 2 - Weakly Reciprocal.
15 --
16 -- This is distributed in the hope that it will be useful,
17 -- but WITHOUT ANY WARRANTY; without even the implied warranty of
18 -- MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
19 -- GNU Lesser General Public License for more details.
20 --
21 -- You should have received a copy of the GNU Lesser General
22 -- Public License and the CERN Open Hardware Licence Version
23 -- 2 - Weakly Reciprocal along with this. If not, see...
24 --     * <https://spdx.org/licenses/LGPL-3.0-or-later.html>; and
25 --     * <https://spdx.org/licenses/CERN-OHL-W-2.0.html>.
26 --
27 -- -----
28
29
30 -- NOTE: All of the source files herein conform to the RFC 678
31 -- plaintext document standard, as well as the Ada 95 Quality and
32 -- Style Guide 2.1.9 (Source Code Line Length); it is good readability
33 -- practice to limit a line of code's columns to 72 characters (which
34 -- include only the printable characters, not line endings or cursors).
35 --     I specifically chose 72 (and not some other limit like 80/132)
36 -- to ensure maximal compatibility with older technology, terminals,
37 -- paper hardcopies, and e-mails. While some other guidelines permit
38 -- more than just 72 characters, it is still important to note that
39 -- American teletypewriters could sometimes write upto only 72, and
40 -- older code (e.g., FORTRAN, Ada, COBOL, Assembler, etc.) used to
41 -- be hand-written on a "code form" in corporations like IBM; said
42 -- code form typically reserved the first 72 columns for statements,
43 -- 8 for serial numbers, and the remainder for comments, which was
44 -- finally turned into a physical punch card with 80 columns.
45 --     Even in modern times, the 72 limit can still be beneficial:
46 -- you can easily quote a 72-character line over e-mail without
47 -- requiring word-wrapping or horizontal scrolling.
48 --     As a sidenote, the reason that some guidelines, like PEP 8
49 -- (Style Guide for Python Code), recommended 79 characters (i.e.,
50 -- not 80) was that the 80th character in a 80x24 terminal might
51 -- have been a bit hard to read.
52
53 -- Thanks to yet another arcane bug within Vivado 2024, in which it
54 -- completely breaks apart when you try to access attributes of this
55 -- constant within the same declaratory region (even though it is
56 -- perfectly valid VHDL and ModelSim also has no problems with it),
57 -- we have no choice but to declare it in a "separate" area:
58 library work;
59     use work.constants.all;
60 library neural;
61     context neural.neural_context;

```

```

62     use      neural.file_parser.all;
63 package attribute_bugfix is
64     constant debug_NETWORK_OBJECT : network_layers :=
65         parse_network_from_dir(c_DAT_PATH);
66 end package attribute_bugfix;
67
68
69 library ieee;
70     use ieee.std_logic_1164.all;
71     use ieee.numeric_std.all;
72
73 library work;
74     use work.constants.all;
75
76 library core;
77
78 library neural;
79     context neural.neural_context;
80     use      neural.file_parser.all;
81
82
83 entity neural_processor is
84     port (
85         i_clk   : in std_ulogic;
86         i_rst   : in std_ulogic;
87         i_uart  : in std_logic;
88         o_uart  : out std_logic;
89         o_led   : out std_logic_vector (3 downto 0)
90     );
91 begin
92 end entity neural_processor;
93
94 architecture structural of neural_processor is
95     alias NETWORK_OBJECT is work.attribute_bugfix.debug_NETWORK_OBJECT;
96
97     -- Number of layers in network (excluding the input data themself)
98     constant NUM_LAYERS : positive :=
99         NETWORK_OBJECT'length;
100    -- Number of neurons in the first (i.e., input) layer
101    constant NUM_INPUTS : positive :=
102        NETWORK_OBJECT(NETWORK_OBJECT'low).dims.cols;
103    -- Number of neurons in the last (i.e., output) layer
104    constant NUM_OUTPUTS : positive :=
105        NETWORK_OBJECT(NETWORK_OBJECT'high).dims.rows;
106
107
108    -- Synchronized/deglitched ports
109    signal i_rst_synced : std_ulogic;
110    signal i_uart_synced : std_ulogic;
111    -- In case the reset button requires inverting
112    signal i_rst_corrected : std_ulogic;
113    -- Synchronized buttons
114    signal i_rst_debounced : std_ulogic;
115
116
117    /* UART signals */
118    signal byte_read_done : std_ulogic := '0';
119    signal byte_read_value : std_logic_vector (7 downto 0) :=
120        (others => '0');
121
122    -- The input array that will be received via UART from a computer.
123    signal input_data : neural_bvector (0 to NUM_INPUTS-1) :=

```

```

124     (others => (others => '0'));
125     signal input_data_count      : natural range 0 to NUM_INPUTS := 0;
126     signal input_data_received : std_ulogic := '0';
127
128     -- UART Transmitter signals
129     signal result_ready : std_ulogic;
130     signal result_byte  : std_logic_vector (7 downto 0);
131
132     signal network_done      : std_ulogic := '0';
133     signal network_outputs   : neural_bvector (0 to NUM_OUTPUTS-1);
134     -- TODO: Have an actual function to binary-encode decimals.
135     signal network_prediction : unsigned (3 downto 0); -- Arg-Max'ed
136 begin
137     -- TODO: Print network metadata here (though Vivado has assert bugs)
138     --assert false
139     --     report natural'image(NETWORK_OBJECT'element.weights'length)
140     --         severity failure;
141
142     /*
143     Port Setup
144     */
145
146     -- Synchronize/deglitch the incoming data, allowing it to be used
147     -- in our own clock domain, avoiding metastability problems.
148     synchronize_input_ports : block
149     begin
150         sync_reset      : entity core.synchronizer
151             generic map (g_NUM_STAGES => c_SYNC_NUM)
152             port map (
153                 i_clk      => i_clk,
154                 i_signal    => i_rst,
155                 o_signal    => i_rst_synced
156             );
157         sync_uart_in : entity core.synchronizer
158             generic map (g_NUM_STAGES => c_SYNC_NUM)
159             port map (
160                 i_clk      => i_clk,
161                 i_signal    => i_uart,
162                 o_signal    => i_uart_synced
163             );
164     end block synchronize_input_ports;
165
166     -- Invert the reset button; even if the FPGA board has a negative
167     -- reset, the FPGA might work "better" with positive resets,
168     -- internally. (more info in config.vhd)
169     -- TODO: Investigate if the reset will have skew, in this case
170     invert_reset : if c_RST_INVT generate
171         i_rst_corrected <= not i_rst_synced;
172     else generate -- Else, don't invert
173         i_rst_corrected <= i_rst_synced;
174     end generate invert_reset;
175
176     -- Remove the bouncing "noise" from input buttons
177     debounce_buttons : block
178     begin
179         debounce_reset : entity core.debounce
180             generic map (g_TIMEOUT_MS => c_DBNC_LIM)
181             port map (
182                 i_clk      => i_clk,
183                 i_button    => i_rst_corrected,
184                 o_button    => i_rst_debounced
185             );

```

```

186     end block debounce_buttons;
187
188
189
190  /*
191      Part Instantiations
192  */
193
194  uart_transceiver : configuration core.uart_xcvr
195      generic map (
196          g_CLK_FREQ => c_CLK_FREQ,
197          g_BAUD      => c_BIT_RATE
198      )
199      port map (
200          i_clk      => i_clk,
201          i_rst      => i_rst_debounced,
202          -- The Receiver Component
203          i_rx_serial => i_uart_synced,
204          o_rx_done   => byte_read_done,
205          o_rx_byte   => byte_read_value,
206          -- TODO: Implement the UART Transmitter, too.
207          i_tx_send   => 'Z',
208          i_tx_byte   => (others => 'Z'),
209          o_tx_active => open, -- Unused
210          o_tx_done   => open, -- Unused
211          o_tx_serial => open
212      );
213
214
215  -- TODO: Eventually, have the option to select between using LUTRAM
216  -- (which is, confusingly, also known as DistRAM or DRAM) and BRAM,
217  -- the dedicated---but single/dual channel---block ram on FPGAs.
218  -- TODO: Also, maybe have this in a separate entity?
219  receive_data : process (i_clk)
220      procedure perform_reset is
221      begin
222          input_data          <= (others => (others => '0'));
223          input_data_count    <= 0;
224          input_data_received <= '0';
225      end procedure perform_reset;
226  begin
227      if not c_RST_SYNC and i_rst_debounced = c_RST_POLE
228          then perform_reset;
229      elsif rising_edge(i_clk) then
230          if c_RST_SYNC and i_rst_debounced = c_RST_POLE
231              then perform_reset;
232          else
233              input_data_received <= '0';
234
235              data_remain : if (input_data_count < NUM_INPUTS) then
236                  byte_read : if (byte_read_done = '1') then
237
238                      input_data(input_data_count) <=
239                          to_ufixed(byte_read_value,
240                              input_data'element'high,
241                              input_data'element'low
242                      );
243
244                      input_data_count <= input_data_count + 1;
245                  end if byte_read;
246              else -- All of the data got received
247                  input_data_received <= '1';

```

```

248         input_data_count    <= 0;
249     end if data_remain;
250
251     end if;
252 end if;
253 end process receive_data;
254
255
256 neural_engine : entity neural.network
257     generic map (
258         g_NETWORK_PARAMS    => NETWORK_OBJECT,
259         g_BATCH_SIZE        => c_BATCH_SIZE,
260         g_PIPELINE_STAGES    => c_PIPE_STAGE
261     )
262     port map (
263         i_inputs    => input_data,
264         o_outputs    => network_outputs,
265         i_clk        => i_clk,
266         i_rst        => i_rst_debounced,
267         i_fire        => input_data_received,
268         o_done        => network_done
269     );
270
271
272     -- NOTE: Unlike popular beliefs, subprograms (like functions)
273     -- CAN be concurrently called outside of processes.
274     -- TODO: Pipeline and clock this later.
275     network_prediction <= to_unsigned(
276         neural.math.arg_max(network_outputs),
277         network_prediction'length
278     );
279
280     -- TODO: Transfer the data back using the UART as opposed to LEDs
281     -- and also check network_done first.
282     o_led(3) <= network_prediction(3);
283     o_led(2) <= network_prediction(2);
284     o_led(1) <= network_prediction(1);
285     o_led(0) <= network_prediction(0);
286     /*
287     transmit_result : process (i_clk)
288         procedure perform_reset is
289         begin
290
291             end procedure perform_reset;
292     begin
293         if not c_RST_SYNC and i_rst_debounced = c_RST_POLE
294             then perform_reset;
295         elsif rising_edge(i_clk) then
296             if c_RST_SYNC and i_rst_debounced = c_RST_POLE
297                 then perform_reset;
298             else
299
300                 end if;
301             end if;
302         end process transmit_result;
303     */
304
305 end architecture structural;
306
307
308 -----
309 -- END OF FILE: main.vhd

```

