

PyGeNN: A Python library for GPU-enhanced neural networks

James C Knight^{1,*}, Anton Komissarov, Thomas Nowotny¹

¹Centre for Computational Neuroscience and Robotics, School of Engineering and Informatics, University of Sussex, Brighton, United Kingdom

Correspondence*:

James C Knight

J.C.Knight@sussex.ac.uk

2 ABSTRACT

3 For full guidelines regarding your manuscript please refer to Author Guidelines.

4 As a primary goal, the abstract should render the general significance and conceptual advance
5 of the work clearly accessible to a broad readership. References should not be cited in the
6 abstract. Leave the Abstract empty if your article does not require one, please see Summary
7 Table for details according to article type.

8 **Keywords:** GPU, high-performance computing, parallel computing, benchmarking, computational neuroscience, spiking neural
9 networks, Python

1 INTRODUCTION

10 A wide range of spiking neural network (SNN) simulators are available, each with their own niche.
11 NEST (Gewaltig and Diesmann, 2007) is widely used for large-scale point neuron simulations on
12 distributed computing systems; NEURON (Carnevale and Hines, 2006) and Arbor (Akar et al., 2019)
13 specialise in the simulation of complex multi-compartmental models; and CARLsim (Chou et al., 2018),
14 NeuronGPU (Golosio et al., 2020) and GeNN (Yavuz et al., 2016) use Graphics Processing Units (GPUs)
15 to accelerate point neuron models. For performance reasons, many of these simulators are written in C++
16 and, especially amongst the older simulators, users describe their models either using a Domain-Specific
17 Language (DSL) or directly in C++. For programming language purists, a DSL may be an elegant way of
18 describing an SNN network model and, for simulator developers, not having to add bindings to another
19 language is convenient. However, both choices act as a barrier to potential users. Therefore, with both the
20 computational neuroscience and machine learning communities gradually coalescing towards a Python-
21 based ecosystem with a wealth of mature libraries for scientific computing (Hunter, 2007; Van Der Walt
22 et al., 2011; Millman and Aivazis, 2011), exposing spiking neural network simulators to Python seems a
23 pragmatic choice. NEST (Eppler et al., 2009), Neuron (Hines et al., 2009) and CARLsim (Balaji et al.,
24 2020) have all taken this route and now offer a Python interface. Furthermore, newer simulators such as
25 Arbor and Brian2 (Stimberg et al., 2019) have been designed from the ground up with a Python interface.

26 While we have recently demonstrated some very competitive performance results (Knight and Nowotny,
27 2018, 2020) using our GeNN simulator (Yavuz2016), it has not been usable directly from Python. GeNN
28 can already be used as a backend for the Python-based Brian2 simulator (Stimberg et al., 2019) but, while
29 Brian2GeNN (Stimberg et al., 2020) allows Brian2 users to harness the performance benefits GeNN

provides, it is not possible to expose all of GeNN's unique features to Python through the Brian2 API. Specifically, GeNN not only allows users to easily define their own neuron and synapse models but, also 'snippets' for offloading the potentially costly initialisation of model parameters and connectivity onto the GPU. Additionally, GeNN provides a lot of freedom for users to integrate their own code into the simulation loop. In this paper we describe the implementation of PyGeNN – a Python package which aims to expose the full range of GeNN functionality with minimal performance overheads. While implementing new neuron and synapse models in the majority of other GPU simulators requires extending the underlying C++ code, using PyGeNN, models can be defined directly from Python. Finally, we demonstrate the flexibility and performance of PyGeNN in two scenarios where minimising performance overheads is particularly critical.

- In a simulation of a large, highly-connected model of a cortical microcircuit (Potjans and Diesmann, 2014) with small simulation timesteps. Here the cost of copying spike data off the GPU from a large number of neurons every timestep can become a bottleneck.
 - In a simulation of a much smaller model of Pavlovian conditioning (Izhikevich, 2007) where learning occurs over 1 h of biological time and stimuli is delivered – following a complex scheme – throughout the simulation. Here any overheads are multiplied by a large number of timesteps and copying stimuli to the GPU can become a bottleneck.
- Using the facilities provided by PyGeNN, we show that both scenarios can be simulated from Python with only minimal overheads over a C++ implementation.

2 MATERIALS AND METHODS

2.1 GeNN

GeNN (Yavuz et al., 2016) is a library for generating CUDA code for the simulation of spiking neural network models. GeNN handles much of the complexity of using CUDA directly as well as automatically performing device-specific optimizations so as to maximize performance.

GeNN consists of a main library – implementing the API used to define models as well as the generic parts of the code generator – and an additional library for each backend (currently there is a reference C++ backend for generating CPU code and a CUDA backend. An OpenCL backend is under development). Users describe their model by implementing a `modelDefinition` function within a C++ file. For example, a model consisting of 4 Izhikevich neurons with heterogeneous parameters, driven by a constant input current might be defined as follows:

```
void modelDefinition(ModelSpec &model)
{
    model.setDT(0.1);
    model.setName("izhikevich");

    NeuronModels::IzhikevichVariable::VarValues popInit(
        -65.0, -20.0, uninitialisedVar(), uninitialisedVar(),
        uninitialisedVar(), uninitialisedVar());

    model.addNeuronPopulation<NeuronModels::IzhikevichVariable>(
        "Pop", 4, {}, popInit);
}
```

```

71     CurrentSourceModels::DC::ParamValues csParams(10.0);
72
73     model.addCurrentSource<CurrentSourceModels::DC>(
74         "CS", "Neurons", csParams, {});
75 }

```

The *genn-buildmodel* command line tool is then used to compile this file; link it against the main GeNN library and the desired backend library; and finally run the resultant executable to generate the source code required to build a simulation dynamic library (a .dll file on Windows or a .so file on Linux and Mac). This dynamic library can then either be statically linked against a simulation loop provided by the user or dynamically loaded by the users simulation code. To demonstrate this latter approach, this example uses the *SharedLibraryModel* helper class supplied with GeNN to dynamically loads the previously defined model, initialise the heterogenous neuron parameters and print each neuron's membrane voltage every timestep:

```

84 #include "sharedLibraryModel.h"
85
86 int main()
87 {
88     SharedLibraryModel<float> model("./", "izhikevich");
89     model.allocateMem();
90     model.initialize();
91     float *aPop = model.getScalar<float>("a");
92     float *bPop = model.getScalar<float>("b");
93     float *cPop = model.getScalar<float>("c");
94     float *dPop = model.getScalar<float>("d");
95     aPop[0] = 0.02; bPop[0] = 0.2; cPop[0] = -65.0; dPop[0] = 8.0; // RS
96     aPop[1] = 0.1; bPop[1] = 0.2; cPop[1] = -65.0; dPop[1] = 2.0; // FS
97     aPop[2] = 0.02; bPop[2] = 0.2; cPop[2] = -50.0; dPop[2] = 2.0; // CH
98     aPop[3] = 0.02; bPop[3] = 0.2; cPop[3] = -55.0; dPop[3] = 4.0; // IB
99     model.initializeSparse();
100
101     float *vPop = model.getScalar<float>("VPop");
102     while(model.getTime() < 200.0f) {
103         model.stepTime();
104         model.pullVarFromDevice("Pop", "V");
105         printf("%f, %f, %f, %f, %f\n", t, vPop[0], vPop[1], vPop[2], vPop[3]);
106     }
107     return EXIT_SUCCESS;
108 }

```

109 2.2 SWIG

In order to use GeNN from Python, both the model creation API and the *SharedLibraryModel* functionality need to be 'wrapped' so they can be called from Python. While this is possible using the API built into Python itself, a wrapper function would need to be manually implemented for each GeNN function to be exposed which would result in a lot of maintenance overhead. Instead, we chose to use SWIG (Beazley, 1996) to automatically generate wrapper functions and classes. SWIG generates

Python modules based on special interface files which can directly include C++ code as well as special ‘directives’ which directly control SWIG:

```
%module(package="package") package
#include "test.h"
```

where the `%module` directive sets the name of the generated module and the package it will be located in and the `%include` directive parses and automatically generates wrapper functions for a C++ header file. We use SWIG in this manner to wrap both the model building and `SharedLibraryModel` APIs described in section 2.1. However, key parts of GeNN’s API such as the `ModelSpec::addNeuronPopulation` method employed in section 2.1, rely on C++ templates which are not directly translatable to Python. Instead, valid template instantiations need to be given a unique name in Python using the `%template` SWIG directive:

```
%template(addNeuronPopulationLIF) ModelSpec::addNeuronPopulation<NeuronModels::LIF>;
```

Having to manually add these directives whenever a model is added to GeNN would be exactly the sort of maintenance overhead we were trying to avoid by using SWIG. Instead, when building the Python wrapper, we search the GeNN header files for the macros used to declare models in C++ and automatically generate SWIG `%template` directives.

As previously discussed, a key feature of GeNN is the ease with which it allows users to define their own neuron and synapse models as well as ‘snippets’ defining how variables and connectivity should be initialised. Beneath the syntactic sugar described in our previous work (Knight and Nowotny, 2018), new models can be defined in C++ by defining a new class derived from, for example, the `NeuronModels::Base` class. The ability to extend this system to Python was a key requirement of PyGeNN and, by using SWIG ‘directors’, C++ classes can be made inheritable from Python using a single SWIG directive:

```
%feature("director") NeuronModels::Base;
```

2.3 PyGeNN

While GeNN *could* be used from Python via the wrapper generated using the techniques described in the previous section, the resultant code would be unpleasant to use directly. For example, rather than being able to specify neuron parameters using a native Python data structure such as a list or dictionary, you have to use a wrapped type such as `DoubleVector([0.25, 10.0, 0.0, 0.0, 20.0, 2.0, 0.5])`. To provide a more user-friendly and pythonic interface, we have built PyGeNN on top of the wrapper generated by SWIG. PyGeNN combines the separate model building and simulation stages of building a GeNN model in C++ into a single API, likely to be more familiar to users of existing Python-based model description languages such as PyNEST (Eppler et al., 2009) or PyNN (Davison et al., 2008). By combining the two stages together, PyGeNN can provide a unified dictionary-based API for initialising homogeneous and heterogeneous parameters as shown in this re-implementation of the previous example:

```
from pygenn import genn_wrapper, genn_model
model = genn_model.GeNNModel("float", "izhikevich")
model.dT = 0.1
izk_init = {"V": -65.0,
            "U": -20.0,
            "a": [0.02, 0.1, 0.02, 0.02],
```

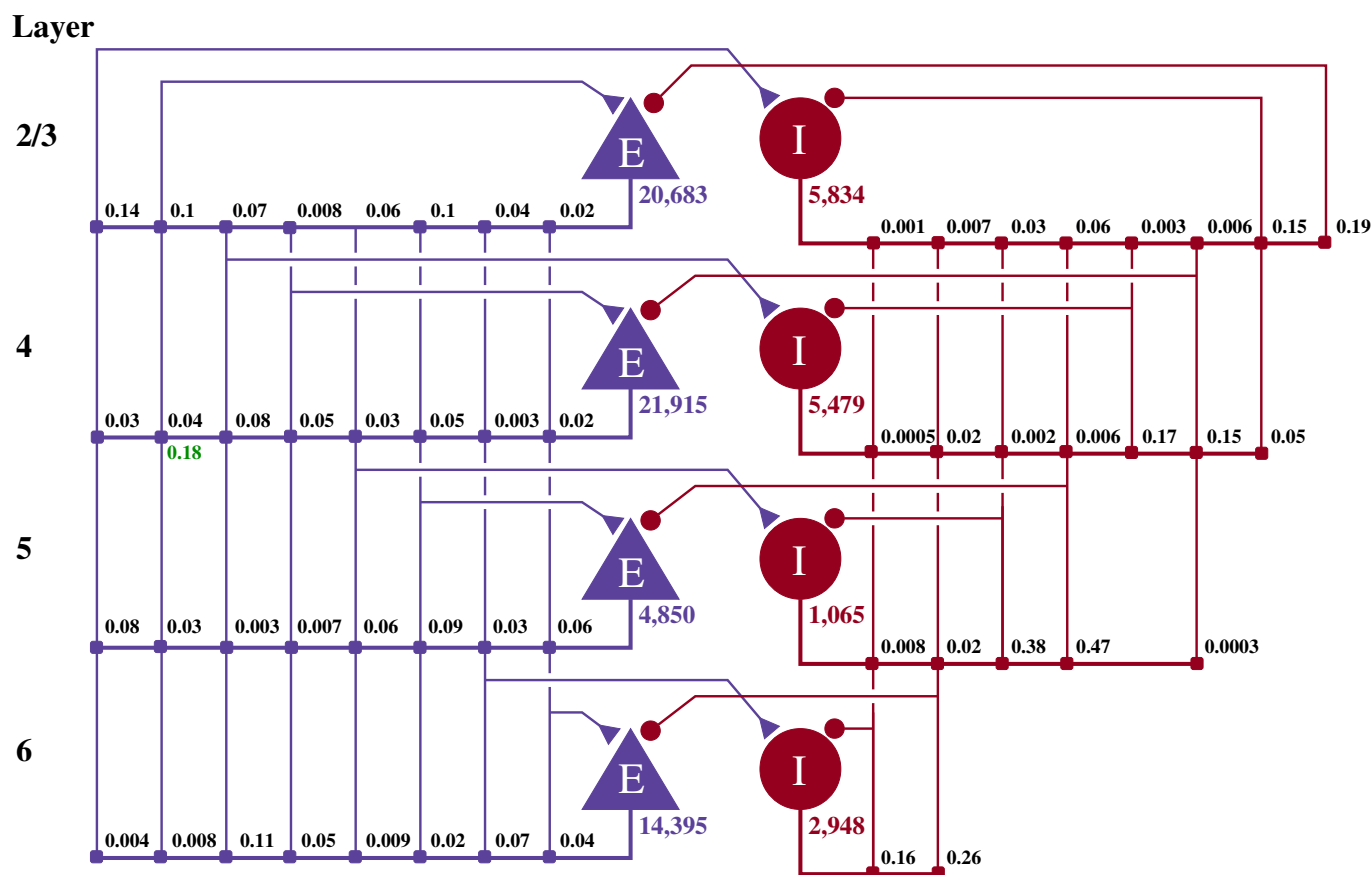


Figure 1. Illustration of the microcircuit model. Blue triangles represent excitatory populations, red circles represent inhibitory populations and the numbers beneath each symbol shows the number of neurons in each population. Connection probabilities are shown in small bold numbers at the appropriate point in the connection matrix. All excitatory synaptic weights are normally distributed with a mean of 0.0878 nA (unless otherwise indicated in green) and a standard deviation of 0.00878 nA. All inhibitory synaptic weights are normally distributed with a mean of 0.3512 nA and a standard deviation of 0.03512 nA.

```

156         "b": [0.2,      0.2,      0.2,      0.2],
157         "c": [-65.0,    -65.0,    -50.0,    -55.0],
158         "d": [8.0,      2.0,      2.0,      4.0]}
159
160 pop = model.add_neuron_population("Pop", 4, "IzhikevichVariable", {}, izk_init)
161 model.add_current_source("CS", "DC", "Neurons", {"amp": 10.0}, {})
162
163 model.build()
164 model.load()
165
166 v = pop.vars["V"].view
167 while model.t < 200.0:
168     model.step_time()
169     model.pull_state_from_device("Neurons")
170     print("%t, %f, %f, %f, %f" % (model.t, v[0], v[1], v[2], v[3]))

```

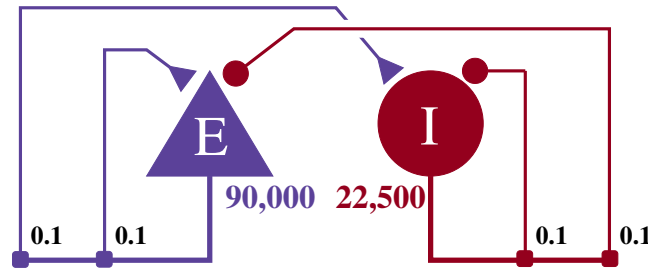


Figure 2. Illustration of the balanced random network model. The blue triangle represents the excitatory population, the red circle represents the inhibitory population, and the numbers beneath each symbol show the number of neurons in each population. Connection probabilities are shown in small bold numbers at the appropriate point in the connection matrix. All excitatory synaptic weights are initialised to 0.045 61 nA and all inhibitory synaptic weights are initialised to 0.228 05 nA. **(TODO: THOMAS: UPDATE THE 800 AND 200 AND PUT SOME SORT OF SYMBOL FOR 3-FACTOR LEARNING ON THE E→E AND E→I.)**

171 Initialisation of variables with homogeneous values – such as the neuron’s membrane voltage – is performed
 172 by GeNN and those with heterogeneous values – such as the a , b and c parameters – are initialised by
 173 PyGeNN when the model is loaded. While the PyGeNN API is more pythonic and, hopefully, more
 174 user-friendly than the C++ interface, it still provides users with the same low-level control over the
 175 simulation. Furthermore, by using SWIG’s numpy (Van Der Walt et al., 2011) interface, the host memory
 176 allocated by GeNN can be accessed directly from Python using the `pop.vars["V"].view` syntax meaning
 177 that no potentially additional expensive copying of data is required. **(TODO: DEFINING NEW NEURON**
 178 **MODELS, PARAMETERS AND VARIABLES)**

179 2.4 Spike recording system

180 Internally, GeNN stores the spikes emitted by a neuron population during one simulation timestep in an
 181 array containing the indices of the neurons that spiked alongside a counter of how many spikes have been
 182 emitted. Previously, recording spikes in GeNN was very similar to the recording of voltages shown in the
 183 previous example code – the array of neuron indices was simply copied from the GPU to the CPU every
 184 timestep. However, especially when simulating models with a small simulation timestep, such frequent
 185 synchronization between the CPU and GPU is costly – especially if a higher-level language such as Python
 186 is involved. Furthermore, biological neurons typically spike at a low rate (in the cortex, the average firing
 187 rate is only around 3 Hz (Buzsáki and Mizuseki, 2014)) meaning that the amount of spike data transferred
 188 every timestep is typically very small. To address both of these sources of inefficiency, we have added a
 189 new data structure to GeNN which stores spike data for many timesteps on device. To reduce the memory
 190 required for this data structure and to make its size independent of neural activity, the spikes emitted by a
 191 population of N neurons in a single simulation timestep are stored in a N bit bitfield where a ‘1’ represents
 192 a spike and a ‘0’ the absence of one. Spiking data over multiple timesteps is then represented by bitfields
 193 stored in a circular buffer. Using this approach, even the spiking output of relatively large models, running
 194 for many timesteps can be stored in a small amount of memory. For example, the spiking output of a model
 195 with 100×10^3 neurons running for 10×10^3 simulation timesteps, required less than 120 MB – a small
 196 fraction of the memory on a modern GPU. While efficiently handling spikes stored in a bitfield is a little
 197 trickier than working with a list of neuron indices, GeNN provides an efficient C++ helper function for
 198 saving the spikes stored in a bitfield to a text file and a numpy-based method for decoding them in PyGeNN.
 199

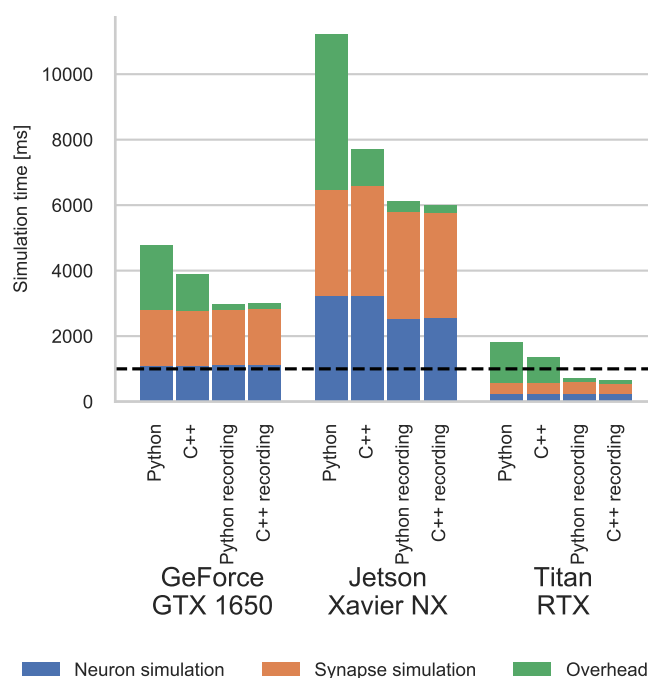


Figure 3. Simulation times of the microcircuit model running on various GPU hardware for 1 s of biological time. ‘Overhead’ refers to time spent in simulation loop but not within CUDA kernels. The dotted horizontal line indicates realtime performance

200 2.5 Cortical microcircuit model

201 Potjans and Diesmann (2014) developed this model of 1 mm^3 of early-sensory cortex. The model
 202 consists of 77 169 LIF neurons, divided into separate populations representing the excitatory and inhibitory
 203 population in each of 4 cortical layers (2/3, 4, 5 and 6) as illustrated by figure 1. Neurons in each population
 204 are connected randomly with numbers of synapses derived from an extensive review of the anatomical
 205 literature. Within each synaptic projection, all synaptic strengths and transmission delays are normally
 206 distributed and, in total, the model has approximately 0.3×10^9 synapses. As well as receiving synaptic
 207 input, each neuron in the network also receives an independent Poisson input current, representing input
 208 from neighbouring un-modelled cortical regions. For a full description of the model parameters, please
 209 refer to Potjans and Diesmann (2014, tables 4 and 5) and for a description of the strategies used by GeNN to
 210 parallelise the initialisation and subsequent simulation of this network, please refer to Knight and Nowotny
 211 (2018, section 2.3). This model requires simulation using a relatively small timestep of 0.1 ms, making the
 212 overheads of copying spikes from the GPU every timestep particularly significant.

213 2.6 Pavlovian conditioning model

214 The cortical microcircuit model described in the previous section is ideal for exploring the performance
 215 of short simulations of relatively large models. However, the performance of longer simulations
 216 of smaller models is equally vital. (TODO: DETERMINE E.G. PERCENTAGE MODELS E.G. ON
 217 OPENSOURCEBRAIN WHICH ARE SMALL). Such models can be particularly troublesome for GPU
 218 simulation as, not only might they not offer enough parallelism to fully occupy the device but, each timestep
 219 can be simulated so quickly that the overheads of launching kernels etc can dominate. Additional overheads
 220 can be incurred when models require injecting external stimuli throughout the simulation. (TODO:
 221 SOMETHING ABOUT NEUROMORPHIC SYSTEMS OFTEN BEING REAL-TIME / BS ACCELERATED

TIME) Longer simulations are particularly useful when exploring synaptic plasticity so, to explore the performance of PyGeNN in this scenario, we simulate a model of Pavlovian conditioning using a three-factor STDP learning rule (Izhikevich, 2007). In this experiment, 100 random groups of 50 neurons (each representing stimuli $S_1 \dots S_{100}$) are chosen from amongst the neurons in a 800 neuron excitatory population and a 200 neuron inhibitory population. Excitatory neurons are modelled as regular-spiking Izhikevich neurons (Izhikevich, 2003) and inhibitory neurons as fast-spiking Izhikevich neurons (Izhikevich, 2003). Stimuli are presented to the network in a random order, separated by intervals sampled from $U(100, 300)$ ms. The neurons associated with an active stimulus are stimulated for a single 1 ms simulation timestep with a current of 40.0 nA, in addition to the random background current of $U(-6.5, 6.5)$ nA, delivered to each neuron throughout the simulation. S_1 is arbitrarily chosen as the Conditional Stimuli (CS) and, whenever this stimuli is presented, a reward in the form of an increase in dopamine is delivered to all the plastic synapses in the network after a delay sampled from $U(0, 1000)$ ms. This delay period is large enough to allow a few irrelevant stimuli to be presented which act as distractors. **(TODO: TALK MORE ABOUT MODEL/LEARNING RULE/REFER BACK TO OLD PAPER)**. The simplest way to implement this stimulation regime is to add a current source to the excitatory and inhibitory neuron populations which adds the uniformly-distributed input current to an externally-controllable per-neuron current. In PyGeNN, the following model can be defined to do just that:

```

239 stim_noise_model = genn_model.create_custom_current_source_class(
240     "stim_noise",
241     param_names=["n"],
242     var_name_types=[("iExt", "scalar", VarAccess_READ_ONLY)],
243     injection_code=
244         """
245         $(injectCurrent, $(iExt) + ($(gennrand_uniform) * $(n) * 2.0) - $(n));
246         """

```

where the `n` parameter sets the magnitude of the background noise, the `$(injectCurrent, I)` function injects a current of I nA into the neuron and `$(gennrand_uniform)` uses cuRAND **(TODO: CITE)** to sample from $U(0, 1)$. Once a current source population using this model has been instantiated and a memory view to `iExt` obtained in the manner described in section 2.3, in timesteps when stimulus injection is required, current can be injected into the list of neurons contained in `stimuli_input_set` with:

```

252 curr_ext_view[stimuli_input_set] = 40.0
253 curr_pop.push_var_to_device("iExt")

```

The same approach can then be used to zero the current afterwards. However, as almost 20 000 stimuli will be injected over the course of a 1 h simulation, in order to reduce potential overheads, we can offload the stimuli delivery entirely to the GPU using the following slightly more complex model:

```

257 stim_noise_model = genn_model.create_custom_current_source_class(
258     "stim_noise",
259     param_names=["n", "stimMagnitude"],
260     var_name_types=[("startStim", "unsigned int"),
261                     ("endStim", "unsigned int", VarAccess_READ_ONLY)],
262     extra_global_params=[("stimTimes", "scalar*")],
263     injection_code=
264         """
265         scalar current = ($(gennrand_uniform) * $(n) * 2.0) - $(n);

```



```

266         if( $(startStim) != $(endStim) && $(t) >= $(stimTimes)[$(startStim)] ) {
267             current += $(stimMagnitude);
268             $(startStim)++;
269         }
270         $(injectCurrent, current);
271         """)

```

272 This model retains the same logic for generating background noise but, additionally, uses a simple sparse
 273 matrix data structure to store the times at which each neuron should have current injected. **(TODO:**
 274 **FIGURE)** The `startStim` and `endStim` variables point to the subset of the `stimTimes` array used by each
 275 neuron's current source and, once the simulation time `$(t)` passes the time pointed to by `startStim`,
 276 current is injected and `startStim` is advanced. This array is stored in a 'extra global parameter' which
 277 is a read-only memory area that can be allocated and populated from PyGeNN, in this case by 'stacking'
 278 together a list of lists of spike times:

```

279 curr_pop.set_extra_global_param("stimTimes", np.hstack(neuron_stimuli_times))

```

3 RESULTS

280 In the following subsection we will analyse the performance of the models introduced in sections 2.5 and 2.6
 281 on a representative selection of NVIDIA GPU hardware:

- 282 • Jetson Xavier NX – a low-power embedded system with a GPU based on the Volta architecture with
 283 8 GB of shared memory.
- 284 • GeForce GTX 1050Ti – a low-end desktop GPU based on the Pascal architecture with 4 GB of
 285 dedicated memory.
- 286 • GeForce GTX 1650 – a low-end desktop GPU based on the Turing architecture with 4 GB of dedicated
 287 memory.
- 288 • Titan RTX – a high-end workstation GPU based on the Turing architecture with 24 GB of dedicated
 289 memory.

290 All of these systems run Ubuntu 18 apart from the system with the GeForce 1050 Ti which runs Windows
 291 10.

292 3.1 Cortical microcircuit model performance

293 Figure 3 shows the simulation times for the full-scale microcircuit mode and, as one might predict, the
 294 Jetson Xavier NX is slower than the two desktop GPUs. However, considering that it only consumes a
 295 maximum of 15 W compared to 75 W or 320 W for the GeForce GTX 1650 and Titan RTX respectively, it
 296 still performs impressively.

297 The time taken to actually simulate the models ('Neuron simulation' and 'Synapse simulation') are the
 298 same when using Python and C++ as all GeNN optimisation options are exposed to PyGeNN. However,
 299 both the PyGeNN and C++ simulations spend a significant amount of every simulation step copying spike
 300 data off the device and storing it in a suitable data structure ('Overhead'). Because Python is an interpreted
 301 language, such operations are inherently slower – this is particularly noticeable on devices with a slower
 302 CPU such as the Jetson Xavier NX. Unlike the desktop GPUs, the Jetson Xavier NX's 8 GB of memory is
 303 shared between the GPU and the CPU meaning that data doesn't have to be copied between GPU and CPU

memory and can instead be accessed by both. **(TODO: RUN XAVIER NX WITHOUT RECORDING AND WITHOUT ZERO COPY)** While, using this shared memory for recording spikes, reduces the overheads associated with copying data off the device, because the GPU and CPU caches are not coherent, caching must be disabled on this memory which reduces the performance of the neuron kernel. However, when the spike recording system described in section 2.4 is used, spike data is kept in GPU memory until the end of the simulation and this overhead is reduced by around a factor of 10. This now means that, on the high-end desktop GPU, this simulation now runs faster than real-time – previously only achievable using a specialised neuromorphic system (Rhodes et al., 2020) and significantly faster than other recently published GPU simulators (Golosio et al., 2020).

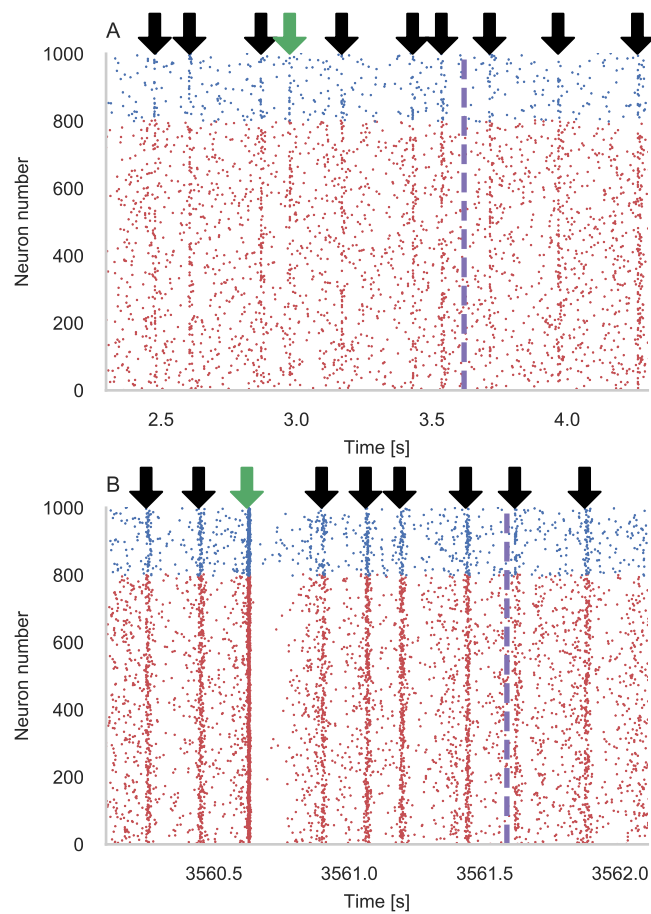


Figure 4. Results of Pavlovian conditioning experiment. Raster plots showing activity centred around first delivery of Conditional Stimulus (CS) during initial (A) and final (B) 50 s of simulation. Downward green arrows indicate times at which CS is delivered and downward black arrows indicate times when other, un-rewarded stimuli are delivered. Vertical dashed lines indicate times at which dopamine is delivered.

312

313 3.2 Pavlovian conditioning performance

Figure 4 shows the results of an example simulation of this model. At the beginning of each simulation (Figure 4A), the neurons representing every stimuli respond equally. However, after 1 h of simulation, the response to the CS becomes much stronger (Figure 4B) – showing that these neurons have been selectively associated even in the presence of the distractors and delayed reward.

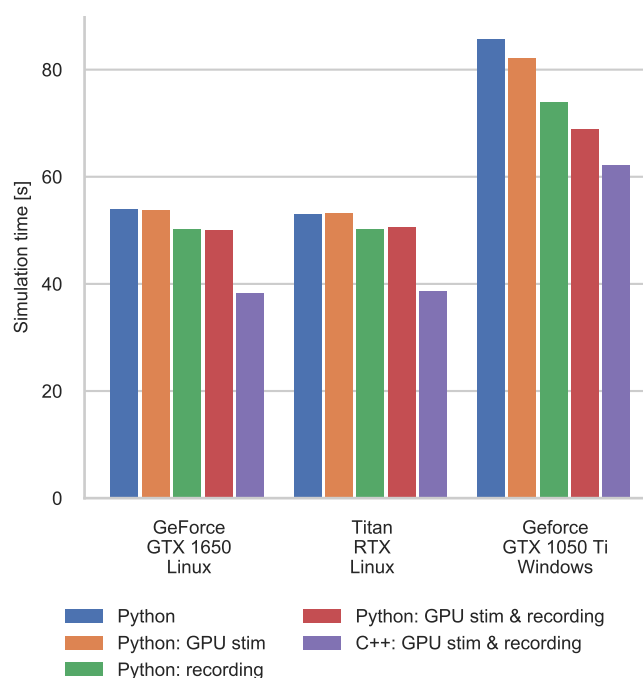


Figure 5. Simulation times of the Pavlovian Conditioning model running on various GPU hardware for 1 h or biological time.

Figure 5 shows the performance of simulations of this model, running on a selection of desktop GPUs using PyGeNN with and without the recording system described in section 2.4 and the optimized stimulus delivery described in section 2.6. These PyGeNN results are compared to a C++ simulation which takes advantage of both optimizations. Interestingly the Titan RTX and GTX 1650 perform identically in this benchmark with speedups ranging from $62\times$ to $72\times$ real-time. This is because, as discussed previously, this model is simply not large enough to fill the 4608 CUDA cores present on the Titan RTX. Therefore, as the two GPUs share the same Turing architecture and have very similar clock speeds (1350 MHz–1770 MHz for the Titan RTX and 1485 MHz–1665 MHz for the GTX 1650), the two GPUs perform very similarly. Furthermore, on these two systems, while using the recording system significantly improves performance, the impact of delivering stimuli on the GPU is minimal. However, the GTX 1050 Ti performs rather differently. Although the clock speed of this device is approximately the same as the other GPUs (1290 MHz–1392 MHz) and it has a similar number of CUDA cores to the GTX 1650, its performance is significantly worse. Furthermore, unlike on the other devices, offloading stimuli delivery to the GPU improves the performance significantly. The difference in performance across all configurations is likely to be due to architectural differences between the older Pascal; and newer Volta and Turing architectures. Specifically, Pascal GPUs have one type of Arithmetic Logic Unit (ALU) which handles both integer and floating point maths whereas, the newer Volta and Turing architectures have equal numbers of dedicated integer and floating point ALUs. This is particularly beneficial for SNN simulations where there is a significant amount of integer maths involved in indexing sparse matrix data structures etc interspersed between the floating point computation. Furthermore, the large performance improvement seen when offloading stimuli delivery to the GPU is likely to be due to overheads relating to the Windows Display Driver Model (WDDM). **(TODO: CONVINCE NSIGHT SYSTEMS TO WORK AND GET WDDM STATS).**

4 DISCUSSION

341 discuss!

- 342 • Turing architecture is great for GeNN! Presented results improve on state-of-the-art.
- 343 • PyGeNN as an intermediate layer - PyNN, ML
- 344 • Cost of C++ - Python calls in models

CONFLICT OF INTEREST STATEMENT

345 The authors declare that the research was conducted in the absence of any commercial or financial
346 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

347 JK and TN wrote the paper. TN is the original developer of GeNN. AK was the original developer of
348 PyGeNN. JK is currently the primary developer of both GeNN and PyGeNN and was responsible for
349 implementing the spike recording system. JK performed the experiments and the analysis of the results that
350 are presented in this work.

FUNDING

351 This work was funded by the EPSRC (Brains on Board project, grant number EP/P006094/1).

ACKNOWLEDGMENTS

352 This is a short text to acknowledge the contributions of specific colleagues, institutions, or agencies that
353 aided the efforts of the authors.

DATA AVAILABILITY STATEMENT

354 The datasets [GENERATED/ANALYZED] for this study can be found in the [NAME OF REPOSITORY]
355 [LINK].

REFERENCES

- 356 Akar, N. A., Cumming, B., Karakasis, V., Kusters, A., Klijn, W., Peyser, A., et al. (2019). Arbor — A
357 Morphologically-Detailed Neural Network Simulation Library for Contemporary High-Performance
358 Computing Architectures. In *2019 27th Euromicro International Conference on Parallel, Distributed
359 and Network-Based Processing (PDP)* (IEEE), 274–282. doi:10.1109/EMPDP.2019.8671560
- 360 Balaji, A., Adiraju, P., Kashyap, H. J., Das, A., Krichmar, J. L., Dutt, N. D., et al. (2020). PyCARL: A
361 PyNN Interface for Hardware-Software Co-Simulation of Spiking Neural Network
- 362 Beazley, D. M. (1996). Using SWIG to control, prototype, and debug C programs with Python. In *Proc.
363 4th Int. Python Conf*
- 364 Buzsáki, G. and Mizuseki, K. (2014). The log-dynamic brain: how skewed distributions affect network
365 operations. *Nature reviews. Neuroscience* 15, 264–78. doi:10.1038/nrn3687
- 366 Carnevale, N. T. and Hines, M. L. (2006). *The NEURON book* (Cambridge University Press)

- Chou, T.-s., Kashyap, H. J., Xing, J., Listopad, S., Rounds, E. L., Beyeler, M., et al. (2018). CARLsim 4: An Open Source Library for Large Scale, Biologically Detailed Spiking Neural Network Simulation using Heterogeneous Clusters. In *2018 International Joint Conference on Neural Networks (IJCNN)* (IEEE), 1–8. doi:10.1109/IJCNN.2018.8489326
- Davison, A. P., Brüderle, D., Eppler, J., Kremkow, J., Müller, E., Pecevski, D., et al. (2008). PyNN: A Common Interface for Neuronal Network Simulators. *Frontiers in neuroinformatics* 2, 11. doi:10.3389/neuro.11.011.2008
- Eppler, J. M., Helias, M., Müller, E., Diesmann, M., and Gewaltig, M. O. (2009). PyNEST: A convenient interface to the NEST simulator. *Frontiers in Neuroinformatics* 2, 1–12. doi:10.3389/neuro.11.012.2008
- Gewaltig, M.-O. and Diesmann, M. (2007). NEST (NEural Simulation Tool). *Scholarpedia* 2, 1430
- Golosio, B., Tiddia, G., De Luca, C., Pastorelli, E., Simula, F., and Paolucci, P. S. (2020). A new GPU library for fast simulation of large-scale networks of spiking neurons, 1–27
- Hines, M. L., Davison, A. P., and Müller, E. (2009). NEURON and Python. *Frontiers in Neuroinformatics* 3, 1–12. doi:10.3389/neuro.11.001.2009
- Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering* 9, 90–95. doi:10.1109/MCSE.2007.55
- Izhikevich, E. M. (2003). Simple model of spiking neurons. *IEEE Transactions on neural networks* 14, 1569–72. doi:10.1109/TNN.2003.820440
- Izhikevich, E. M. (2007). Solving the Distal Reward Problem through Linkage of STDP and Dopamine Signaling. *Cerebral Cortex* 17, 2443–2452. doi:10.1093/cercor/bhl152
- Knight, J. C. and Nowotny, T. (2018). GPUs Outperform Current HPC and Neuromorphic Solutions in Terms of Speed and Energy When Simulating a Highly-Connected Cortical Model. *Frontiers in Neuroscience* 12, 1–19. doi:10.3389/fnins.2018.00941
- Knight, J. C. and Nowotny, T. (2020). Larger GPU-accelerated brain simulations with procedural connectivity. *bioRxiv* doi:10.1101/2020.04.27.063693
- Millman, K. J. and Aivazis, M. (2011). Python for scientists and engineers. *Computing in Science and Engineering* 13, 9–12. doi:10.1109/MCSE.2011.36
- Potjans, T. C. and Diesmann, M. (2014). The Cell-Type Specific Cortical Microcircuit: Relating Structure and Activity in a Full-Scale Spiking Network Model. *Cerebral Cortex* 24, 785–806. doi:10.1093/cercor/bhs358
- Rhodes, O., Peres, L., Rowley, A. G. D., Gait, A., Plana, L. A., Brenninkmeijer, C., et al. (2020). Real-time cortical simulation on neuromorphic hardware. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 378, 20190160. doi:10.1098/rsta.2019.0160
- Stimberg, M., Brette, R., and Goodman, D. F. (2019). Brian 2, an intuitive and efficient neural simulator. *eLife* 8, 1–41. doi:10.7554/eLife.47314
- Stimberg, M., Goodman, D. F., and Nowotny, T. (2020). Brian2GeNN: accelerating spiking neural network simulations with graphics hardware. *Scientific Reports* 10, 1–12. doi:10.1038/s41598-019-54957-7
- Van Der Walt, S., Colbert, S. C., and Varoquaux, G. (2011). The NumPy array: A structure for efficient numerical computation. *Computing in Science and Engineering* 13, 22–30. doi:10.1109/MCSE.2011.37
- Yavuz, E., Turner, J., and Nowotny, T. (2016). GeNN: a code generation framework for accelerated brain simulations. *Scientific reports* 6, 18854. doi:10.1038/srep18854