# PyGeNN: A Python library for GPU-enhanced neural networks

**James C Knight** [1,*]**, Anton Komissarov, Thomas Nowotny** [1]

[1]*Centre 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

## ABSTRACT

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

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

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

## 1 INTRODUCTION

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

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

30  provides, it is not possible to expose all of GeNN's unique features to Python through the Brian2 API.
31  Specifically, GeNN not only allows users to easily define their own neuron and synapse models but, also
32  'snippets' for offloading the potentially costly initialisation of model parameters and connectivity onto
33  the GPU. Additionally, GeNN provides a lot of freedom for users to integrate their own code into the
34  simulation loop. In this paper we describe the implementation of PyGeNN – a Python package which aims
35  to expose the full range of GeNN functionality with minimal performance overheads. We then demonstrate
36  this performance using two large-scale models from the literature.

## 2 MATERIALS AND METHODS

### 2.1 GeNN

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

41  GeNN consists of a main library – implementing the API used to define models as well as the generic
42  parts of the code generator – and an additional library for each backend (currently there is a reference C++
43  backend for generating CPU code and a CUDA backend. An OpenCL backend is under development).
44  Users describe their model by implementing a `modelDefinition` function within a C++ file. For example,
45  a model consisting of 4 Izhikevich neurons with heterogeneous parameters, driven be a constant input
46  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);

    CurrentSourceModels::DC::ParamValues csParams(10.0);

    model.addCurrentSource<CurrentSourceModels::DC>(
        "CS", "Neurons", csParams, {});
}
```

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

```
72  #include "sharedLibraryModel.h"
73
74  int main()
75  {
76      SharedLibraryModel<float> model("./", "izhikevich");
77      model.allocateMem();
78      model.initialize();
79      float *aPop = model.getScalar<float>("a");
80      float *bPop = model.getScalar<float>("b");
81      float *cPop = model.getScalar<float>("c");
82      float *dPop = model.getScalar<float>("d");
83      aPop[0] = 0.02; bPop[0] = 0.2;  cPop[0] = −65.0;   dPop[0] = 8.0;   // RS
84      aPop[1] = 0.1;  bPop[1] = 0.2;  cPop[1] = −65.0;   dPop[1] = 2.0;   // FS
85      aPop[2] = 0.02; bPop[2] = 0.2;  cPop[2] = −50.0;   dPop[2] = 2.0;   // CH
86      aPop[3] = 0.02; bPop[3] = 0.2;  cPop[3] = −55.0;   dPop[3] = 4.0;   // IB
87      model.initializeSparse();
88
89      float *vPop = model.getScalar<float>("VPop");
90      while(model.getTime() < 200.0f) {
91          model.stepTime();
92          model.pullVarFromDevice("Pop", "V");
93          printf("%f, %f, %f, %f, %f\n", t, VPop[0], VPop[1], VPop[2], VPop[3]);
94      }
95      return EXIT_SUCCESS;
96  }
```

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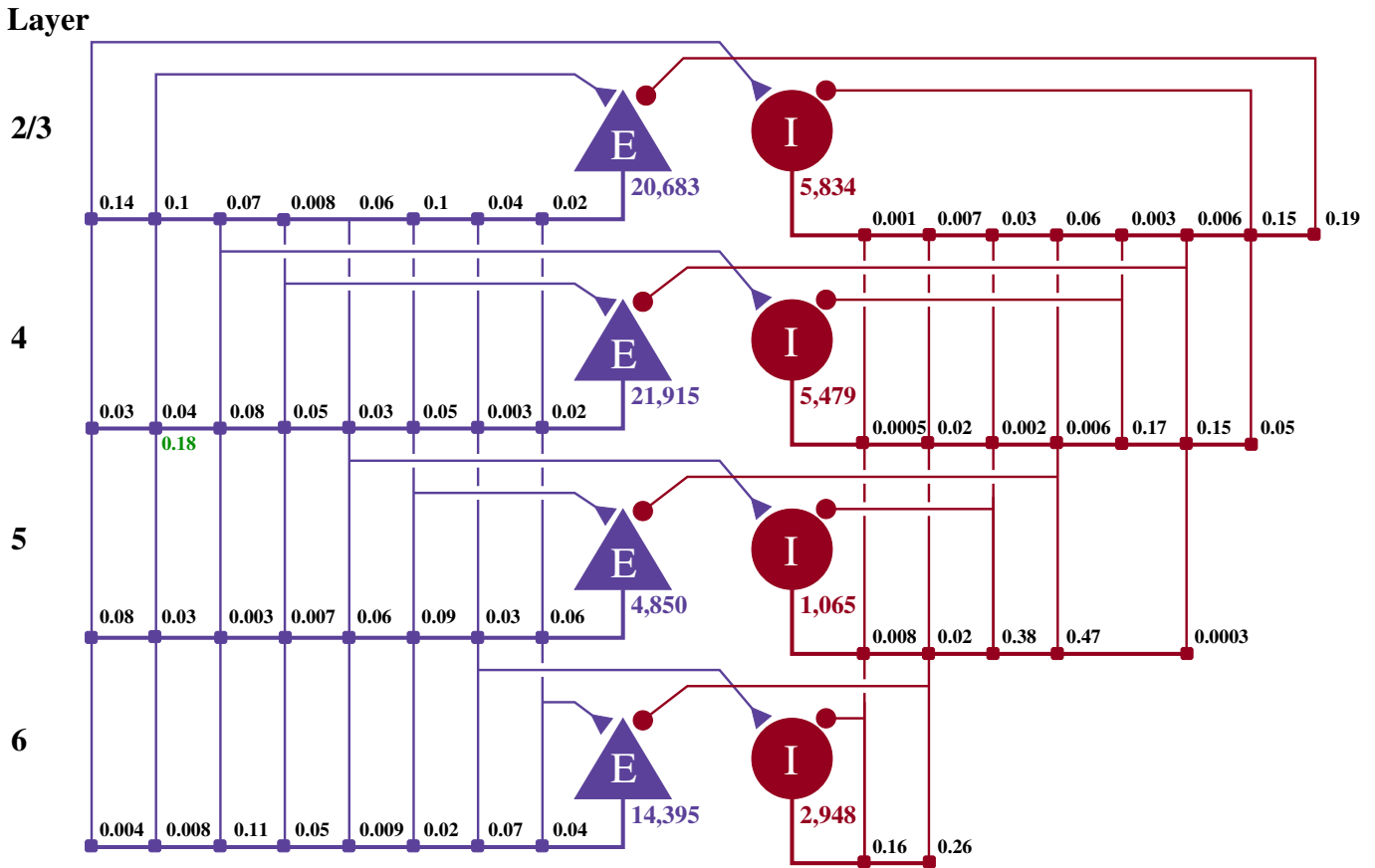
```
%module(package="package") package
%include "test.h"
```

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,

**Layer**



**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 $0.0878\,\text{nA}$ (unless otherwise indicated in green) and a standard deviation of $0.008\,78\,\text{nA}$. All inhibitory synaptic weights are normally distributed with a mean of $0.3512\,\text{nA}$ and a standard deviation of $0.035\,12\,\text{nA}$.

116 we search the GeNN header files for the macros used to declare models in C++ and automatically generate
117 SWIG `%template` directives.

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

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

## 2.3  PyGeNN

126    While GeNN *could* be used from Python via the wrapper generated using the techniques described in the
127 previous section, the resultant code would be unpleasant to use directly. For example, rather than being
128 able to specify neuron parameters using a native Python data structure such as a list or dictionary, you
129 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). This allows PyGeNN to provide a unified dictionary-based API for initialising homogeneous and heterogeneous parameters as shown in this re-implementation of the previous example:
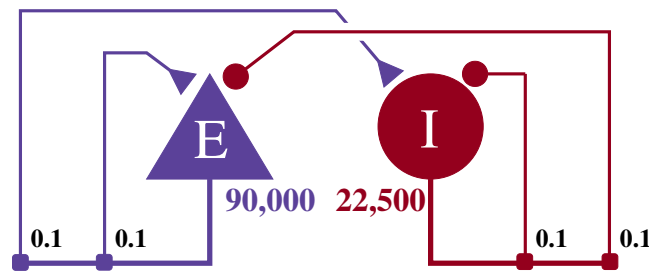
```python
from pygenn import genn_wrapper, genn_model

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],
            "b": [0.2,     0.2,    0.2,     0.2],
            "c": [-65.0,  -65.0,  -50.0,  -55.0],
            "d": [8.0,     2.0,    2.0,     4.0]}

pop = model.add_neuron_population("Pop", 4, "IzhikevichVariable", {}, izk_init)
model.add_current_source("CS", "DC", "Neurons", {"amp": 10.0}, {})

model.build()
model.load()

v = pop.vars["V"].view
while model.t < 200.0:
    model.step_time()
    model.pull_state_from_device("Neurons")
    print("%t, %f, %f, %f, %f" %(model.t, v[0], v[1], v[2], v[3]))
```

While the PyGeNN API is more pythonic and, hopefully, more user-friendly than the C++ interface, it still provides users with the same low-level control over the simulation. Furthermore, by using SWIG's numpy (Van Der Walt et al., 2011) interface, the host memory allocated by GeNN can be accessed directly from Python using the `pop.vars["V"].view` syntax meaning that no potentially additional expensive copying of data is required. **(TODO: DEFINING NEW NEURON MODELS)**
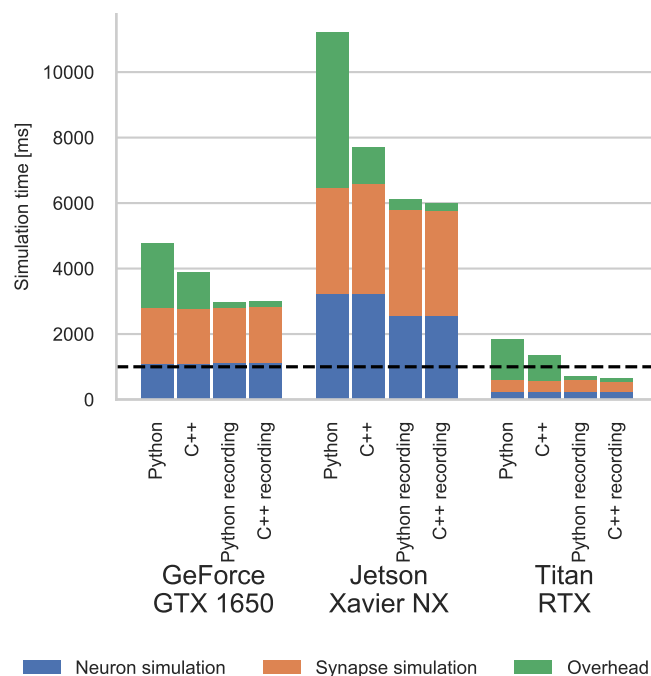
## 2.4 Spike recording system

Internally, GeNN stores the spikes emitted by a neuron population during one simulation timestep in an array containing the indices of the neurons that spiked alongside a counter of how many spikes have been emitted. Previously, recording spikes in GeNN was very similar to the recording of voltages shown in the previous example code – the array of neuron indices was simply copied from the GPU to the CPU every timestep. However, especially when simulating models with a small simulation timestep, such frequent synchronization between the CPU and GPU is costly – especially if a higher-level language such as Python is involved. Furthermore, biological neurons typically spike at a low rate (in the cortex, the average firing rate is only around $3\,\text{Hz}$ (Buzsáki and Mizuseki, 2014)) meaning that the amount of spike data transferred every timestep is typically very small. To address both of these sources of inefficiency, we have added a new data structure to GeNN which stores spike data for many timesteps on device. To reduce the memory

**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 $0.045\,61\,\text{nA}$ and all inhibitory synaptic weights are initialised to $0.228\,05\,\text{nA}$.

175 required for this data structure and to make its size independent of neural activity, the spikes emitted by a
176 population of $N$ neurons in a single simulation timestep are stored in a $N\text{bit}$ bitfield where a '1' represents
177 a spike and a '0' the absence of one. Spiking data over multiple timesteps is then represented by bitfields
178 stored in a circular buffer. Using this approach, even the spiking output of relatively large models, running
179 for many timesteps can be stored in a small amount of memory. For example, the spiking output of a model
180 with $100 \times 10^3$ neurons running for $10 \times 10^3$ simulation timesteps, required less than $120\,\text{MB}$ – a small
181 fraction of the memory on a modern GPU. While efficiently handling spikes stored in a bitfield is a little
182 trickier than working with a list of neuron indices, GeNN provides an efficient C++ helper function for
183 saving the spikes stored in a bitfield to a text file and a numpy-based method for decoding them in PyGeNN.



**Figure 3.** Simulation times of the microcircuit model running on various GPU hardware for $1\,\text{s}$ of biological time. 'Overhead' refers to time spent in simulation loop but not within CUDA kernels. The dotted horizontal line indicates realtime performance

## 2.5  Cortical microcircuit model

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

## 2.6  Random balanced network HPC benchmark model

sdasd

## 3  RESULTS

In the following subsection we will analyse the performance of the models introduced in sections 2.5 and 2.6 on a representative selection of NVIDIA GPU hardware:
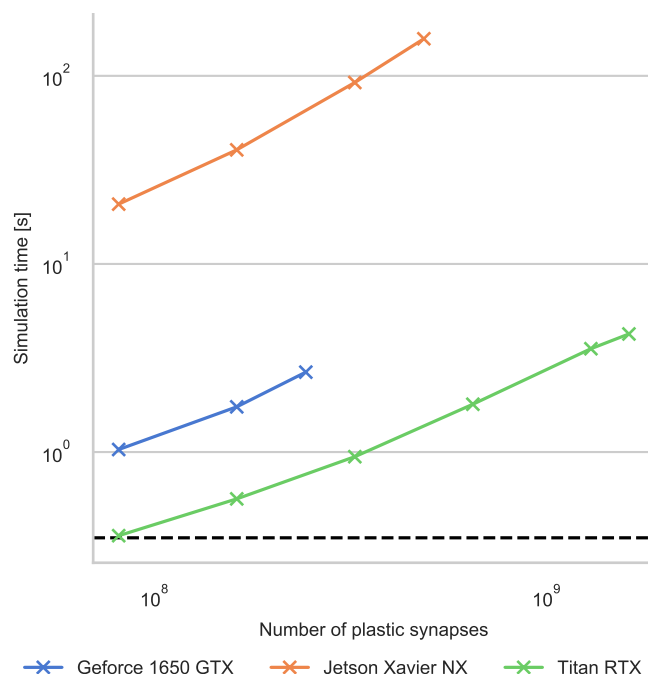
- Jetson Xavier NX – a low-power embedded system with $8\,\mathrm{GB}$ of shared memory.
- GeForce GTX 1650 – a low-end desktop GPU with $4\,\mathrm{GB}$ of dedicated memory.
- Titan RTX – a high-end workstation GPU with $24\,\mathrm{GB}$ of dedicated memory.

### 3.1  Cortical microcircuit model performance

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

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

221 end of the simulation and this overhead is reduced by around a factor of 10. This now means that, on
222 the high-end desktop GPU, this simulation now runs faster than real-time – previously only achievable
223 using a specialised neuromorphic system (Rhodes et al., 2020) and significantly faster than other recently
published GPU simulators (Golosio et al., 2020).



**Figure 4.**

224

## 4  DISCUSSION

225 discuss!

226   • PyGeNN as an intermediate layer - PyNN, ML
227   • Cost of C++ - Python calls in models

## CONFLICT OF INTEREST STATEMENT

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

## AUTHOR CONTRIBUTIONS

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

## FUNDING

## ACKNOWLEDGMENTS

## DATA AVAILABILITY STATEMENT

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

## REFERENCES

239 Akar, N. A., Cumming, B., Karakasis, V., Kusters, A., Klijn, W., Peyser, A., et al. (2019). Arbor — A
240    Morphologically-Detailed Neural Network Simulation Library for Contemporary High-Performance
241    Computing Architectures. In *2019 27th Euromicro International Conference on Parallel, Distributed*
242    *and Network-Based Processing (PDP)* (IEEE), 274–282. doi:10.1109/EMPDP.2019.8671560

243 Balaji, A., Adiraju, P., Kashyap, H. J., Das, A., Krichmar, J. L., Dutt, N. D., et al. (2020). PyCARL: A
244    PyNN Interface for Hardware-Software Co-Simulation of Spiking Neural Network

245 Beazley, D. M. (1996). Using SWIG to control, prototype, and debug C programs with Python. In *Proc.*
246    *4th Int. Python Conf*

247 Buzsáki, G. and Mizuseki, K. (2014). The log-dynamic brain: how skewed distributions affect network
248    operations. *Nature reviews. Neuroscience* 15, 264–78. doi:10.1038/nrn3687

249 Carnevale, N. T. and Hines, M. L. (2006). *The NEURON book* (Cambridge University Press)

250 Chou, T.-s., Kashyap, H. J., Xing, J., Listopad, S., Rounds, E. L., Beyeler, M., et al. (2018). CARLsim 4:
251    An Open Source Library for Large Scale, Biologically Detailed Spiking Neural Network Simulation
252    using Heterogeneous Clusters. In *2018 International Joint Conference on Neural Networks (IJCNN)*
253    (IEEE), 1–8. doi:10.1109/IJCNN.2018.8489326

254 Davison, A. P., Brüderle, D., Eppler, J., Kremkow, J., Muller, E., Pecevski, D., et al. (2008). PyNN: A
255    Common Interface for Neuronal Network Simulators. *Frontiers in neuroinformatics* 2, 11. doi:10.3389/
256    neuro.11.011.2008

257 Eppler, J. M., Helias, M., Muller, E., Diesmann, M., and Gewaltig, M. O. (2009). PyNEST: A convenient
258    interface to the NEST simulator. *Frontiers in Neuroinformatics* 2, 1–12. doi:10.3389/neuro.11.012.2008

259 Gewaltig, M.-O. and Diesmann, M. (2007). NEST (NEural Simulation Tool). *Scholarpedia* 2, 1430

260 Golosio, B., Tiddia, G., De Luca, C., Pastorelli, E., Simula, F., and Paolucci, P. S. (2020). A new GPU
261    library for fast simulation of large-scale networks of spiking neurons , 1–27

262 Hines, M. L., Davison, A. P., and Muller, E. (2009). NEURON and Python. *Frontiers in Neuroinformatics*
263    3, 1–12. doi:10.3389/neuro.11.001.2009

264 Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering* 9,
265    90–95. doi:10.1109/MCSE.2007.55

266 Knight, J. C. and Nowotny, T. (2018). GPUs Outperform Current HPC and Neuromorphic Solutions
267    in Terms of Speed and Energy When Simulating a Highly-Connected Cortical Model. *Frontiers in*
268    *Neuroscience* 12, 1–19. doi:10.3389/fnins.2018.00941

269  Knight, J. C. and Nowotny, T. (2020).  Larger GPU-accelerated brain simulations with procedural
270    connectivity. *bioRxiv* doi:10.1101/2020.04.27.063693
271  Millman, K. J. and Aivazis, M. (2011). Python for scientists and engineers. *Computing in Science and*
272    *Engineering* 13, 9–12. doi:10.1109/MCSE.2011.36
273  Potjans, T. C. and Diesmann, M. (2014). The Cell-Type Specific Cortical Microcircuit: Relating Structure
274    and Activity in a Full-Scale Spiking Network Model. *Cerebral Cortex* 24, 785–806. doi:10.1093/cercor/
275    bhs358
276  Rhodes, O., Peres, L., Rowley, A. G. D., Gait, A., Plana, L. A., Brenninkmeijer, C., et al. (2020). Real-time
277    cortical simulation on neuromorphic hardware. *Philosophical Transactions of the Royal Society A:*
278    *Mathematical, Physical and Engineering Sciences* 378, 20190160. doi:10.1098/rsta.2019.0160
279  Stimberg, M., Brette, R., and Goodman, D. F. (2019). Brian 2, an intuitive and efficient neural simulator.
280    *eLife* 8, 1–41. doi:10.7554/eLife.47314
281  Stimberg, M., Goodman, D. F., and Nowotny, T. (2020). Brian2GeNN: accelerating spiking neural network
282    simulations with graphics hardware. *Scientific Reports* 10, 1–12. doi:10.1038/s41598-019-54957-7
283  Van Der Walt, S., Colbert, S. C., and Varoquaux, G. (2011). The NumPy array: A structure for efficient
284    numerical computation. *Computing in Science and Engineering* 13, 22–30. doi:10.1109/MCSE.2011.37
285  Yavuz, E., Turner, J., and Nowotny, T. (2016). GeNN: a code generation framework for accelerated brain
286    simulations. *Scientific reports* 6, 18854. doi:10.1038/srep18854