

LARGER GPU-ACCELERATED BRAIN SIMULATIONS WITH PROCEDURAL CONNECTIVITY

JAMES C KNIGHT AND THOMAS NOWOTNY

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

ABSTRACT. Large-scale simulations of spiking neural network models are an important tool for improving our understanding of the dynamics and ultimately the function of brains. However, even small mammals such as mice have on the order of 1×10^{12} synaptic connections which, in simulations, are each typically characterized by at least one floating-point value. This amounts to several terabytes of data – an unrealistic memory requirement for a single desktop machine. Large models are therefore typically simulated on distributed supercomputers which is costly and limits large-scale modelling to a few privileged research groups. In this work, we describe extensions to GeNN – our Graphical Processing Unit (GPU) accelerated spiking neural network simulator – that enable it to ‘procedurally’ generate connectivity and synaptic weights ‘on the go’ as spikes are triggered, instead of storing and retrieving them from memory. We find that GPUs are well-suited to this approach because of their raw computational power which, due to memory bandwidth limitations, is often under-utilised when simulating spiking neural networks. We demonstrate the value of our approach with a recent model of the Macaque visual cortex consisting of 4.13×10^6 neurons and 24.2×10^9 synapses. Using our new method, it can be simulated on a single GPU – a significant step forward in making large-scale brain modelling accessible to many more researchers.

1. INTRODUCTION

The brain of a mouse has around 70×10^6 neurons, but this number is dwarfed by the 1×10^{12} synapses which connect them (Herculano-Houzel et al., 2006). In computer simulations of spiking neural networks, propagating spikes involves adding the synaptic input from each spiking presynaptic neuron to the postsynaptic neurons. The information describing which neurons are synaptically connected and with what weight is typically generated before a simulation is run and stored in large arrays. For large-scale brain models this creates high memory requirements, so that they can typically only be simulated on large distributed computer systems using software such as NEST (Gewaltig and Diesmann, 2007) or NEURON (Carnevale and Hines, 2006). By careful design, these simulators can keep the memory requirements for each node almost constant, even when a simulation is distributed across tens of thousands of nodes (Jordan et al., 2018). However, high performance computer (HPC) systems are bulky, expensive and consume a lot of power and are hence typically shared resources, only accessible to a limited number of researchers and for time-limited investigations.

Neuromorphic systems (Frenkel et al., 2018; Furber et al., 2014; Merolla et al., 2014; Qiao et al., 2015; Schemmel et al., 2017) take inspiration from the brain and have been developed specifically for simulating large spiking neural networks more efficiently. One particular relevant feature of the brain is that its memory elements – the synapses – are co-located with the computing elements – the neurons. In neuromorphic systems, this often translates to dedicating a large proportion of each chip to memory. This allows some classes of spiking neural networks to be simulated very efficiently, but reducing the degree of connectivity to fit within the constraints of current neuromorphic systems inevitably changes the dynamics of brain simulations (van Albada et al., 2015). However, while such on-chip memory is fast, it can only be fabricated at relatively low density

E-mail address: J.C.Knight@sussex.ac.uk.

so that many of these systems economize – either by reducing the maximum number of synapses per neuron to as few as 256 or by reducing the precision of the synaptic weights to 6 bit (Schemmel et al., 2017), 4 bit (Frenkel et al., 2018) or even 1 bit (Merolla et al., 2014). Unlike most other neuromorphic systems, the SpiNNaker (Furber et al., 2014) neuromorphic supercomputer is fully programmable and combines large on-chip and external memories, distributed across the system, which enables real-time simulation of large-scale models (Rhodes et al., 2020). This is promising for the future but, due to its prototype nature, the availability of SpiNNaker hardware is limited and even moderately-sized simulations still require a physically large system (9 boards for a model with around 100×10^3 neurons and 300×10^6 synapses (Rhodes et al., 2020)).

Modern GPUs have relatively little on-chip memory and, instead, dedicate the majority of their silicon area to arithmetic logic units. GPUs use dedicated hardware to rapidly switch between tasks so that the latency of accessing external memory can be ‘hidden’ behind computation, as long as there is sufficient computation to be performed. For example, the memory latency of a typical modern GPU can be completely hidden if each CUDA core performs approximately 10 arithmetic operations per byte of data accessed from memory. Unfortunately, propagating a spike in a spiking neural network simulation typically requires reading a synaptic weight (≥ 2 B) and postsynaptic index (≥ 2 B) from memory and then writing back the synaptic input (≥ 4 B), while performing many fewer than the corresponding 80 instructions. This makes spike propagation highly memory bound. Nonetheless, we have shown in previous work (Knight and Nowotny, 2018) that, as GPUs have significantly higher total memory bandwidth than even the fastest CPU, moderately sized models of around 100×10^3 neurons and 1×10^9 synapses can be simulated on a single GPU with competitive speed and energy consumption. However, individual GPUs do not have enough memory to simulate larger brain models and, although small numbers of GPUs can be connected using the high-speed NVLink (NVIDIA Corporation, 2020b) interconnect, larger GPU clusters suffer from the same communication overheads as any other distributed HPC system.

In this work, we present a novel approach that uses the large amount of computational power available on a GPU to reduce both memory and memory bandwidth requirements and enable large-scale brain simulations on a single GPU workstation.

2. RESULTS

In the following subsections, we first present two recent innovations in our GeNN simulator (Yavuz et al., 2016) which enable simulations of very large models on a GPU. We then demonstrate the power of the new features by simulating a recent model of the Macaque visual cortex (Schmidt et al., 2018b) with 4.13×10^6 neurons and 24.2×10^9 synapses.

2.1. Procedural connectivity. The first crucial innovation that enables large-scale simulations on a GPU is what we call ‘procedural connectivity’. In a brain simulation, neurons and synapses can be described by a variety of mathematical models but these are eventually all translated into time or event-driven update algorithms (Brette et al., 2007). Our GeNN simulator (Yavuz et al., 2016) uses code generation to convert neuron and synapse update algorithms – described using ‘snippets’ of C-like code – into CUDA code for efficient GPU simulation. Before a simulation can be run, its parameters, in particular the state variables and the synaptic connectivity, need to be initialised. Traditionally, this is done by running initialisation algorithms on the main CPU prior to the simulation. The results are stored in CPU memory, uploaded to GPU memory and then used during the simulation. We have recently extended GeNN to use code generation from code snippets to also generate efficient, parallel code for model initialisation (Knight and Nowotny, 2018). Offloading initialisation to the GPU in this way made it around $20\times$ faster on a desktop PC (Knight and Nowotny, 2018), demonstrating that initialisation algorithms are well-suited for GPU acceleration. Here, we are going one step further. We realised that, if each synaptic connection can be re-initialised in less than the 80 operations required to hide the latency incurred when fetching its 8 B of state from memory, it could be faster and vastly more memory efficient to regenerate synaptic connections on demand rather than storing them in memory. This is the concept of procedural connectivity. It is applicable whenever



FIGURE 1. Scaling of 1s balanced random network simulation using different algorithms on a range of modern GPUs. **A** Jetson TX2. **B** GeForce MX130. **C** GeForce GTX 1650. **D** Titan RTX.

synapses are static – plastic synapses which change their weights during a simulation will have to be simulated in the traditional way. Although a similar approach was used by Eugene Izhikevich for simulating an extremely large thalamo-cortical model with 1×10^{11} neurons and 1×10^{15} synapses on a modest PC cluster in 2005 (Izhikevich, 2005) – an incredible achievement – it has not been subsequently applied to modern hardware.

We implemented procedural connectivity in GeNN by repurposing our previously developed parallel initialisation methods. Instead of running them once for all synapses at the beginning of the simulation, we rerun the methods during the simulation to regenerate the outgoing synapses of each neuron that fires a spike and immediately use the identified connections and weights to run the post-synaptic code which calculates the effect of the spike onto other neurons. This is possible because the outgoing synaptic connections from each neuron are typically largely independent from those of other neurons as we shall see from typical examples below.

In the absence of knowledge of the exact microscopic connectivity in the brain, there are a number of typical connectivity schemes that are used in brain models. We will now discuss two typical examples and how they can be implemented efficiently on a GPU. One very common connectivity scheme is the ‘fixed probability connector’ for which each neuron in the presynaptic population is connected to each neuron in the postsynaptic population with fixed probability P_{conn} . The postsynaptic targets of any presynaptic neuron can hence be sampled from a Bernoulli process with success probability P_{conn} . One simple way of sampling from the Bernoulli process is to repeatedly draw samples from the uniform distribution $\text{Unif}[0, 1]$ and generate a synapse if the sample is less than P_{conn} . However, for sparse connectivity ($P_{\text{conn}} \ll 1$), it is much more efficient to sample from the geometric distribution $\text{Geom}[P_{\text{conn}}]$ which governs the number of Bernoulli trials until the next success (i.e. a synapse). The geometric distribution can be sampled in constant time by inverting the cumulative density function of the equivalent continuous distribution (the exponential distribution) to obtain $\frac{\log(\text{Unif}[0, 1])}{\log(1 - P_{\text{conn}})}$ (Devroye, 2013, p499). Note that, if we were to directly draw from the uniform distribution, the sampling for each potential synapse would be independent from any other potential synapse and all these operations could be performed in parallel. However, for the more efficient ‘geometric sampling’ employed here, the sampling for the post-synaptic targets of a presynaptic neuron must be done serially, but is still independent from the sampling for any other presynaptic neuron.

Another common scheme for defining connectivity is the ‘fixed total number connector’ in which a fixed total number N_{syn} of synapses is placed between randomly chosen partners from the pre- and postsynaptic populations. In order to initialise this connectivity in parallel, the number of synapses that originate from each of the N_{pre} presynaptic neurons must first be calculated by sampling from the multinomial distribution $\text{Mult}[N_{\text{syn}}, \{P_n, P_n, \dots, P_n\}]$, where $P_n = \frac{1}{N_{\text{pre}}}$, on the host CPU up front because these

numbers need to add to N_{syn} and are hence not independent. However, once the numbers of outgoing synapses are determined, the postsynaptic targets for a presynaptic neuron can be generated very efficiently in parallel by sampling from the discrete uniform distribution $\text{Unif}[0, N_{\text{post}}]$ where N_{post} is the size of the postsynaptic population. Note, that this can only be done because the targets of each presynaptic neuron are independent from those of any other presynaptic neuron. Where synaptic weights and delays are not constant across synapses, but are described by some statistical distribution, they can also be sampled independently from each other and hence in parallel.

In order to use these parallel initialisation schemes for procedural connectivity, we require reproducible pseudorandom numbers that can be generated independently for each presynaptic neuron. In principle this could be done with ‘conventional’ pseudorandom number generators (PRNGs), but each presynaptic neuron would need to maintain its own PRNG state which would lead to a significant memory overhead. Instead, we use the ‘counter-based’ Philox4×32-10 PRNG (Salmon et al., 2011). Counter-based PRNGs are designed for parallel applications and essentially consist of a pseudo-random bijective function which takes a counter as an input (for Philox4×32-10 a 128 bit number) and outputs a random number. In contrast to conventional PRNGs, this means that generating the n^{th} random number in a stream has exactly the same cost as generating the ‘next’ random number, allowing us to trivially divide up the random number stream between multiple parallel processes (in this case presynaptic neurons).

For an initial demonstration of the performance and scalability of procedural connectivity, we simulated a balanced random network initially designed to investigate signal propagation through cortical networks (Vogels and Abbott, 2005) but subsequently widely used as a scalable benchmark (Brette et al., 2007). The network consists of N integrate-and-fire neurons, partitioned into $\frac{4N}{5}$ excitatory and $\frac{N}{5}$ inhibitory neurons. The two populations of neurons are connected to each other and with themselves with fixed probability $P_{\text{conn}} = 10\%$.

We ran simulations of this network at scales ranging from 1×10^3 to 1×10^6 neurons (100×10^3 to 100×10^9 synapses respectively) on a representative selection of NVIDIA GPU hardware: Jetson TX2, a low-power embedded system with 8 GB (shared memory); Geforce MX130, a laptop GPU with 2 GB; Geforce GTX 1650, a low-end desktop GPU with 4 GB; and Titan RTX, a high-end workstation GPU with 24 GB. Fig. 1 shows the duration of these simulations using our new procedural approach or using the standard approach of storing synaptic connections in memory employing two different data structures. Both data structures are described in more detail in our previous work (Knight and Nowotny, 2018) but briefly, in the ‘sparse’ data structure, a presynaptic neuron’s postsynaptic targets are represented as an array of indices whereas, in the ‘bitfield’ data structure, they are represented as a N_{post} array of bits where a ‘1’ at position i indicates that there is a connection to postsynaptic neuron i and a ‘0’ that there is not. None of our devices have enough memory to store the 100×10^9 synapses required for the largest scale using either data structure but, at the 100×10^3 neuron scale, the bitfield data structure allows the model to fit into the memory of several devices it otherwise would not. However, not only is the new procedural approach the *only* way of simulating models at the largest scales but, as Fig. 1 illustrates, even at smaller scales the performance of the procedural approach is competitive with and sometimes better than the standard approach. All of the synapses in this model have the same synaptic weight meaning that they can be hard-coded into the procedural connectivity kernels. However, if weights vary across synapses, the ‘bitfield’ cannot be used and the memory constraints for the ‘sparse’ representation become even more severe.

2.2. Kernel merging. NVIDIA GPUs are typically programmed in CUDA using a Single Instruction Multiple Thread (SIMT) paradigm where programmers write ‘kernel’ functions containing serial C-like code which is then executed in parallel across many virtual threads. We call our second innovation ‘kernel merging’ and it relates to the way these kernels are implemented. While the procedural connectivity presented in the previous section allows simulating models which would otherwise not fit into the memory of a GPU, there are additional problems when using code generation for models with a large number

of neuron and synapse populations. GeNN and other SNN simulators which use code generation to generate all of their simulation code (Blundell et al., 2018) (as opposed to, for example NESTML (Plotnikov et al., 2016), which uses code generation only to generate neuron simulation code) generate separate pieces of code for each population of neurons and synapses. This allows optimizations such as hard-coding constant parameters and, although generating code for models with many populations will result in large code size, C++ CPU code can easily be divided between multiple modules and compiled in parallel, minimizing the effects on build time. However, GPUs can only run a small number of kernels – which are equivalent to modules in this context – simultaneously (128 on the latest NVIDIA GPUs (NVIDIA Corporation, 2019, p278)). Therefore, GeNN employs a “Kernel Fusion” approach (Wang et al., 2010) where multiple neuron populations are simulated within each kernel. CUDA threads are grouped into ‘thread blocks’ and to maximise performance, the code running on threads within a thread block should not ‘diverge’. To minimise this divergence, we add ‘padding’ threads around each population of neurons so population and block boundaries are aligned. Therefore, in the following pseudocode we simulating 3 populations of 100 neurons each in a single kernel and – assuming a thread block size of 32 – we pad each population to 128 threads:

```
void updateNeurons() {
    if(thread < 128) {
        if(thread < 100) {
            // Update neuron population A
        }
    } else if(thread >= 128 && thread < 256) {
        if((thread - 128) < 100) {
            // Update neuron population B
        }
    } else if(thread >= 256 && thread < 384) {
        if((thread - 256) < 100) {
            // Update neuron population C
        }
    }
}
```

This approach works well for a small number of populations but, as Fig. 2A illustrates, when we partition a model consisting of a single population of 1 000 000 LIF neurons, stimulated by a Gaussian current input into an increasingly large number of ever smaller populations, compilation time increases super-linearly and quickly becomes impractical. Furthermore, the simulation also runs slower with a large number of populations (Fig. 2B). Normally, we would expect this model to be memory bound as each thread in the model reads 32 B of data (8 B of neuron state and 24 B of RNG state) and, as discussed above, hiding the latency of these memory accesses would require approximately 320 arithmetic operations – many more than required to sample an input current from the normal distribution and update a LIF neuron. Fig. 2C – obtained using data from the NVIDIA Nsight compute profiler (NVIDIA Corporation, 2020a) – shows that this is true for small numbers of populations. In this case, the memory system is around 90% utilised. However, when the model is partitioned into larger numbers of smaller populations, the memory is used less efficiently and the kernel becomes latency bound, i.e. neither memory *nor* compute are used efficiently. Investigating further, we found that this drop in performance was accompanied by an increasing number of “No instruction” stalls (Fig. 2D) which are events that prevent the GPU from doing any work during a clock cycle. These particular events are likely to be caused by “Excessively jumping across large blocks of assembly code” (NVIDIA Corporation, 2020a, p47), which makes sense as we are generating kernels with hundreds of thousands of lines of code.

Several neural modelling tools including Brian2 (Stimberg et al., 2019) provide modellers with tools to work with ‘slices’ of neuron populations, allowing models to be defined with fewer populations. However, if a model is defined by connecting these slices together,

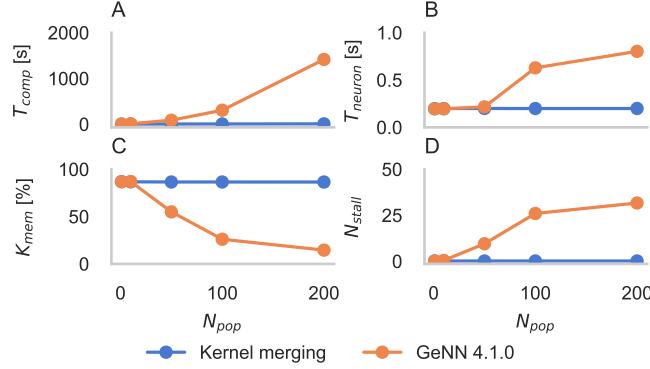


FIGURE 2. Performance of a 1 s simulation of 1×10^6 LIF neurons driven by a Gaussian input current, partitioned into varying numbers (N_{pop}) of populations and running on a workstation equipped with a Titan RTX GPU. **A** Compilation time (T_{comp}) using GCC 7.5.0. **B** Neuron kernel time (T_{neuron}) for a 1 s simulation. **C** Memory throughput (K_{mem}) reported by NVIDIA Nsight compute profiler “Speed of light” metric. **D** Number of “No instruction” stalls reported by NVIDIA Nsight compute profiler (N_{stall}).

the resulting connectivity is the result of *multiple* simple connection rules of the type discussed in the previous section, making it much more difficult to apply our procedural connectivity approach. Furthermore, such an approach places the responsibility for structuring a model in such a way that it can be simulated efficiently onto the modellers, who often prefer to concentrate on the science and organise populations according to anatomy or physiology.

To address the issue of too many populations, we developed a new code generator for GeNN which applies a Single Program Multiple Data (SPMD) approach and ‘merges’ the model description, grouping together populations which can be simulated using the same generated code. From this merged description, structures are generated to store the pointers to state variables and parameter values which are still allowed to differ between merged populations:

```
struct NeuronUpdateGroup {
    unsigned int numNeurons;
    float* V;
};
```

An array of these structures is then declared for each merged population and each element is initialised with pointers to state variables and parameter values:

```
NeuronUpdateGroup neuronUpdateGroup[3];
neuronUpdateGroup[0] = {100, VA};
neuronUpdateGroup[1] = {100, VB};
neuronUpdateGroup[2] = {100, VC};
```

where VA is a pointer to the array containing the state variable ‘V’ of populations ‘A’ and so on. In order for a thread to determine which neuron in which population it should simulate, we generate an additional data structure – an array containing a cumulative sum of threads used for each population:

```
unsigned int startThread[3] = {0, 128, 256};
```

Each thread performs a simple binary search within this array to find the index of the neuron and population it should simulate. As Fig. 2 shows, this approach solves the observed issues with compilation time and simulation performance.

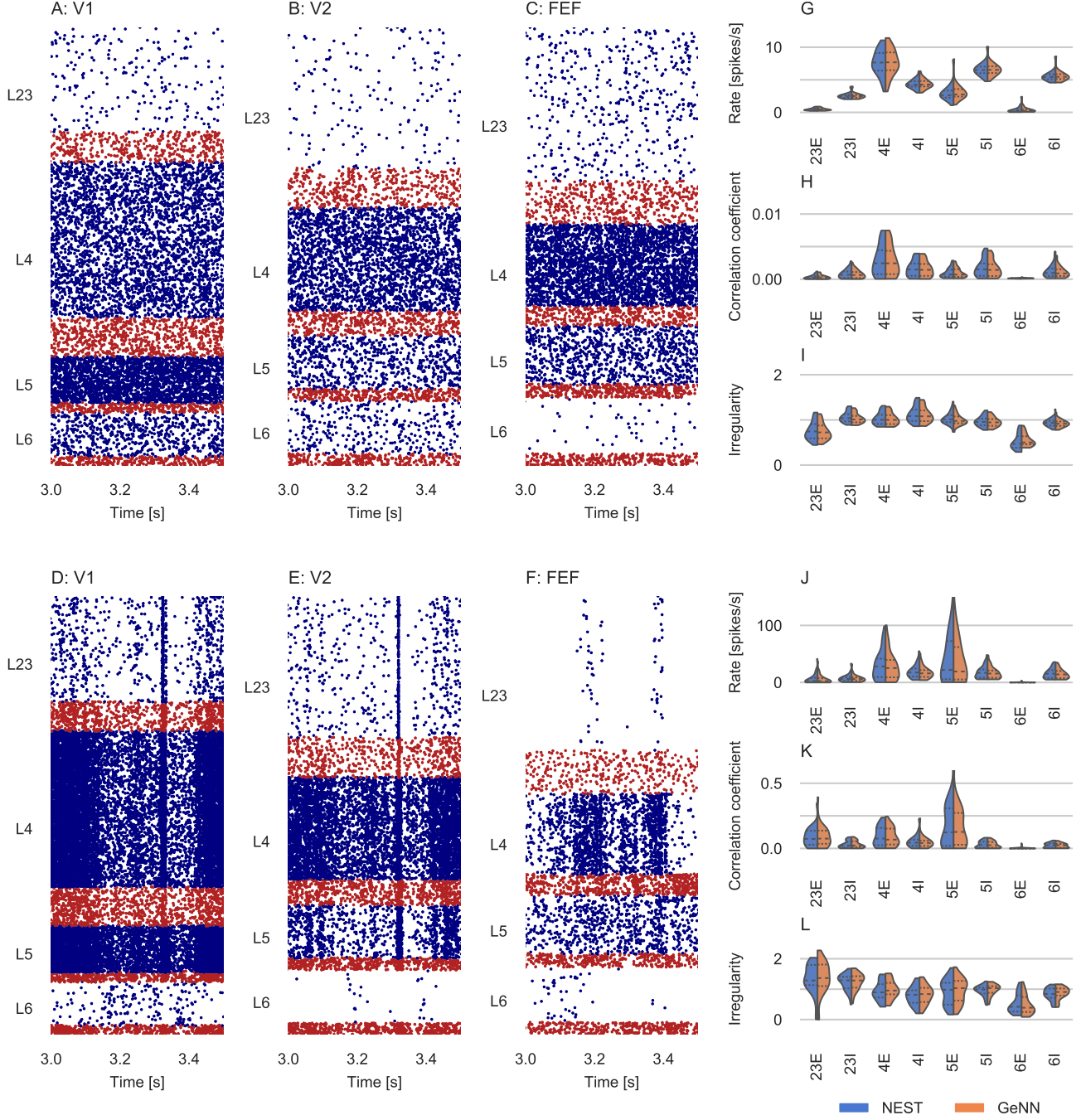


FIGURE 3. Results of full-scale multi-area model simulation in ground resting states. **A-F** Raster plots from GeNN simulation showing spiking activity of 3% of the neurons in area V1 (A,D), V2 (B,E), and FEF (C,F). Blue: excitatory neurons, red: inhibitory neurons. **G-L** Spiking statistics for each population across all 32 areas simulated using GeNN and NEST shown as split violin plots. Solid lines: medians, Dashed lines: Interquartile range. **G,J** Population-averaged firing rates. **H,K** Average pairwise correlation coefficients of spiking activity. **I,L** Irregularity measured by revised local variation LvR (Shinomoto et al., 2009) averaged across neurons.

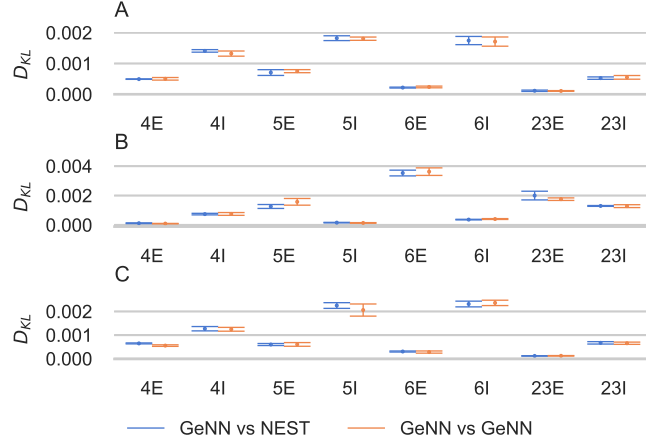


FIGURE 4. Comparison of full-scale multi-area model spike statistics between multiple GeNN simulations with different seeds; and NEST simulations. Comparisons calculated using the Kullback-Leibler (KL) divergence. **A** Neuron firing rates **B** Pairwise correlation coefficients of spiking activity. **C** Irregularity measured by revised local variation LvR (Shinomoto et al., 2009).

2.3. The multi-area model. Due to lack of computing power and sufficiently detailed connectivity data, previous models of the cortex have either focussed on modelling individual local microcircuits at the level of individual cells (Izhikevich and Edelman, 2008; Potjans and Diesmann, 2014) or modelling multiple connected areas at a higher level of abstraction (Cabral et al., 2014). However, recent data (Belitski et al., 2008) has shown that cortical activity has distinct features at both the global and local levels which can only be captured by modelling interconnected microcircuits at the level of individual cells. The recent multi-area model (Schmidt et al., 2018a,b) is an example of such multi-scale modeling. It uses scaled versions of a previous, 4 layer microcircuit model (Potjans and Diesmann, 2014) to implement 1 mm^2 ‘patches’ for 32 areas of the macaque visual cortex. The patches are connected together according to inter-area axon tracing data from the CoCoMac (Bakker et al., 2012) database, further refined using additional anatomical data (Markov et al., 2014a) and heuristics (Ercsey-Ravasz et al., 2013) to obtain estimates for the number of synapses between areas. The synapses are distributed between populations in the source and target area using layer-specific tracing data (Markov et al., 2014b) and cell-type-specific dendritic densities (Binzegger et al., 2004). Individual populations are connected by the fixed number connectors described above. For a full description of the multi-area model please see Schmidt et al. (2018a,b). In 2018, this model was simulated using NEST (Gewaltig and Diesmann, 2007) on an IBM Blue Gene/Q supercomputer. Because each Blue Gene/Q compute node only had a small amount of memory, these simulation were distributed across an entire rack (1024 compute nodes) of the system and simulating 1 s of biological time took approximately 12 min (Schmidt et al., 2018b). However, on the newer ‘JURECA’ supercomputer – where each node has 128 GB of memory – this model can be simulated on as few as 11 compute nodes. Furthermore, by distributing the model across 160 compute nodes, it can be initialized in approximately 8 min and simulating 1 s of biological takes only 31 s (van Albada et al., *ress*).

The multi-area model consists of 4.13×10^6 neurons in 254 populations and 24.2×10^9 synapses in 64516 projections. Without kernel merging, it would therefore be unlikely that the model would compile or simulate at a workable speed using GeNN. Additionally, unlike the model we benchmarked previously, each synapse in this model has an independent weight and synaptic delay sampled from a normal distribution so the bitfield data structure cannot be used. Even if we assume that 16 bit floating-point would provide sufficient weight precision, that delays could be expressed as 8 bit integers and that neuron

populations are all small enough to be indexed using 16 bit indices, our sparse data structure would still require 5 B per synapse, such that the complete synaptic data would need over 100 GB of GPU memory. While a cluster of GPUs connected using NVLink could be built with this much memory, it is more than any single GPU has available. However, using procedural connectivity, we are able to simulate this model on a single workstation with a Titan RTX GPU.

In order to validate our GeNN simulations, we ran a 10.5 s simulation of the multi-area model in a ‘ground state’ where inter-area connections have the same strength as intra-area connections and a 100.5 s simulation in a ‘resting state’ where inter-area connections are $1.9\times$ stronger. Initialization of our model took 6 min (3 min of which was spent generating and compiling code) and simulation of each biological second took 7.9 min in the ground state and 8.5 min in the resting state (averaged over 3 simulation runs). While this is significantly faster than the older Blue Gene/Q supercomputer, it is around $15\times$ slower than the newer JURECA system. However, because individual synapses simulated using procedural connectivity do not require initialization, in shorter simulations this large performance difference can be somewhat overcome by the faster initialisation time of the GeNN simulations. For example a 1 s simulation would take just under 11 min using GeNN (if the model is not recompiled) and around 8.5 min using 160 JURECA nodes – only 2.5 min faster.

Fig. 3A-C shows some example spike rasters from three of the modelled areas, illustrating the asynchronous irregular nature of the model’s ground state whereas Fig. 3D-F illustrate the characteristic irregular activity and population bursts of the same areas in the resting state. Next, we calculated the per-layer distributions of rates, spike-train irregularity and cross-correlation coefficients across all areas (disregarding the first 500 ms of simulation) and compared them to the same measures obtained from spike trains generated by the supercomputer simulations. We calculated irregularity using the revised local variation LvR (Shinomoto et al., 2009), averaged over a subsample of 2000 neurons and cross-correlation from spike histograms with 1 ms bins, calculated from a subset of 2000 non-silent neurons. The distributions of these values, obtained from the NEST and GeNN simulations, are shown as violin plots in Fig. 3G-L. Upon visual inspection, the distributions are very similar but, to assess how meaningful the differences between these distributions are, we compared the distributions obtained from 3 GeNN simulations with different random number generator seeds against those obtained from NEST. We performed this comparison by computing histograms of the three measures – using bin sizes determined with the Freedman-Diaconis rule (Freedman and Diaconis, 1981) – and comparing pairs of distributions using the Kullback-Leibler divergence D_{KL} . The results of these comparisons are shown in Fig. 4 and suggest that the influence of the random seeds is comparable to the influence of random seeds and simulator combined – indicating that the choice of simulator has little effect on the model dynamics.

3. DISCUSSION

In this work we have presented a novel approach for large-scale brain simulation on GPU devices which entirely removes the need to store connectivity data in memory. We have shown that this approach allows us to simulate a cortical model with 4.13×10^6 neurons and 24.2×10^9 synapses (Schmidt et al., 2018a,b) on a single modern GPU. While this represents a significant step forward in terms of making large-scale brain modelling tools accessible to a large community of brain researchers, this model still has around $20\times$ fewer neurons and $40\times$ fewer synapses than the brain of even a small mammal such as a mouse (Herculano-Houzel et al., 2006). Our implementation of the multi-area model requires a little over 12 GB of GPU memory, with the majority (8.5 GB) being used for the circular dendritic delay buffers (see Knight and Nowotny (2018)). These are a per-neuron (rather than per-synapse) data structure but, because the inter-area connections in the model have delays of up to 500 simulation timesteps, the delay buffers become quite large. However, for models without large delays such as the balanced random network simulated in section 2.1, only around 20 B is required per neuron meaning that, theoretically, over 1×10^9 neurons could be simulated on a 24 GB GPU.

Parameter	Procedural connectivity benchmark	Merging benchmark	Multi-area model
τ_m [ms]	20	20	2
V_{rest} [mV]	-60.0	-70.0	-65
V_{th} [mV]	-50.0	-51.0	-50
R_m [M Ω]	20	20	40
τ_{syn} [ms]	5/10 ¹	—	0.5
τ_{ref} [ms]	5	2	2
I_{ext_j} [nA]	0.55	1.00 \pm 0.25	Poisson ²
w_{ij} [nA]	$\frac{3.2}{N} / \frac{40.8}{N}$ ¹	—	Various ²

TABLE 1. Model parameters.

¹Excitatory/Inhibitory.²Please refer to Schmidt et al. (2018b, Table 1,2)

One important aspect of large-scale brain simulations not addressed in this work is synaptic plasticity and its role in learning. As discussed by Knight and Nowotny (2018), GeNN supports a wide variety of synaptic plasticity rules. In order to modify synaptic weights, they need to be stored in memory rather than generated procedurally. However, connectivity could still be generated procedurally, potentially halving the memory requirements of models with synaptic plasticity. This would be sufficient for synaptic plasticity rules that only require access to presynaptic spikes and postsynaptic neuron states (Brader et al., 2007; Clopath et al., 2010) but, for many Spike-Timing-Dependent Plasticity (STDP) rules, access to *postsynaptic* spikes is also required. GeNN supports such rules by automatically generating a lookup table structure (see Knight and Nowotny (2018)). While this process could be adapted to generate a lookup table from procedural connectivity, this would further erode memory savings. However, typically not all synapses in a simulation are plastic and those that are not could be simulated fully procedurally.

In this work, we have discussed the idea of procedural connectivity in the context of GPU hardware but, we believe that there is also potential for developing new types of neuromorphic hardware built from the ground up for procedural connectivity. Key components such as the random number generator could be implemented directly in hardware leading to truly game-changing compute time improvements.

4. METHODS

In all experiments presented in this work, neurons are modelled as leaky integrate-and-fire (LIF) units with the parameters listed in Table 1. The membrane voltage V_i of neuron i is modelled as

$$(1) \quad \tau_m \frac{dV_i}{dt} = (V_i - V_{\text{rest}}) + R_m (I_{\text{syn}_i} + I_{\text{ext}_i}),$$

where τ_m and R_m represent the time constant and resistance of the neuron’s cell membrane, V_{rest} defines the resting potential, I_{syn_i} represents the synaptic input current and I_{ext_i} represents an external input current. When the membrane voltage crosses a threshold V_{th} a spike is emitted, the membrane voltage is reset to V_{rest} and updating of V is suspended for a refractory period τ_{ref} . In the models where there are synaptic connections, pre-synaptic spikes lead to exponentially-decaying input currents I_{syn_i}

$$(2) \quad \tau_{\text{syn}} \frac{dI_{\text{syn}_i}}{dt} = -I_{\text{syn}_i} + \sum_{i=0}^n w_{ij} \sum_{t_j} \delta(t - t_j),$$

where τ_{syn} represents the decay time constant and t_j are the arrival times of incoming spikes from n presynaptic neurons. The continuous terms of the Eq. 1 and 2 are separately solved algebraically so that the input currents I_{syn_i} and I_{ext_i} are treated as constants throughout each simulation timestep.

5. ACKNOWLEDGEMENTS

We would like to thank Jari Pronold, Sacha van Albada, Agnes Korcsak-Gorzo and Maximilian Schmidt for their help with the multi-area model data; and Dan Goodman and Mantas Mikaitis for their feedback on the manuscript. This work was funded by the EPSRC (Brains on Board project, grant number EP/P006094/1).

6. AUTHOR CONTRIBUTIONS

J.K. and T.N. wrote the paper. T.N. is the original developer of GeNN. J.K. is currently the primary GeNN developer and was responsible for extending the code generation approach to the procedural simulation of synaptic connectivity. J.K. performed the experiments and the analysis of the results that are presented in this work.

7. DATA AVAILABILITY

The raw data used to produce Fig. 1 and Fig. 2; and the pre-processed data used to produce Fig. 3 are available at https://github.com/BrainsOnBoard/procedural_paper. The raw spiking data from the multi-area model simulations is available upon request.

8. CODE AVAILABILITY

All experiments were carried out using the GeNN 4.3.0, available at <https://github.com/genn-team/genn/releases/tag/4.3.0>. A GeNN port of the multi-area model is available at <https://github.com/neworderofjamie/multi-area-model>. The models used to produce Fig. 1 and Fig. 2 are all available at https://github.com/BrainsOnBoard/procedural_paper.

REFERENCES

- Bakker, R., Wachtler, T., and Diesmann, M. (2012). CoCoMac 2.0 and the future of tract-tracing databases. *Frontiers in Neuroinformatics*, 6(DEC):1–6.
- Belitski, A., Gretton, A., Magri, C., Murayama, Y., Montemurro, M. A., Logothetis, N. K., and Panzeri, S. (2008). Low-frequency local field potentials and spikes in primary visual cortex convey independent visual information. *Journal of Neuroscience*, 28(22):5696–5709.
- Binzegger, T., Douglas, R. J., and Martin, K. A. (2004). A quantitative map of the circuit of cat primary visual cortex. *Journal of Neuroscience*, 24(39):8441–8453.
- Blundell, I., Brette, R., Cleland, T. A., Close, T. G., Coca, D., Davison, A. P., Diaz-Pier, S., Fernandez Musoles, C., Gleeson, P., Goodman, D. F. M., Hines, M., Hopkins, M. W., Kumbhar, P., Lester, D. R., Marin, B., Morrison, A., Müller, E., Nowotny, T., Peyser, A., Plotnikov, D., Richmond, P., Rowley, A., Rumpe, B., Stimberg, M., Stokes, A. B., Tomkins, A., Trensch, G., Woodman, M., and Eppler, J. M. (2018). Code Generation in Computational Neuroscience: A Review of Tools and Techniques. *Frontiers in Neuroinformatics*, 12(November).
- Brader, J. M., Senn, W., and Fusi, S. (2007). Learning real-world stimuli in a neural network with spike-driven synaptic dynamics. *Neural computation*, 19(11):2881–912.
- Brette, R., Rudolph, M., Carnevale, T., Hines, M., Beeman, D., Bower, J. M., Diesmann, M., Morrison, A., Goodman, P. H., Harris, F. C., Zirpe, M., Natschläger, T., Pecevski, D., Ermentrout, B., Djurfeldt, M., Lansner, A., Rochel, O., Vieville, T., Muller, E., Davison, A. P., El Boustani, S., and Destexhe, A. (2007). Simulation of networks of spiking neurons: a review of tools and strategies. *Journal of computational neuroscience*, 23(3):349–98.
- Cabral, J., Kringelbach, M. L., and Deco, G. (2014). Exploring the network dynamics underlying brain activity during rest. *Progress in Neurobiology*, 114:102–131.
- Carnevale, N. T. and Hines, M. L. (2006). *The NEURON book*. Cambridge University Press.
- Clopath, C., Büsing, L., Vasilaki, E., and Gerstner, W. (2010). Connectivity reflects coding: a model of voltage-based STDP with homeostasis. *Nature neuroscience*, 13(December 2009):344–352.

- Devroye, L. (2013). *Non-uniform random variate generation*. Springer-Verlag New York, New York.
- Ercsey-Ravasz, M., Markov, N. T., Lamy, C., VanEssen, D. C., Knoblauch, K., Toroczka, Z., and Kennedy, H. (2013). A Predictive Network Model of Cerebral Cortical Connectivity Based on a Distance Rule. *Neuron*, 80(1):184–197.
- Freedman, D. and Diaconis, P. (1981). On the histogram as a density estimator: L 2 theory. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 57(4):453–476.
- Frenkel, C., Lefebvre, M., Legat, J.-D., and Bol, D. (2018). A 0.086-mm² 12.7-pJ/SOP 64k-Synapse 256-Neuron Online-Learning Digital Spiking Neuromorphic Processor in 28nm CMOS. *IEEE Transactions on Biomedical Circuits and Systems*, PP(XX):1–1.
- Furber, S. B., Galluppi, F., Temple, S., and Plana, L. A. (2014). The SpiNNaker Project. *Proceedings of the IEEE*, 102(5):652–665.
- Gewaltig, M.-O. and Diesmann, M. (2007). NEST (NEural Simulation Tool). *Scholarpedia*, 2(4):1430.
- Herculano-Houzel, S., Mota, B., and Lent, R. (2006). Cellular scaling rules for rodent brains. *Proceedings of the National Academy of Sciences*, 103(32):12138–12143.
- Izhikevich, E. M. (2005). Large-Scale Simulation of the Human Brain.
- Izhikevich, E. M. and Edelman, G. M. (2008). Large-scale model of mammalian thalamocortical systems. *Proceedings of the National Academy of Sciences of the United States of America*, 105(9):3593–8.
- Jordan, J., Ippen, T., Helias, M., Kitayama, I., Sato, M., Igarashi, J., Diesmann, M., and Kunkel, S. (2018). Extremely Scalable Spiking Neuronal Network Simulation Code: From Laptops to Exascale Computers. *Frontiers in Neuroinformatics*, 12(February):2.
- 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(December):1–19.
- Markov, N. T., Ercsey-Ravasz, M. M., Ribeiro Gomes, A. R., Lamy, C., Magrou, L., Vezoli, J., Misery, P., Falchier, A., Quilodran, R., Gariel, M. A., Sallet, J., Gamanut, R., Huissoud, C., Clavagnier, S., Giroud, P., Sappey-Mariniere, D., Barone, P., Dehay, C., Toroczka, Z., Knoblauch, K., Van Essen, D. C., and Kennedy, H. (2014a). A weighted and directed interareal connectivity matrix for macaque cerebral cortex. *Cerebral Cortex*, 24(1):17–36.
- Markov, N. T., Vezoli, J., Chameau, P., Falchier, A., Quilodran, R., Huissoud, C., Lamy, C., Misery, P., Giroud, P., Ullman, S., Barone, P., Dehay, C., Knoblauch, K., and Kennedy, H. (2014b). Anatomy of hierarchy: Feedforward and feedback pathways in macaque visual cortex. *Journal of Comparative Neurology*, 522(1):225–259.
- Merolla, P. A., Arthur, J. V., Alvarez-Icaza, R., Cassidy, A. S., Sawada, J., Akopyan, F., Jackson, B. L., Imam, N., Guo, C., Nakamura, Y., Brezzo, B., Vo, I., Esser, S. K., Appuswamy, R., Taba, B., Amir, A., Flickner, M. D., Risk, W. P., Manohar, R., and Modha, D. S. (2014). A million spiking-neuron integrated circuit with a scalable communication network and interface. *Science*, 345(6197):668–673.
- NVIDIA Corporation (2019). *CUDA C++ Programming Guide*.
- NVIDIA Corporation (2020a). *Nsight Compute*.
- NVIDIA Corporation (2020b). *NVLink Fabric Multi-GPU Processing*.
- Plotnikov, D., Blundell, I., Ippen, T., Eppler, J. M., Rumpe, B., and Morrison, A. (2016). NESTML: a modeling language for spiking neurons. In *Lecture Notes in Informatics (LNI)*, volume P-254, pages 93–108. Modellierung 2016, Karlsruhe (Germany), 17 Mar 2016 - 19 Mar 2016, Gesellschaft für Informatik e.V. (GI).
- 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(3):785–806.
- Qiao, N., Mostafa, H., Corradi, F., Osswald, M., Stefanini, F., Sumislawska, D., and Indiveri, G. (2015). A reconfigurable on-line learning spiking neuromorphic processor comprising 256 neurons and 128K synapses. *Frontiers in Neuroscience*, 9(APR):1–17.
- Rhodes, O., Peres, L., Rowley, A. G. D., Gait, A., Plana, L. A., Brenninkmeijer, C., and Furber, S. B. (2020). Real-time cortical simulation on neuromorphic hardware. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering*

- Sciences*, 378(2164):20190160.
- Salmon, J. K., Moraes, M. A., Dror, R. O., and Shaw, D. E. (2011). Parallel random numbers: As Easy as 1, 2, 3. In *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis on - SC '11*, volume 81, page 1, New York, New York, USA. ACM Press.
- Schemmel, J., Kriener, L., Muller, P., and Meier, K. (2017). An accelerated analog neuromorphic hardware system emulating NMDA- and calcium-based non-linear dendrites. *Proceedings of the International Joint Conference on Neural Networks*, 2017-May:2217–2226.
- Schmidt, M., Bakker, R., Hilgetag, C. C., Diesmann, M., and van Albada, S. J. (2018a). Multi-scale account of the network structure of macaque visual cortex. *Brain Structure and Function*, 223(3):1409–1435.
- Schmidt, M., Bakker, R., Shen, K., Bezgin, G., Diesmann, M., and van Albada, S. J. (2018b). A multi-scale layer-resolved spiking network model of resting-state dynamics in macaque visual cortical areas. *PLoS Computational Biology*, 14(10):1–38.
- Shinomoto, S., Kim, H., Shimokawa, T., Matsuno, N., Funahashi, S., Shima, K., Fujita, I., Tamura, H., Doi, T., Kawano, K., Inaba, N., Fukushima, K., Kurkin, S., Kurata, K., Taira, M., Tsutsui, K. I., Komatsu, H., Ogawa, T., Koida, K., Tanji, J., and Toyama, K. (2009). Relating neuronal firing patterns to functional differentiation of cerebral cortex. *PLoS Computational Biology*, 5(7).
- Stimberg, M., Brette, R., and Goodman, D. F. (2019). Brian 2, an intuitive and efficient neural simulator. *eLife*, 8:1–41.
- van Albada, S. J., Helias, M., and Diesmann, M. (2015). Scalability of Asynchronous Networks Is Limited by One-to-One Mapping between Effective Connectivity and Correlations. *PLoS Computational Biology*, 11(9):1–37.
- van Albada, S. J., Pronold, J., van Meegen, A., and Diesmann, M. (in press). Usage and Scaling of an Open-Source Spiking Multi-Area Model of Monkey Cortex. In Amunts, K., Grandinetti, L., Lippert, T., and Petkov, N., editors, *Brain-Inspired Computing*. Springer.
- Vogels, T. P. and Abbott, L. F. (2005). Signal Propagation and Logic Gating in Networks of Integrate-and-Fire Neurons. *The Journal of Neuroscience*, 25(46):10786–10795.
- Wang, G., Lin, Y. S., and Yi, W. (2010). Kernel fusion: An effective method for better power efficiency on multithreaded GPU. *Proceedings - 2010 IEEE/ACM International Conference on Green Computing and Communications, GreenCom 2010, 2010 IEEE/ACM International Conference on Cyber, Physical and Social Computing, CPSCom 2010*, pages 344–350.
- Yavuz, E., Turner, J., and Nowotny, T. (2016). GeNN: a code generation framework for accelerated brain simulations. *Scientific reports*, 6(November 2015):18854.