**frontiers**

# CloudBrain: Online neural computation in the cloud

**Leon Bonde Larsen** [1,*], **Rasmus Karnøe Stagsted** [1], **Beck Strohmer** [1] and **Anders Lyhne Christensen** [1]

[1]*SDU Biorobotics, Maersk McKinney Moller Institute, University of Southern Denmark*

Correspondence*:
Leon Bonde Larsen
lelar@mmmi.sdu.dk

## 1 ABSTRACT

Neuromorphic computing currently relies heavily on complicated hardware design to implement asynchronous, parallel and very large-scale brain simulations. This dependency slows down the migration of biological insights into technology. It typically takes several years from idea to finished hardware and once developed the hardware is not broadly available to the community. In this contribution, we present the CloudBrain research platform, an alternative based on modern cloud computing and event stream processing technology. Typical neuromorphic design goals, such as small form factor and low power consumption, are traded for 1) no constraints on the model elements, 2) access to all events and parameters during and after the simulation, 3) online reconfiguration of the network, and 4) real-time simulation. We explain principles for how neuron, synapse and network models can be implemented and we demonstrate that our implementation can be used to control a physical robot in real-time. CloudBrain is open source and can run on commodity hardware or in the cloud, thus providing the community a new platform with a different set of features supporting research into, for example, neuron models, structural plasticity and three-factor learning.

Keywords: Neuromorphic, Cloud, Robotics, Bio-inspired, Event-based, Stream-processing, Structural plasticity

## 1 INTRODUCTION

In traditional artificial neural networks (ANNs), the activation of neurons is represented as a scalar and information is propagated through the network in discrete steps. While this model can be computed efficiently on standard CPUs and GPUs, it is a very simplistic abstraction of biological neural networks. Biological neurons primarily communicate asynchronously through action potentials or *spikes* (Sterling and Laughlin, 2015). The timing of spikes can encode crucial information, for example in the auditory system where temporal information is used to infer direction of a sound source (Carr and Konishi, 1990; Haessig et al., 2020). Spiking neural networks (SNN) are a class of ANNs in which the temporal aspects of inter-neuron communication are explicitly considered: neurons asynchronously produce and communicate via discrete events (spikes), and SNNs thus allow encoding of information in the timing of the events.

Different spiking models have been described in literature ranging from the relatively simple integrate-and-fire model (Keat et al., 2001; Jolivet et al., 2004; Paninski et al., 2004) to the more complex Hodgkin-Huxley model (Hodgkin and Huxley, 1952). To improve biological fidelity there is, however, still a need for experimenting with new models. For example biological neurons can be non-spiking (Sterling and Laughlin,

30  2015) and there could be advantages of combining spiking and non-spiking models (Woźniak et al., 2020).
31  Integrating non-spiking neurons in SNNs can be a biologically plausible way to interface analogue sensors
32  and can provide more control over network behaviour (Strohmer et al., 2020).

33  Current implementations of learning, both in ANNs and SNNs, are based almost exclusively on adapting
34  parameters in the synapses connecting the neurons. Such synaptic learning also plays a crucial role in
35  biological learning, but in ensemble with for example structural adaptations of the network and influence
36  from different neuromodulators providing reward signals or adapting neuron behaviour (Sterling and
37  Laughlin, 2015; Price et al., 2017). In a neuromorphic engineering context, structural plasticity has been
38  shown to improve facilitation of the hardware (Qi et al., 2018), increase success of learning a task in
39  reinforcement learning (Spüler et al., 2015), and in unsupervised classification tasks (Roy and Basu,
40  2017). Online adaptation of network structure is, however, not directly supported in current neuromorphic
41  hardware, thus restricting research to pure simulations.

42  In this paper, we present the CloudBrain platform for simulating SNNs. CloudBrain utilises modern
43  cloud technology to create an infrastructure capable of executing SNNs in a computer cluster. This gives
44  several advantages: 1) There are practically no constraints on the model elements. If the concept can be
45  described in code it will also run in the cluster. 2) It allows access to all events and parameters both online
46  and offline, making it easier to monitor, develop and test solutions. 3) The structure of the network can be
47  reconfigured online allowing model elements to affect connectivity. 4) The network can run online and
48  control a robot through the *robot as a service* principle (Kuffner, 2010) to interact with the environment. 5)
49  It runs on standard computers and is based on well-documented, field-tested, free, open source software.
50  We present the architecture of CloudBrain, demonstrate the advantages of the approach, and deploy it in
51  closed-loop control of a robot.

## 1.1   SNN simulators

53  SNNs can be simulated on PCs or supercomputers using specialised software such as the GENESYS
54  (Bower et al., 2003) and NEURON (Carnevale and Hines, 2006) simulators or the more computationally
55  tractable NEST (Gewaltig and Diesmann, 2007), BRIAN (Stimberg et al., 2019) and CARLsim (Chou
56  et al., 2018). These simulators are very useful for investigating the behaviour of networks and of their
57  constituent parts. Their limitation lies in not being able to embody the neural simulation for instance to
58  control a physical robot. It has also been suggested that they are less suited for evolving models (Nowke
59  et al., 2018) because the model requires external control while evolving.

60  An alternative to the software simulations is neuromorphic hardware. Since a spiking neuron only needs
61  to do work whenever it receives a spike, it can operate asynchronously. That observation has been the
62  basis for developing non-von-Neumann computer chips (Furber, 2016) leading to small, scalable, fast and
63  energy-efficient devices for researching and deploying SNNs, such as SpiNNaker (Furber et al., 2014) and
64  BrainScales (Schemmel et al., 2010) developed in the Human Brain Project, IBM's True North (Essera
65  et al., 2016), Loihi from Intel (Davies et al., 2018), and the analog DYNAPs (Qiao et al., 2015; Moradi et al.,
66  2018) developed at ETH Zurich. Such neuromorphic hardware is well suited for embodied experiments
67  since the neuronal computations can run in real time, for example controlling a physical robot.

68  Each chip represents a trade-off between features and limitations. Common for all of them is that the
69  interface to and from the chip is a bottleneck and does not allow the user to export all the spike events
70  happening in the chip. This can complicate monitoring and makes it harder to analyse a network. The
71  SpiNNaker platform (Furber et al., 2014) is available for loan but otherwise gaining access to neuromorphic
72  hardware can be challenging because only few units exist or because intellectual property rights restrict its

73  use. Support in the form of software frameworks and documentation can also be limited as is support for
74  special features, for example to investigate new neuron models.

75    More flexible hardware implementations of SNNs have been demonstrated on Field Programmable Gate
76  Arrays (FPGAs) and Graphics Processing Units (GPUs). FPGAs are programmable devices consisting of
77  numerous logic blocks that can be almost arbitrarily connected. Once programmed, the FPGA's performance
78  is comparable to specialised chips. The use of FPGAs to simulate SNNs has been found to be highly
79  scalable (Moore et al., 2012; Wang and van Schaik, 2018) and some work suggests vector processing
80  implemented in FPGAs can help mitigate the memory bottleneck problem that reduces access to spikes and
81  parameters (Naylor et al., 2013). However, the lack of hardware support for floating-point arithmetic limits
82  FPGAs to the simpler neuron models and they do not easily support online changes to network structure.

83    FPGAs are still quite uncommon and programming them is very different from computer programming
84  so their use requires special training. GPUs, on the other hand, are common in most modern PCs and
85  implementations of SNNs on GPUs have been demonstrated to be highly scalable (Hoang et al., 2013;
86  Chou et al., 2018). GeNN, a GPU-enhanced simulation software based on NVIDIA CUDA technology,
87  even out-performed some state-of-the-art specialised chips with regards to speed and power-consumption
88  (Knight and Nowotny, 2018). The availability of embedded GPU platforms, such as NVIDIA's Jetson TX2
89  also enables GeNN to be used interactively to control a robot. The strong commercial development of
90  GPUs is constantly moving the boundaries for what is possible but generally, moving data to and from the
91  GPU memory is a bottleneck limiting access to spikes and parameters.

92    Sometimes hybrid systems can enable new features or remove limitations. For example the SpiNNaker
93  million-core machine is available through the Human Brain Project's portal (Human Brain Project, 2017)
94  for running even very large simulations and Intel Labs developed a cloud-based platform for research
95  community access to scalable Loihi-based infrastructure (Intel, 2019). However, none of them support
96  online experimentation, for example with robots. Brian2GeNN (Stimberg et al., 2020) is a software
97  package that uses GeNN to accelerate simulations defined in Brian on GPU hardware. GPU acceleration of
98  simulation software has also been demonstrated to improve performance on supercomputers and enable
99  larger simulations on single computers (Hoang et al., 2013; Chou et al., 2018).

## 1.2 Cloud infrastructure

101    In recent years, cloud-based technology has seen rapid development driven by the demand for distributed
102  and highly scalable IT-solutions (**?**). When executed in the cloud, a computer program often runs in a
103  virtual environment called a *container*. Seen from the program the container is like any computer with
104  resources such as CPU, memory, disk and an operating system while in fact the containers share these
105  resources. Asynchronous, event-based architectures in particular have excelled in order to handle millions
106  of social media users or e-commerce transactions. Programs are often asynchronous, meaning that they are
107  waiting for input for example from a user requesting a website or completing a purchase. While waiting the
108  program needs no computing resources and thus other programs in other containers can use the hardware.
109  This fits well with the SNN model, where the neuron only does work when an input event is present.

110    Modern IT solutions generate a lot of data and it is common to handle it as event-streams (**?**). Event-
111  streams are ordered in topics such that nodes subscribing to a topic receive events that are published
112  on that topic. There can be many publishers and many subscribers to a topic and there can be many
113  topics. Copying and distribution of the events are handled by highly optimised and extremely scalable
114  infrastructure software making it easy to interface programs with the event-stream. Such an architecture is
115  well suited for handling spike events since many neurons need to receive the same events.

## 2 ARCHITECTURE

116 An SNN simulation in CloudBrain is essentially a collection of small programs implementing mathematical
117 models and communicating the resulting events as they happen. CloudBrain is the platform that runs the
118 programs in a scalable and modular way. In CloudBrain, both neurons, synapses and any other models
119 are user-defined programs, while spikes, parameters, and any other messages are events consisting of a
120 timestamp and an arbitrary payload.

### 2.1 Programs

122 The *NeuronProgram* and *SynapsePrograms* run under the *ControlProgram* in order to hide the complexity
123 of communication and OS-specific interfaces. One *ControlProgram* can run one *NeuronProgram* and
124 multiple *SynapsePrograms* and to ensure scalability, the *ControlProgram* is executed in a container. Thus
125 the *ControlProgram* can run on any host within the cluster and multiple *ControlPrograms* can run on the
126 same host (figure 1). The user controls if the *ControlProgram* runs synchronously updating the neuron
127 model at a specific rate or asynchronously only updating the model when a spike is received. The same
128 goes for the *SynapseProgram* which is responsible for keeping information about the connections (for
129 example weight and delay) and attaching it to the payload of received spikes before handing them to the
130 *NeuronProgram*. The connection information is provided when the connection is first made but can also
131 be updated during execution. *Events* are implemented as asynchronous messages transmitted following
132 the publish-subscribe pattern (Birman and Joseph, 1987) so a *ControlProgram* receives messages only
133 from topics it has subscribed to. Topics contain either *ControlEvents*, handled by the *ControlProgram*
134 or *NeuroEvents*, passed first through a *SynapseProgram* and then handled by a *NeuronProgram*. Each
135 *ControlProgram* has an individual *ControlTopic* that it always subscribes to while all other subscriptions
136 are set up at run-time.

### 2.2 Global control

138 A *GlobalController* is responsible for scaling the number of containers in the cluster and for configuring
139 the *ControlPrograms* by emitting *ControlEvents*. Each *ControlProgram* in the cluster has a unique ID
140 known to itself and the *GlobalController*. To set up an experiment, the user writes a *GlobalController*
141 program that tells each of the *ControlPrograms* which *NeuronProgram* and *SynapsePrograms* to run, which
142 parameters to use and how to connect. It also defines any groupings of *NeuronPrograms*, for example
143 into populations. If non-standard neuron models are used, they have to be implemented in code and either
144 provisioned to the cluster before running the *GlobalController* or sent to the individual container using
145 control events, during execution. The *GlobalController* can be run from any computer on the same network
146 as the cluster.

### 2.3 Simulation method

148 Ideally, the *NeuronPrograms* and *SynapsePrograms* should be asynchronous, meaning that they only
149 use processing resources when they receive or transmit an *Event*. This greatly improves performance but
150 is not practical for all neuron models and thus they can also be periodic. To run asynchronously, every
151 time a neuron receives a spike, it must predict when it will spike based on the mathematical model of the
152 neuron and its internal parameters. The neuron then registers a timer to wake it at that time and removes
153 any previously registered timers (figure 1).

154 Communication in the cluster is faster than in biological neurons. Because the *NeuronPrograms* calculate
155 behaviour based on timestamps and not on the actual time of arrival, the cluster essentially spends the time

156  accounted for in the biological delays to complete the required calculations and communicate the results. If
157  that time is insufficient, the post-synaptic neuron will know that a deadline was missed and can report it.
158  The *SynapseProgram* receives incoming *SpikeEvents* and attaches the connection information (typically
159  weight and delay) before passing them on to the *NeuronProgram*. When the *NeuronProgram* emits a spike,
160  it also alerts the *SynapseProgram* thus allowing it to update the connection parameters according to its
161  learning rule.

## 3 METHODS

162  A proof-of-concept implementation based on Python and open-source software was built in order to validate
163  the proposed architecture. The code is available on a git repository along with a demo of CloudBrain
164  running on a single PC and instructions on how to set up a cluster with a cloud provider. All is available at
165  sdu.dk/cloudbrain.

### 3.1 Implementation

167  *ControlProgram* and *NeuronProgram* are written in Python. A custom *NeuronProgram* inherits from
168  a base class so the user only needs to override the functions used and need not care about the inner
169  workings of the *ControlProgram*. The functionality exposed by the *NeuronProgram* base class include
170  methods to emit events and to register callback functions for event reception or the expiration of timers.
171  The Python code is integrated in a minimalistic docker image (**?**) and executed in docker swarm (**?**). Using
172  the continuous integration tools included in gitlab (**?**), the process of deploying the code on the docker
173  swarm is automated. Log activity from the containers is collected using Filebeat (**?**), allowing the logs to
174  be searched and viewed from a web-based interface. Events are handled by Apache Kafka (**?**), an open-
175  source stream-processing software platform, providing high-throughput and low-latency communication.
176  Messages are JSON encoded and consist of a timestamp, sender ID and an arbitrary payload. Since Kafka
177  is agnostic to the payload, it can be seamlessly changed to fit any computational model, to set parameters in
178  a *NeuronProgram* or to retrieve arbitrary information. One of the great advantages of an event-based system
179  is the availability of tools to view, search and aggregate data. We use Elasticsearch and the visualisation
180  dash-board Kibana. This allows for fast, online and virtually unconstrained visualisation of activity in the
181  network, for example, to monitor spiking rates at population level or for individual neurons, to analyse
182  behavioural patterns or to monitor the flow of *ControlEvents*.

### 3.2 Hardware

184  The on-premises cluster consists of 15 PCs, each with 2 Intel Xeon 2.55GHz cores, 8 GB RAM, Gigabit
185  network and with Debian9 installed on a solid state drive (figure 5 shows a photo of the cluster). An
186  additional PC with Intel Xeon 2.67GHz quad-core, 12 GB RAM and SSD is used to run Kibana, Kafka,
187  Elasticsearch connector and Zookeeper (used by kafka). Another PC with Intel I5 3.30 GHz quad-core, 12
188  GB RAM and SSD is used to run Elasticsearch. Elasticsearch runs on a separate host to keep peak CPU
189  and memory usage from interfering with the performance of Kafka. Finally a PC with Intel I5 2.90 GHz
190  quad-core, 16 GB RAM and SSD was used to monitor the cluster using Grafana and InfluxDB. On each of
191  the hosts in the cluster, Telegraf was installed to collect information about the utilisation of the nodes. The
192  machines are connected to a gigabit managed switch using Cat5e cables and communication to the robot is
193  provided by a VPN tunnel through a Wi-Fi access point. As an alternative to procuring an on-premises
194  cluster, we repeated the experiments with Google Cloud Platform (GCP) executing CloudBrain. 14 virtual
195  machines were configured each with 4 vCPUs running at 2 GHz, 15 GB RAM and 100GB disk. 10 VMs

196 were used to run the neurons and the rest were used as VPN gateway, Kafka broker, Elasticsearch and
197 connectors to Elasticserch. All VMs were located in Finland while the robot was in Denmark.

## 3.3 Robot platform

199     The mobile robot described in Larsen et al. (2013) is optimised for rapid prototyping and consists of a
200 wooden board with two rear wheels connected to motors and a castor wheel in front. In this work it was
201 fitted with two custom bumper sensors, one on each side of the front (figure 5). An on-board RaspberryPi 3
202 model B connected to the cluster via Wi-Fi is responsible for generating PWM signals for the motors based
203 on received events and for emitting events based on the state of the sensors. A piece of software running in
204 the cloud translates spikes into motor messages and sensor messages into spikes. An H-bridge supplies
205 the current to the motors based on the generated PWM signals. The two sensors are implemented with
206 micro-switches and each sensor emits spikes on its own topic. When the sensor is activated, it emits spikes
207 with a frequency of 100Hz and when the sensor is not activated it emits spikes with a frequency of 10Hz.
208 Each motor has its own topic and the output is controlled by a running average of the number of spikes
209 received within the past 50ms. The average is then linearly mapped from 1-4 spikes to $\pm$ 2 wheel rotations
210 per second, controlled by a standard PID controller.

211     During experiments the robot was enclosed in a 2.5m x 1.5m environment with slanted 25cm corners (see
212 centre insert of figure 2). The environment can be configured as *corridor* or *box* by adding or removing the
213 rectangle in the centre. The experiments were repeated several times in each configuration. The controller
214 for the robot is based on the Braitenberg principle (Dennett and Braitenberg, 1986) where the left sensor has
215 an inhibitory effect on the right motor and vice-versa. The network consists of six populations, each with
216 five neurons (figure 2). Apart from the two sensor populations (A and B) and the two motor populations (C
217 and D), an excitatory bias (E) is provided to both motor populations to make the robot move, and the last
218 population (F) injects noise into the system to make the neurons fire out of phase.

## 3.4 Neuronal models

220     We implemented and tested three popular neuron models: Integrate-and-Fire (IF) (Keat et al., 2001;
221 Jolivet et al., 2004; Paninski et al., 2004), Leaky-Integrate-and-Fire (LIF) (Stein, 1967; Tuckwell, 1989) and
222 Adaptive-Exponential-Integrate-and-Fire (AEIF) (Brette and Gerstner, 2005; Gerstner and Brette, 2009).
223 The IF and LIF neurons are implemented fully asynchronously so the neurons only use computational power
224 when input spikes arrive or the neuron emits a spike. The AEIF model consists of two coupled differential
225 equations that need to be integrated over time. To keep it asynchronous would make it computationally
226 heavy because it would need to recalculate next spike time and set the timer accordingly every time a new
227 spike arrives. Instead the model was implemented synchronously using the looping function with an update
228 frequency of 1ms. We implemented the Spike-Timing-Dependent Plasticity (STDP) rule, such that it runs
229 asynchronously and updates the weight of the synapse.

## 4 RESULTS

230 To support the claims that CloudBrain provides 1) no constraints on the model elements, 2) online
231 reconfiguration of the network, 3) online operation, and 4) access to all information, we provide three
232 demonstrations. We demonstrate 1) the implementations of three popular neuron models and a synapse
233 model, 2) that the morphology of the network can be changed during operation, and 3) that a robot can be
234 controlled online by CloudBrain running on a cluster. Furthermore, to demonstrate how all information in

235 the SNN can be monitored online, we provide the live plots from Kibana. Finally, we evaluate timing and
236 load both on our own cluster and with a cloud provider.

## 4.1 Model implementations

238 To demonstrate the implemented neuron models, two different experiments were made. In the first one,
239 a spike source was connected to the IF and LIF, respectively and their voltage potential and firing was
240 observed (figure 3, top). The voltage potential of the IF neuron increases linearly until it fires, whereas
241 the voltage potential of the LIF neuron charges exponentially. The AEIF model was tested in the regular
242 bursting mode with constant current input (Naud et al., 2008) using parameters from NEST (Gewaltig and
243 Diesmann, 2007) while voltage potential, adaptation variable and spiking pattern were observed (figure 3,
244 middle). The voltage potential of the AEIF neuron depends on the adaptation variable and, as it falls to a
245 certain level, the neuron bursts. This confirms normal behaviour for all three neuron models.

246 Synaptic learning is demonstrated in a setup where a noisy spike source has excitatory connections to two
247 LIF neurons. The two LIF neurons are connected with a synapse using STDP on all pre- and post synaptic
248 events. The noisy spike source emits events with a random time difference between 100 ms and 1 s. Every
249 time one of the LIF neurons spikes, the STDP synapse will update its weight based on the time since the
250 other neuron last spiked. Figure 4 shows the evolution of synaptic weight in an STDP synapse. From the
251 zoomed-in figure, we see that when a pre-synaptic spike occurs (green vertical line) the weight increases by
252 an amount inversely proportional to the time since last post-synaptic spike. Similarly at the post-synaptic
253 spike times the weight is reduced inversely proportional to the time since last pre-synaptic spike.

## 4.2 Morphology

255 To demonstrate that the network can be reconfigured online, an experiment was made with a 10Hz spike
256 source (A) and two IF neurons (B and C). After 5 seconds, the source is connected to neuron B thus making
257 it fire. After another 10 seconds, the connection is removed and the same procedure is repeated for neuron
258 C. We plot the spiking pattern of the three neurons in figure 3. The morphology changes shown here are
259 very basic, but demonstrate that connections can be created, updated or removed by sending ControlEvents.

## 4.3 Online operation

261 To demonstrate online operation, two experiments with the robot are reported, each in a different
262 environment. The spiking rate of the populations were plotted online in Kibana and the resource use in the
263 cluster were plotted using Grafana. Here we provide screenshots from both and a video of the experiment is
264 available in the supplementary material. The experiments were run on the on-premises cluster and repeated
265 in the GCP cluster.

266 To evaluate the total delay in the robot experiments, we used videos to analyse reaction times when the
267 robot is colliding with the wall. We also estimated the delay from a change in spike rate of the sensor
268 population until the corresponding change in the opposite motor population. This represents the total delay
269 in the neuronal system without the robot. The delay of a single event through Kafka was measured by
270 subscribing to a topic and publishing 100 events on it. This was measured both over WiFi and within the
271 clusters. We report the median of the 100 measurements along with the 1st and 99th percentiles.

272 The robot was able to negotiate both corridor and box, reacting quickly and displaying the expected
273 Braitenberg behaviour. The total delay from sensor activation to motor reaction was estimated based
274 on several collisions with the wall, recorded on video. Upon collision the robot was thrown back a few
275 centimetres and the controller reaction was fast enough to prevent the robot from hitting the wall again. It

took on average 5 frames (min. 4 frames, max. 8 frames), corresponding to 166 ms (min. 133 ms, max. 266 ms) to react. There are different contributions to this delay, mainly the characteristics of the motor. The total delay from a change in a sensor population to the corresponding change in the motor population could not be determined precisely due to noise but was estimated to 30–80 ms based on the plots from Kibana (figure 5). The delay of a single event through Kafka was measured over 100 messages. The time from an event was sent until it was received was 24 ms via Wi-Fi (median: 24.0 ms, 1st: 8.0 ms, 99th: 320.4 ms) and 4.2 ms within the cluster (median: 4.2 ms, 1st: 2.7 ms, 99th: 5.5 ms).

Figure 5 also shows the load on the cluster while the robot is navigating the corridor. The memory usage is constant because once the neurons are running they do not make additional allocations. The network usage is less than 2%. For comparison, the idle cluster had memory usage of 6%, CPU usage 3% and network 0%. Under the experiment an average of 1930 spikes per second were passing through Kafka. Benchmarking of the Kafka broker showed that it could handle approximately 400.000 spikes per second.

The same experiment was conducted on GCP for a duration of 15 minutes. Similar to the delay measured with the on-premises cluster, the latency was measured from the robot to the GCP. The delay measured to 47.3 ms (median: 46.7 ms, 1st: 46.2 ms, 99th: 52.6 ms). The spike delay within the cluster was measured to 3.5 ms (median: 3.2 ms, 1st: 2.7 ms, 99th: 5.5 ms).

The *ControlProgram* allows the user to toggle live logging of internal variables and parameters. The values are sent, as events, only when they change. The plots in figure 3 and 5 are all live screenshots from Kibana demonstrating that parameters are available online. Data can also be accessed online without Kibana for more advanced plots or extracted for offline generation of plots better suited for publication.

## 5 DISCUSSION AND CONCLUSION

In this contribution, we propose a novel architecture for simulating spiking neural networks on readily available cloud infrastructure. We describe a proof-of-concept implementation, demonstrating how neuron and synapse models can be implemented and how the network connections can be updated while the system is running. Finally we demonstrate how a Braitenberg-controlled differential drive robot could be operated online from both our on-premises cluster and from a Google Cloud Platform cluster.

The numbers provided are highly dependent on hardware, application and other implementation specific factors. We therefore provide the numbers observed in our proof-of-concept implementation without expecting them to generalise perfectly to other implementations. Our on-premises cluster built from 15 refurbished PCs over ten years old was able to support the experiments in this paper and the total cost of running the robot experiment for 15 minutes on GCP was approximately 2 USD. This shows that hardware for this type of simulation is readily available.

We did experience some jitter in the delay times and occasionally an event was late by a factor of hundreds which means it missed the deadline and thus essentially was lost. In our applications it was not a problem but means of mitigation should be further investigated. The Kafka broker can handle approximately 400,000 spikes per second and this number scales linearly with the number of brokers. There are no limits in principle to the size of a computer cluster and cloud providers put huge clusters at our disposal. We simulated networks with up to 30 neurons for this paper but much larger networks should be possible and will be the focus of further work. In conclusion, we demonstrated key features that make CloudBrain especially suited for some types of experiments and argue that trading small form factor and low power consumption for such extra features can be sensible for research purposes.
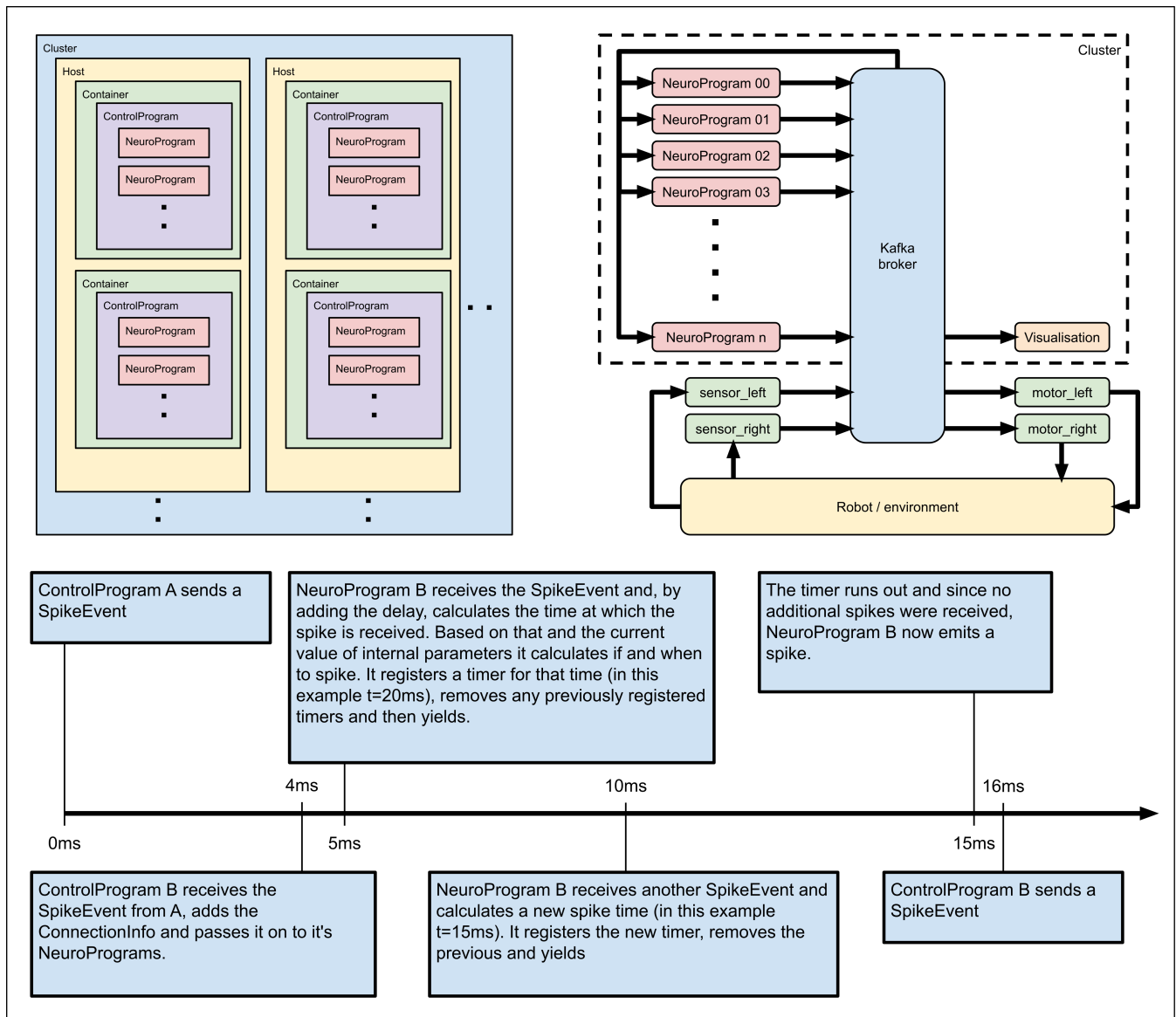
## 6 ACKNOWLEDGEMENTS

## REFERENCES

Birman, K. and Joseph, T. (1987). Exploiting virtual synchrony in distributed systems. *ACM SIGOPS Operating Systems Review* 21, 123–138. doi:10.1145/37499.37515

Bower, J. M., Beeman, D., and Hucka, M. (2003). The GENESIS simulation system. *The Handbook of Brain Theory and Neural Networks* , 475–478

Brette, R. and Gerstner, W. (2005). Adaptive exponential integrate-and-fire model as an effective description of neuronal activity. *Journal of Neurophysiology* 94, 3637–3642. doi:10.1152/jn.00686.2005

Carnevale, N. T. and Hines, M. L. (2006). *The NEURON book*. doi:10.1017/CBO9780511541612

Carr, C. E. and Konishi, M. (1990). A circuit for detection of interaural time differences in the brain stem of the barn owl. *Journal of Neuroscience* 10, 3227–3246. doi:10.1523/jneurosci.10-10-03227.1990

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 *Proceedings of the International Joint Conference on Neural Networks* (Institute of Electrical and Electronics Engineers Inc.), vol. 2018-July. doi:10.1109/IJCNN.2018. 8489326

Davies, M., Srinivasa, N., Lin, T.-H., Chinya, G., Cao, Y., Choday, S. H., et al. (2018). Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro* 38, 82–99

Dennett, D. C. and Braitenberg, V. (1986). Vehicles: Experiments in Synthetic Psychology. *The Philosophical Review* 95, 137. doi:10.2307/2185146

Essera, S. K., Merollaa, P. A., Arthura, J. V., Cassidya, A. S., Appuswamya, R., Andreopoulosa, A., et al. (2016). Convolutional networks for fast energy-efficient neuromorphic computing. *Proc. Nat. Acad. Sci. USA* 113, 11441–11446

[Dataset] Furber, S. (2016). Large-scale neuromorphic computing systems. doi:10.1088/1741-2560/13/5/ 051001

Furber, S. B., Galluppi, F., Temple, S., and Plana, L. A. (2014). The SpiNNaker project. *IEEE. Proceedings* 102, 652–665. doi:10.1109/JPROC.2014.2304638

Gerstner, W. and Brette, R. (2009). Adaptive exponential integrate-and-fire model. *Scholarpedia* 4, 8427. doi:10.4249/scholarpedia.8427

Gewaltig, M.-O. and Diesmann, M. (2007). NEST (NEural Simulation Tool). *Scholarpedia* 2, 1430. doi:10.4249/scholarpedia.1430

Haessig, G., Milde, M., Aceituno, P. V., Oubari, O., Knight, J. C., van Schaik, A., et al. (2020). Event-based computation for touch localization based on precise spike timing. *Frontiers in Neuroscience* 14, 1–29. doi:10.3389/fnins.2020.00420

Hoang, R. V., Tanna, D., Jayet Bray, L. C., Dascalu, S. M., and Harris, F. C. (2013). A novel CPU/GPU simulation environment for large-scale biologically realistic neural modeling. *Frontiers in Neuroinformatics* 7, 19. doi:10.3389/fninf.2013.00019

Hodgkin, A. L. and Huxley, A. F. (1952). A quantitative description of membrane current and its application to conduction and excitation in nerve. *The Journal of Physiology* 117, 500–544. doi:10.1113/jphysiol. 1952.sp004764
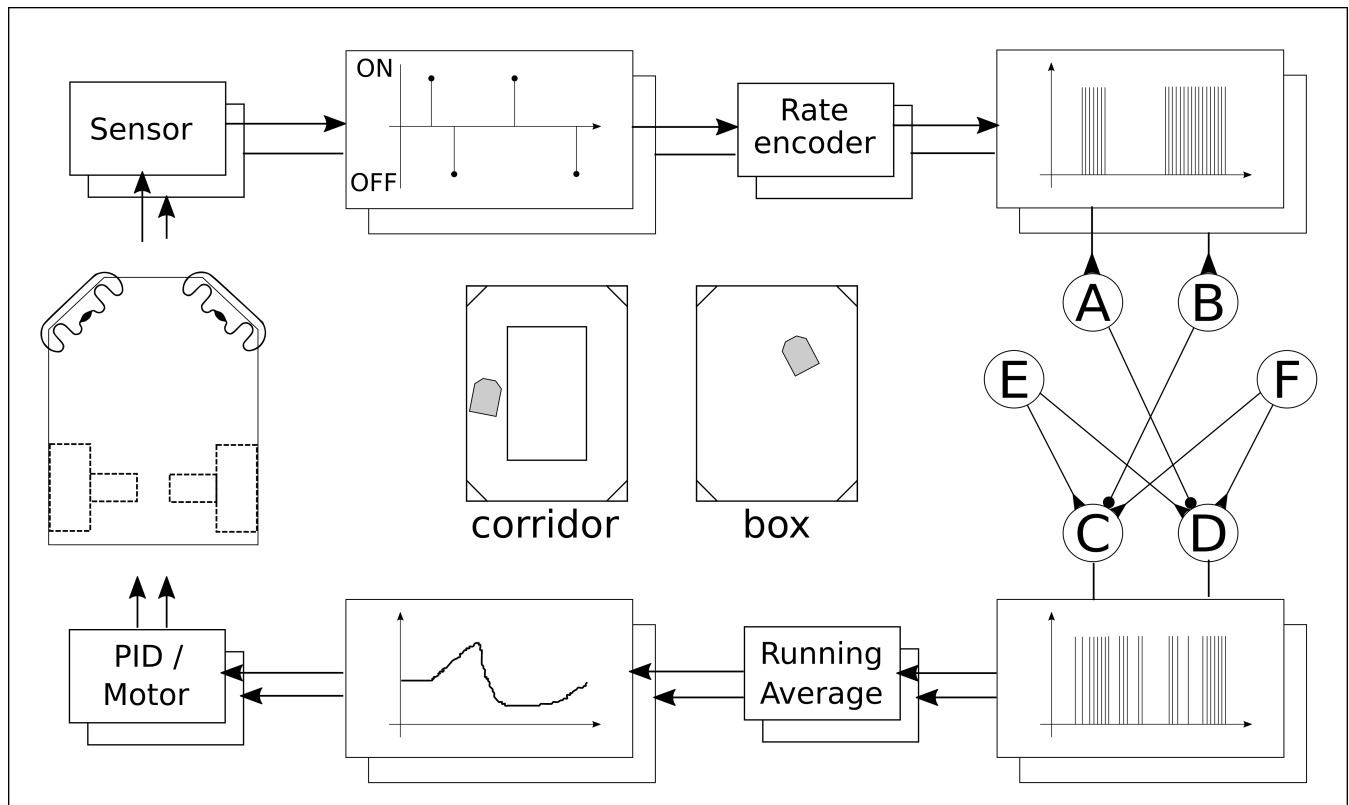
357 [Dataset] Human Brain Project (2017). Silicon Brains

358 [Dataset] Intel (2019). Beyond Today's AI

359 Jolivet, R., Lewis, T. J., and Gerstner, W. (2004). Generalized integrate-and-fire models of neuronal activity
360 approximate spike trains of a detailed model to a high degree of accuracy. *Journal of Neurophysiology*
361 92, 959–976. doi:10.1152/jn.00190.2004

362 Keat, J., Reinagel, P., Reid, R. C., and Meister, M. (2001). Predicting every spike: A model for the
363 responses of visual neurons. *Neuron* 30, 803–817. doi:10.1016/S0896-6273(01)00322-1

364 Knight, J. C. and Nowotny, T. (2018). GPUs Outperform Current HPC and Neuromorphic Solutions
365 in Terms of Speed and Energy When Simulating a Highly-Connected Cortical Model. *Frontiers in*
366 *Neuroscience* 12, 941. doi:10.3389/fnins.2018.00941

367 Kuffner, J. (2010). Cloud-enabled robots In: IEEE-RAS international conference on humanoid robots.
368 *Piscataway, NJ: IEEE*

369 Larsen, L. B., Olsen, K. S., Ahrenkiel, L., and Jensen, K. (2013). Extracurricular Activities Targeted
370 towards Increasing the Number of Engineers Working in the Field of Precision Agriculture . *XXXV*
371 *CIOSTA & CIGR V Conference* , 1–12

372 Moore, S. W., Fox, P. J., Marsh, S. J., Markettos, A. T., and Mujumdar, A. (2012). Bluehive - A field-
373 programable custom computing machine for extreme-scale real-time neural network simulation. In
374 *Proceedings of the 2012 IEEE 20th International Symposium on Field-Programmable Custom Computing*
375 *Machines, FCCM 2012*. 133–140. doi:10.1109/FCCM.2012.32

376 Moradi, S., Qiao, N., Stefanini, F., and Indiveri, G. (2018). A Scalable Multicore Architecture with
377 Heterogeneous Memory Structures for Dynamic Neuromorphic Asynchronous Processors (DYNAPs).
378 *IEEE Transactions on Biomedical Circuits and Systems* 12, 106–122. doi:10.1109/TBCAS.2017.
379 2759700

380 Naud, R., Marcille, N., Clopath, C., and Gerstner, W. (2008). Firing patterns in the adaptive exponential
381 integrate-and-fire model. *Biological Cybernetics* 99, 335–347. doi:10.1007/s00422-008-0264-7

382 Naylor, M., Fox, P. J., Markettos, A. T., and Moore, S. W. (2013). Managing the FPGA memory wall:
383 Custom computing or vector processing? In *2013 23rd International Conference on Field Programmable*
384 *Logic and Applications, FPL 2013 - Proceedings* (IEEE Computer Society). doi:10.1109/FPL.2013.
385 6645538

386 Nowke, C., Diaz-Pier, S., Weyers, B., Hentschel, B., Morrison, A., Kuhlen, T. W., et al. (2018). Toward
387 rigorous parameterization of underconstrained neural network models through interactive visualization
388 and steering of connectivity generation. *Frontiers in Neuroinformatics* 12, 32. doi:10.3389/fninf.2018.
389 00032

390 [Dataset] Paninski, L., Pillow, J. W., and Simoncelli, E. P. (2004). Maximum likelihood estimation of a
391 stochastic integrate-and-fire neural encoding model. doi:10.1162/0899766042321797

392 Price, D. J., Jarman, A. P., Mason, J. O., and Kind, P. C. (2017). *Building brains - An introduction to*
393 *neural development*. doi:10.1002/9781119293897

394 Qi, Y., Shen, J., Wang, Y., Tang, H., Yu, H., Wu, Z., et al. (2018). Jointly learning network connections
395 and link weights in spiking neural networks. In *IJCAI International Joint Conference on Artificial*
396 *Intelligence*. vol. 2018-July, 1597–1603. doi:10.24963/ijcai.2018/221

397 Qiao, N., Mostafa, H., Corradi, F., Osswald, M., Stefanini, F., Sumislawska, D., et al. (2015). A
398 reconfigurable on-line learning spiking neuromorphic processor comprising 256 neurons and 128K
399 synapses. *Frontiers in Neuroscience* 9. doi:10.3389/fnins.2015.00141

400  Roy, S. and Basu, A. (2017). An online unsupervised structural plasticity algorithm for spiking neural
401    networks. *IEEE Transactions on Neural Networks and Learning Systems* 28, 900–910. doi:10.1109/
402    TNNLS.2016.2582517
403  Schemmel, J., Briiderle, D., Griibl, A., Hock, M., Meier, K., and Millner, S. (2010). A wafer-
404    scale neuromorphic hardware system for large-scale neural modeling. In *Proceedings of 2010 IEEE
405    International Symposium on Circuits and Systems* (IEEE), 1947–1950
406  Spüler, M., Nagel, S., and Rosenstiel, W. (2015). A spiking neuronal model learning a motor control task
407    by reinforcement learning and structural synaptic plasticity. In *Proceedings of the International Joint
408    Conference on Neural Networks*. vol. 2015-Septe. doi:10.1109/IJCNN.2015.7280521
409  Stein, R. B. (1967). Some Models of Neuronal Variability. *Biophysical Journal* 7, 37–68. doi:10.1016/
410    S0006-3495(67)86574-3
411  Sterling, P. and Laughlin, S. (2015). *Principles of neural design*. doi:10.7551/mitpress/9780262028707.
412    001.0001
413  Stimberg, M., Brette, R., and Goodman, D. F. (2019). Brian 2, an intuitive and efficient neural simulator.
414    *eLife* 8. doi:10.7554/eLife.47314
415  Stimberg, M., Goodman, D. F., and Nowotny, T. (2020). Brian2GeNN: accelerating spiking neural network
416    simulations with graphics hardware. *Scientific Reports* 10, 1–12. doi:10.1038/s41598-019-54957-7
417  Strohmer, B., Manoonpong, P., and Larsen, L. B. (2020). Integrating Non-Spiking Interneurons in Spiking
418    Neural Networks. *bioRxiv* , 2020.08.13.249375doi:10.1101/2020.08.13.249375
419  Tuckwell, H. (1989). Introduction to theoretical neurobiology volume 2, nonlinear and stochastic theories.
420    *Comparative Biochemistry and Physiology Part A: Physiology* 92, 268. doi:10.1016/0300-9629(89)
421    90177-1
422  Wang, R. and van Schaik, A. (2018). Breaking Liebig's law: An advanced multipurpose neuromorphic
423    engine. *Frontiers in Neuroscience* 12, 593. doi:10.3389/fnins.2018.00593
424  Woźniak, S., Pantazi, A., Bohnstingl, T., and Eleftheriou, E. (2020). Deep learning incorporating
425    biologically inspired neural dynamics and in-memory computing. *Nature Machine Intelligence* 2,
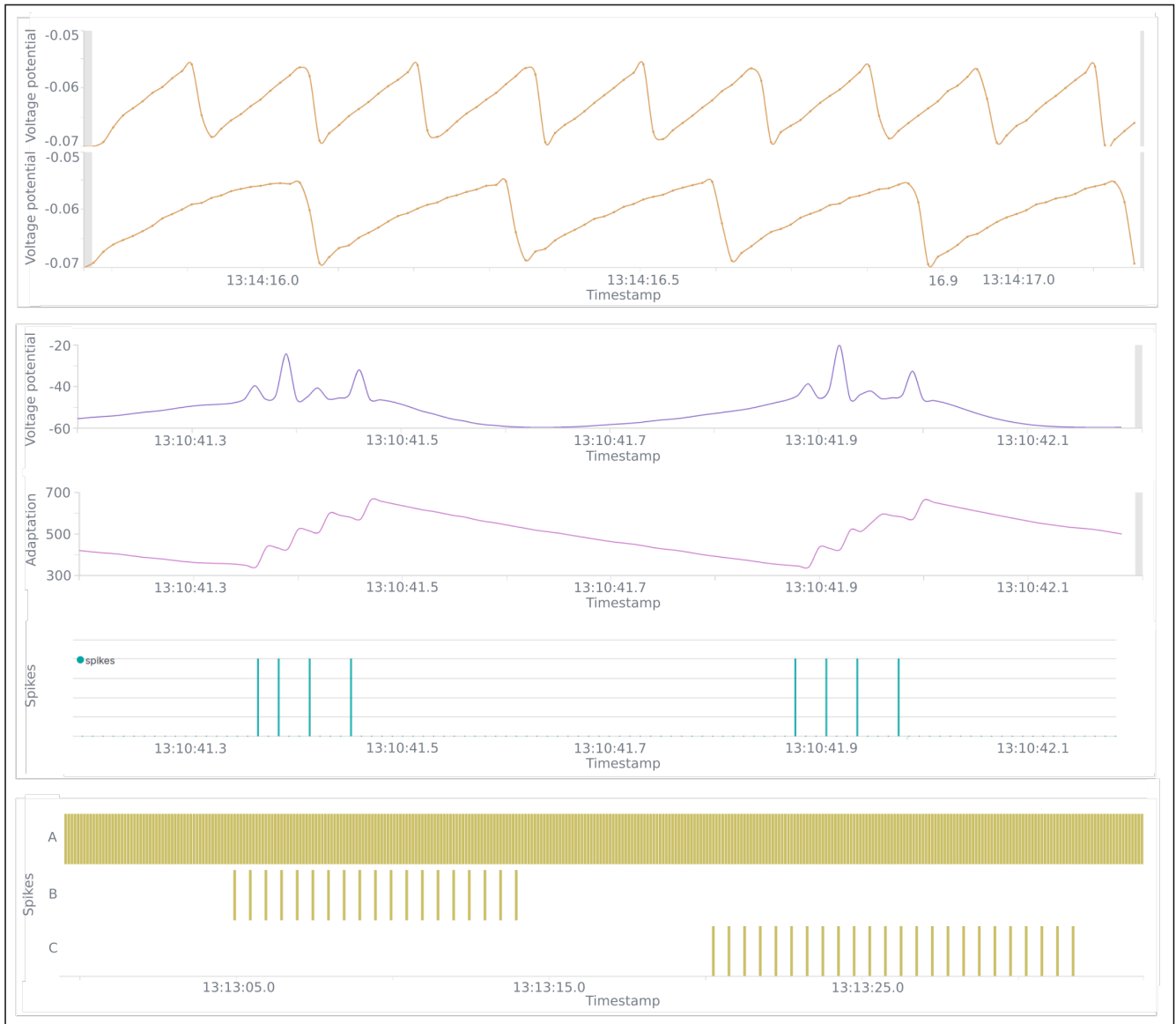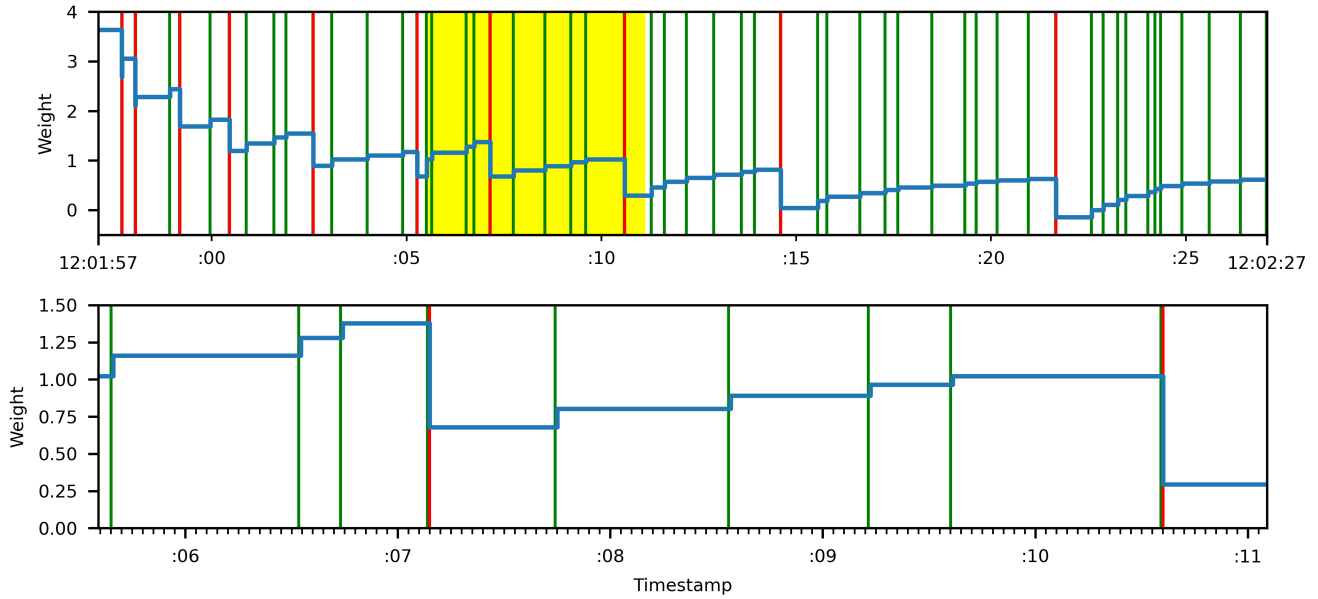426    325–336. doi:10.1038/s42256-020-0187-0

## FIGURE CAPTIONS

**Figure 1. Top left**: Shows how CloudBrain scales. Several *NeuroPrograms* can run under one *ControlProgram*. Several containers, each running one *ControlProgram*, can run on one host and the cluster is made up of any number of hosts. **Top right**: Conceptual overview of the robot experiment. The *NeuroPrograms* communicate through the kafka broker with the sensor and motor populations on the robot. **Bottom**: Simplified scenario demonstrating how asynchronous neuron models are implemented in CloudBrain.
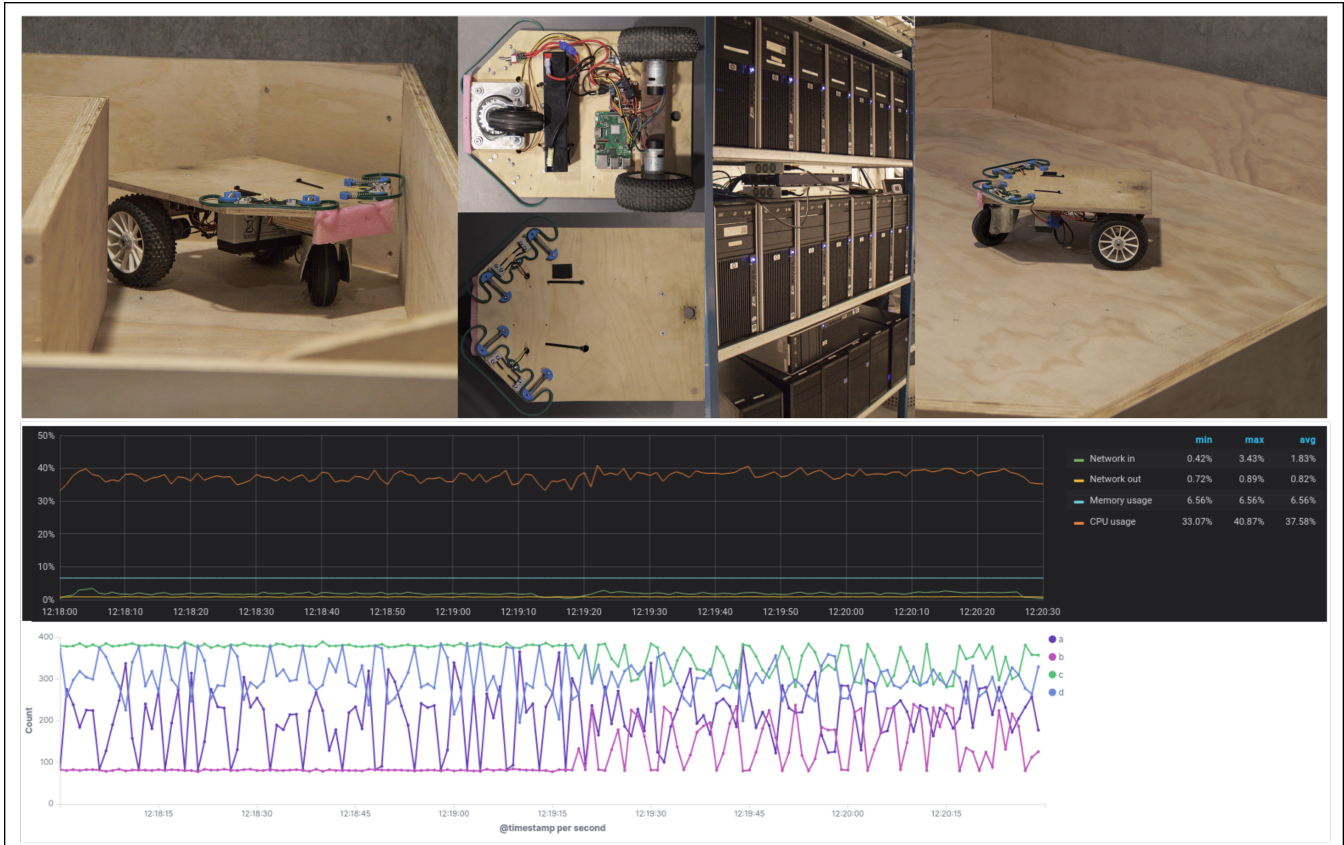
**Figure 2.** The control loop of the robot. Onset/offset events from the sensors are rate encoded and fed into populations A and B, respectively. E is injecting a bias (to keep the robot moving) and F injects noise. The motor neuron populations C and D are decoded by a running average and used as set-point in the PID for the motors. The centre insert shows the corridor track and the box track used for testing the robot.

**Figure 3.** Screenshots from live monitoring in Kibana. **Top**: Voltage potential from experiment with Integrate-and-fire and Leaky-integrate-and-fire. **Middle**: Voltage potential, adaptation variable and spikes from experiment with Adaptive-exponential-integrate-and-fire. **Bottom**: Spikes from morphology experiment showing spikes from the spike source (top line) and from the two Integrate-and-fire neurons. The spike source is in turn connected to and disconnected from the neurons, causing them to spike.

**Figure 4.** Synapse learning using STDP. Green vertical bars represent pre-synaptic spikes, while red vertical bars represent post-synaptic spikes. The blue line represents the synaptic weight. The bottom plot is a magnification of the yellow area in the top plot.

**Figure 5. Left image**: The robot in the corridor. **Right image**: The robot in the box. **Middle images**: The robot seen from below, the robot seen from above and the on-premise cluster. **Top plot**: Load on the cluster during operation. All values are taken as the percentage of the cluster's total capacity. The plot is taken directly from Grafana as displayed live. **Bottom plot**: Spike rates on the interface populations (a/b for left/right sensor and c/d for left/right motor). The first half is in the box and the second half is in the corridor. While driving in the box, the robot displays a wall following behaviour thus always activating the same sensor and reacting with the same motor. The plot is taken directly from Kibana as displayed live.