

A scalable approach to modeling on accelerated neuromorphic hardware

Eric Müller^{*†1}, Elias Arnold^{†1}, Oliver Breitwieser^{†1}, Milena Czierlinski^{†1}, Arne Emmel^{†1}, Jakob Kaiser^{†1}, Christian Mauch^{†1}, Sebastian Schmitt^{†2}, Philipp Spilger^{†1}, Raphael Stock^{†1}, Yannik Stradmann^{†1}, Johannes Weis^{†1}, Andreas Baumbach^{1,3}, Sebastian Billaudelle¹, Benjamin Cramer¹, Falk Ebert¹, Julian Göltz^{3,1}, Joscha Ilmberger¹, Vitali Karasenko¹, Mitja Kleider¹, Aron Leibfried¹, Christian Pehle¹, Johannes Schemmel¹

¹*Kirchhoff-Institute for Physics, Heidelberg University, Germany*

²*Third Institute of Physics, University of Göttingen, Germany*

³*Department of Physiology, University of Bern, Switzerland*

Correspondence*:

Eric Müller

mueller@kip.uni-heidelberg.de

1 ABSTRACT

2 Neuromorphic systems open up opportunities to enlarge the explorative space for computational
3 research. However, it is often challenging to unite efficiency and usability. This work presents
4 the software aspects of this endeavor for the BrainScaleS-2 system, a hybrid accelerated
5 neuromorphic hardware architecture based on physical modeling. We introduce key aspects of
6 the BrainScaleS-2 Operating System: experiment workflow, API layering, software design, and
7 platform operation. We present use cases to discuss and derive requirements for the software
8 and showcase the implementation. The focus lies on novel system and software features such as
9 multi-compartmental neurons, fast re-configuration for hardware-in-the-loop training, applications
10 for the embedded processors, the non-spiking operation mode, interactive platform access, and
11 sustainable hardware/software co-development. Finally, we discuss further developments in terms
12 of hardware scale-up, system usability and efficiency.

13 Keywords: hardware abstraction, neuroscientific modeling, accelerator, analog computing, neuromorphic,
14 embedded operation, local learning

1 INTRODUCTION

15 The feasibility and scope of neuroscientific research projects is often limited due to long simulation runtimes
16 and therefore long wall-clock runtimes, especially for large-scale networks (van Albada et al., 2021). Other
17 areas of neuromorphic research—such as lifelong learning in robotic applications— inherently rely on
18 very long network runtimes to capture physical transformations of their embodiment on the one hand and
19 evolutionary processes on the other. Furthermore, training mechanisms relying on iterative reconfiguration
20 benefit from low execution latencies.

[†]contributed equally.

21 Traditional software-based simulations typically still often rely on general-purpose high-performance
22 computing (HPC) hardware. While some efforts towards GPU-based accelerators provide an intermediate
23 step to improve scalability and runtimes (Abi Akar et al., 2019; Yavuz et al., 2016), domain-specific
24 accelerators —a subset of which are neuromorphic hardware architectures—, have come more and more
25 into the focus of HPC (Dally et al., 2020). Such systems specifically aim to improve on performance and
26 scalability issues — both, in the strong and in the weak scaling cases. Particularly, the possibility to achieve
27 high throughput at low execution latencies can pose a crucial advantage compared to massively parallel
28 simulations.

29 The BrainScaleS (BSS) neuromorphic architecture is an accelerator for spiking neural networks based on
30 a physical modeling approach. It provides a neuromorphic substrate for neuroscientific modeling as well as
31 neuro-inspired machine learning. Earlier work shows its scalability in wafer-scale applications, emulating
32 up to 200k neurons and 40M synapses (Schmitt et al., 2017; Göltz et al., 2021; Kungl et al., 2019; Müller
33 et al., 2020b), as well as its energy-efficient application as standalone system with 512 neurons and 128k
34 synapses in use cases related to edge computing (Stradmann et al., 2021; Pehle et al., 2022). Compared to
35 the biological time domain, the model dynamics evolve on a 1000-fold accelerated time scale making the
36 system interesting for iterative and long-running experiments. Constant model emulation speed is attractive
37 for hardware users. However, it often comes with algorithmic challenges. Similar to other neuromorphic
38 systems based on the physical modeling concept, neuroscientific modeling on the BrainScaleS-2 (BSS-2)
39 system requires a translation from a user-defined neural network experiment to a corresponding hardware
40 configuration. BSS-2 operates in continuous time and does not support pausing or resuming of model
41 dynamics. The algorithmic problem statement is global for the user-defined experiment. Therefore, the
42 complexity of the translation process cannot be reduced by partitioning the problem. Many neuromorphic
43 systems have been providing software solutions to solve this problem and enable higher-level experiment
44 descriptions. We developed a software stack for the wafer-scale BrainScaleS-1 (BSS-1) system covering
45 the translation of user-defined experiments from the PyNN high-level domain-specific description language
46 to a hardware configuration (Müller et al., 2020b). While the BSS-2 neuromorphic architecture hasn't
47 been scaled to full wafer size yet, other feature additions such as structured and non-linear neuron models
48 as well as single instruction, multiple data (SIMD) processors make BSS-2 an appealing substrate for
49 modeling of smaller network sizes. In particular, a new challenge is posed by the introduction of SIMD
50 processors in BSS-2 as programmable elements with real-time vectorized access to many observables from
51 the physical modeling substrate. Observables such as correlation sensors are implemented in the synapse
52 circuits, yielding an immense computational power by offloading computational tasks into the analog
53 substrate. Moreover, the configuration space increases significantly: in addition to a static configuration of
54 network topology, the processors allow for flexible handling of dynamic aspects such as structural plasticity,
55 homeostatic behavior, and virtual environments enabling robotic or other closed-loop applications. This
56 “hybrid” approach requires modeling support in the software stack integrating code generation for the
57 processors as well as mechanisms to parameterize plasticity algorithms and other code parts running on the
58 embedded processors.

59 We present recent modeling advances on the substrate showcasing new features of the system: complex
60 neurons (section 3.1), neuro-inspired machine-learning experiments (section 3.2), closed-loop sensor-
61 motor interaction (section 3.3) and non-spiking operation (section 3.4). We demonstrate network-attached
62 accelerator operation as well as standalone operation. We argue that for successful and sustainable advances
63 in the usage of neuromorphic systems a deep integration between hardware and software is crucial on
64 all layers. The complete system —software together with hardware— needs to be explicitly designed to
65 support access with varying abstraction levels: high-level modelers, expert users and component developers

66 possess different perceptions of the system; in order for a modeling substrate to be successful, it has to
 67 deliver on all of these aspects.

68 1.1 The BrainScaleS-2 hardware

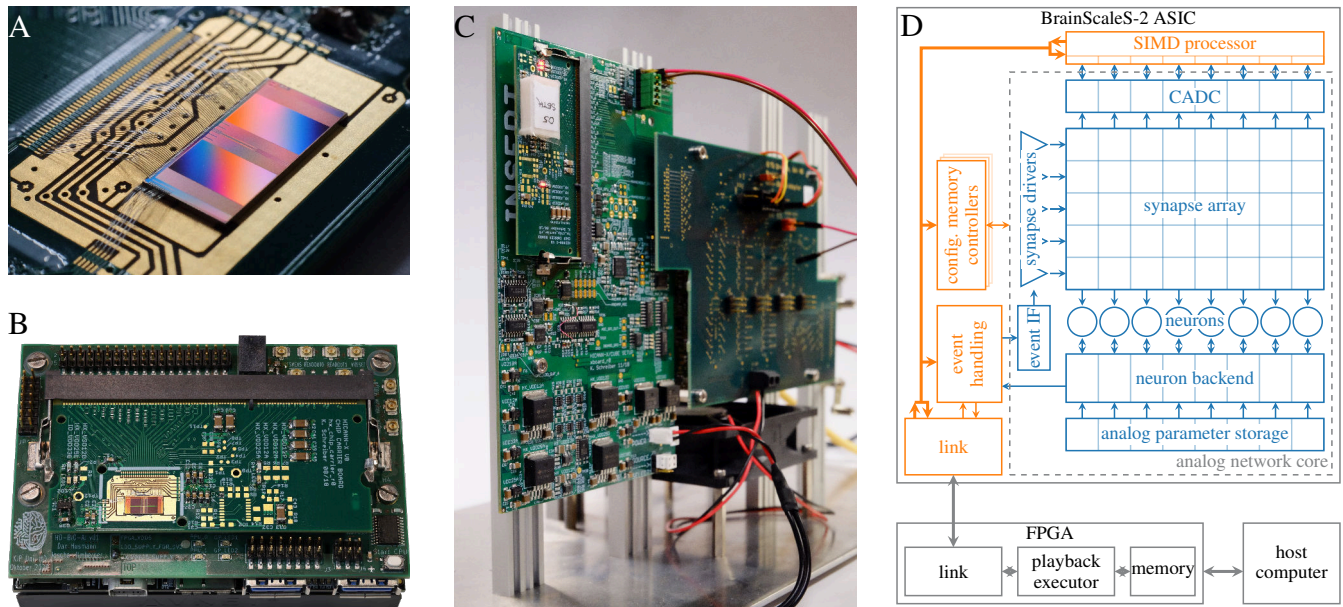


Figure 1. Overview of the BSS-2 system. (A) BSS-2 application-specific integrated circuit (ASIC) bonded to a carrier board. The ASIC is organized in two hemispheres each hosting 256 neurons and the accompanying synapse matrix, cf. (D). (B) Portable BSS-2 system. (C) Laboratory setup. (D) Overview over the signal flow in the BSS-2 system. The depicted analog neural network core and SIMD processor represent one of the two hemispheres visible in (A), which are mirrored vertically below the neurons.

69 In this section we introduce the BSS-2 system and highlight the basic hardware design which is guiding
 70 the development of the accompanying software stack. For a more in depth description of the hardware
 71 aspects of the BSS-2 system refer to Pehle et al. (2022); Schemmel et al. (2020); Aamir et al. (2018).

72 BrainScaleS is a family of mixed-signal neuromorphic accelerators; analog circuits emulate neuron as
 73 well as synapse dynamics in continuous time, while communication of spike events and configuration data
 74 is handled in the digital domain. In this paper we focus on the single chip BSS-2 system with 512 neurons
 75 and 131 072 synapses circuits, see fig. 1 A. Due to the intrinsic properties of the silicon substrate, the
 76 physical emulation of neuron dynamics is 1000× faster than in biological real time. Currently, the BSS-2
 77 ASIC is integrated in a stationary laboratory setup, fig. 1 C, as well as in a portable system, fig. 1 B.

78 The high configurability of the BSS-2 system facilitates many different applications, see section 3. For
 79 example, the neuron circuits replicate the dynamics of the adaptive exponential integrate-and-fire (AdEx)
 80 neuron model (Brette and Gerstner, 2005) and are individually configurable by a number of analog and
 81 digital parameters. By connecting several neuron circuits together to form one logical neuron, more
 82 complex multi-compartmental neuron models can be formed and the synaptic fan-in of individual neurons
 83 can be increased; a single neuron circuit on its own has access to 256 synapses (fig. 1 D). In addition to
 84 the emulation of biologically plausible neural networks, BSS-2 also supports non-spiking artificial neural
 85 networks (ANNs). This is facilitated by disabling spiking as well as the exponential, the adaptive and the
 86 leak current of the AdEx neuron model, turning the neuron circuits into simple integrators. Furthermore,

87 the high configurability allows countering device-specific deviations between analog circuits which result
88 from imperfections during the manufacturing process, see section 2.3.6.

89 The digital handling of spike events enables the implementation of various network topologies. All spikes,
90 including external spikes as well as spikes generated in the neuron circuits, are collected in the “event
91 handling” block and subsequently routed off chip for recording or via the synapse drivers and synapses
92 to post-synaptic on-chip partners, cf. fig. 1 D. One of the key challenges during experiment setup is the
93 translation of neural networks to valid hardware configurations. This includes assigning specific neuron
94 circuits to the different neurons in the network as well as routing events between neurons, cf. sections 2.3.2
95 and 2.3.3.

96 Apart from forwarding spikes, the synapse circuits are also equipped with analog correlation sensors
97 which measure the causal and anti-causal correlation between pre- and post-synaptic spikes. The measured
98 correlation can be accessed by two columnar ADCs (CADCs), which measure correlations row-wise in
99 parallel and can be used in the formulation of plasticity rules, cf. sections 3.2 and 3.3. An additional
100 analog-to-digital converter (ADC), the so-called membrane ADC (MADC), offers the possibility to record
101 single neurons with a higher temporal and value resolution.

102 Aside the analog neural network core, two embedded SIMD processors, based on the Power™ architec-
103 ture (PowerISA, 2010), which allow for arbitrary calculations and reconfigurations of the BSS-2 ASIC
104 during hardware runtime and are the experiment master in standalone operation. They are equipped with
105 16 KiB static random-access memory (SRAM) memory each and feature a weakly-coupled vector unit
106 (VU), which can access the hemisphere-local synapse matrix as well as the CADC.

107 Communication to the BSS-2 ASIC as well as real-time runtime control is handled by a field-
108 programmable gate array (FPGA). It provides memory buffers for data received from a host computer or
109 from the chip, with which it orchestrates experiment executions in real time, see section 2.1. To allow for
110 more complex programs and larger data storage, the on-chip processors can access memory connected to
111 the FPGA.

112 The software stack covered in this paper handles all the necessary steps to turn high-level experiment
113 descriptions into configuration data, spike stimuli or programs for the on-chip SIMD processor.

114 In the following we will at first describe the BSS-2 Operating System (BSS-2 OS) in section 2 before
115 showcasing several applications in section 3. We conclude the paper with a discussion in section 4.

2 BRAINSCALES-2 OPERATING SYSTEM

116 This section introduces key concepts and software components that are essential for the operation of
117 BrainScaleS-2 systems. First, we introduce the workflow of experiments incorporating BSS-2, derive an
118 execution model and specify common modes of operation in section 2.1. Continuing, we give a structural
119 overview of the complete software stack including the foundation developed in (Müller et al., 2020a)
120 in section 2.2. Following this, we motivate key design decisions and show their incorporation into the
121 development of the software stack in section 2.3. Finally, we describe advancements in platform operation
122 towards seamless integration of BSS-2 as an accelerator resource in multi-site compute environments in
123 section 2.4.

124 Higher abstraction layers scale down the level of required hardware detail knowledge. Naturally, such
 125 abstractions impose constraints on and reduce the flexibility of system usage introducing tradeoffs. There-
 126 fore, there are tradeoffs between abstraction level and the flexibility to exploit system capabilities. In the
 127 following, we explain existing tradeoffs at their occurrence.

128 2.1 Experiment Workflow

129 Unlike numerical simulations, which are orchestrated as number-crunching on traditional computers,
 130 experiments on BSS-2 are more akin to physical experiments in a traditional lab. Just like for these
 131 there is an *initialization* phase, which ensures the correct configuration of the system for this particular
 132 experiment and a *real-time* section, where the network dynamics are recorded and the actual emulation
 133 happens. If multiple emulations share (parts of) the configuration, those experiments can be composited by
 134 concatenating the trigger commands for both input and recording (see fig. 2).

135 The fundamental physical nature of the emulation on BSS-2 requires these control commands to be issued
 136 with very high temporal precision as the dynamics of the on-chip circuitry can neither be interrupted nor
 137 exactly repeated. To achieve this, the accompanying FPGA is used to play-back a sequence of instructions
 138 with clock-precise timing, in the order of 10 ns. In order to limit the FPGA firmware complexity, the play-
 139 back unit is restricted to sequential execution, which includes blocking instructions (used for times without
 140 explicit interaction), but excludes branching instructions. Concurrently to the FPGA-based instruction
 141 sequence execution, the embedded single instruction, multiple data central processing units (SIMD CPUs)
 142 can be configured to perform readout of observables and arbitrary alterations to the hardware configuration.
 143 This means that conditional decisions, e.g. the issuance of rewards, can be performed either via the SIMD
 144 CPU if they are not computationally too complex or via synchronization with the executing host computer
 145 which in the current setup has no guaranteed timing.

146 The initialization phase typically includes time-consuming write operations to provide an initial state
 147 of the complete hardware configuration. This is due to both, the amount of data to be transmitted,
 148 e.g. for the synapse matrix, and required settling-time for the analog parameters. Since this can take
 149 macroscopic amounts of time, at least around 100 μ s due to round-trip latency, around 100 ms for a
 150 complete reconfiguration, back-to-back concatenation of real-time executions is needed to keep their
 151 timeshare high and therefor the configuration overhead low.

152 Due to the hardware's analog speed-up factor compared to typical biological processes, a single real-time
 153 section can be short compared to the initialization phase. Therefore, we concatenate multiple real-time
 154 sections after a single initialization phase to increase the real-time executions' timeshare. In the following,
 155 this composition is called execution instance and is depicted in fig. 2.

156 Alternatively, instead of this asynchronous high-throughput operation, the low minimal latency allows for
 157 fast iterative workflows with partial reconfiguration, e.g., iterative reconfiguration of a small set of synaptic
 weights.

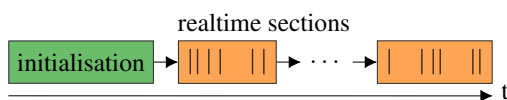


Figure 2. Time evolution of a single execution instance. The initialization is followed by possibly multiple real-time executions with input spike-trains represented by vertical lines.

158

159 Based on this we differentiate between three modes of operation. First, in batch-like operation one
 160 or multiple execution instances are predefined and run on hardware. Second, in the so-called hardware

161 in-the-loop case hardware runs are executed iteratively where the results of previous runs determine the
162 parameters of successive runs. Last, in closed-loop operation is characterized by tightly coupling the
163 network dynamics of the analog substrate to the experiment controller, either the SIMD CPU or the control
164 host.

165 2.2 Software Stack Overview

166 Structuring software into well-defined layers is vital for keeping it maintainable and extendable. The
167 layers are introduced and implemented via a bottom-up approach matching the order of requirements in the
168 current stage of the hardware development and commissioning process. This means, that first raw data
169 exchange and transport from and to the hardware via the communication layer is established. Subsequently,
170 the hardware abstraction layer implements translation of typed configuration, e.g. enabling a neuron's
171 event output, to and from this raw data. On this level, the calibration layer allows to programmatically
172 configure the analog hardware to a desired working point. Then, hardware-intrinsic relations between
173 configurables and their interplay in experiments, cf. section 2.1, is encapsulated in a graph structure. Lastly,
174 automated generation of hardware configuration from an abstract network specification enables embedding
175 into modelling frameworks for high-level usage. Figure 3 gives a graphical overview of this software
176 architecture¹.

177 2.2.1 Communication

178 From the software point of view, the first step to utilize hardware systems is the ability to exchange
179 data. With proper abstraction the underlying transport protocol and technology are interchangeable.
180 Communication is therefore structured into a common *connection* interface *hxcomm*² that supports various
181 back-ends.

182 For most hardware setups, we use a custom, reliable regarding data integrity, transport protocol on top of
183 the user datagram protocol (UDP), *Host-ARQ* provided by *sctrltp*³. Additionally, we support connection to
184 hardware design simulations via *flange*⁴, compare section 3.6 for both the use during debugging of current
185 and unit testing of future chip generations. Multi-site workflows are transparently enabled already at this
186 level via the micro scheduler *quiggeldy*⁵.

187 2.2.2 Hardware Abstraction

188 A major aspect of any system configuration software is *hardware abstraction*, which encapsulates
189 knowledge about the raw bit configuration, e.g. that bit i at address j corresponds to enabling neuron k 's
190 event output. It therefore decouples hardware usage and detailed knowledge about its memory layout,
191 which is an important step towards providing hardware access beyond the group of developers of the
192 hardware. Responsibility of this layer can be compared to device drivers. The layers provide an abstract
193 software representation of various hardware components, such as synaptic weights on the chip or values of
194 supply voltages on the periphery board, as well as their control flow.

¹ All the repositories mentioned in the following are available at <https://github.com/electronicvisions> under the *GNU Lesser General Public License v2/v3*.

² *hxcomm* is available at <https://github.com/electronicvisions/hxcomm>

³ *sctrltp* is available at <https://github.com/electronicvisions/sctrltp>

⁴ *flange* is available at <https://github.com/electronicvisions/flange>

⁵ *quiggeldy* is available at <https://github.com/electronicvisions/hxcomm>

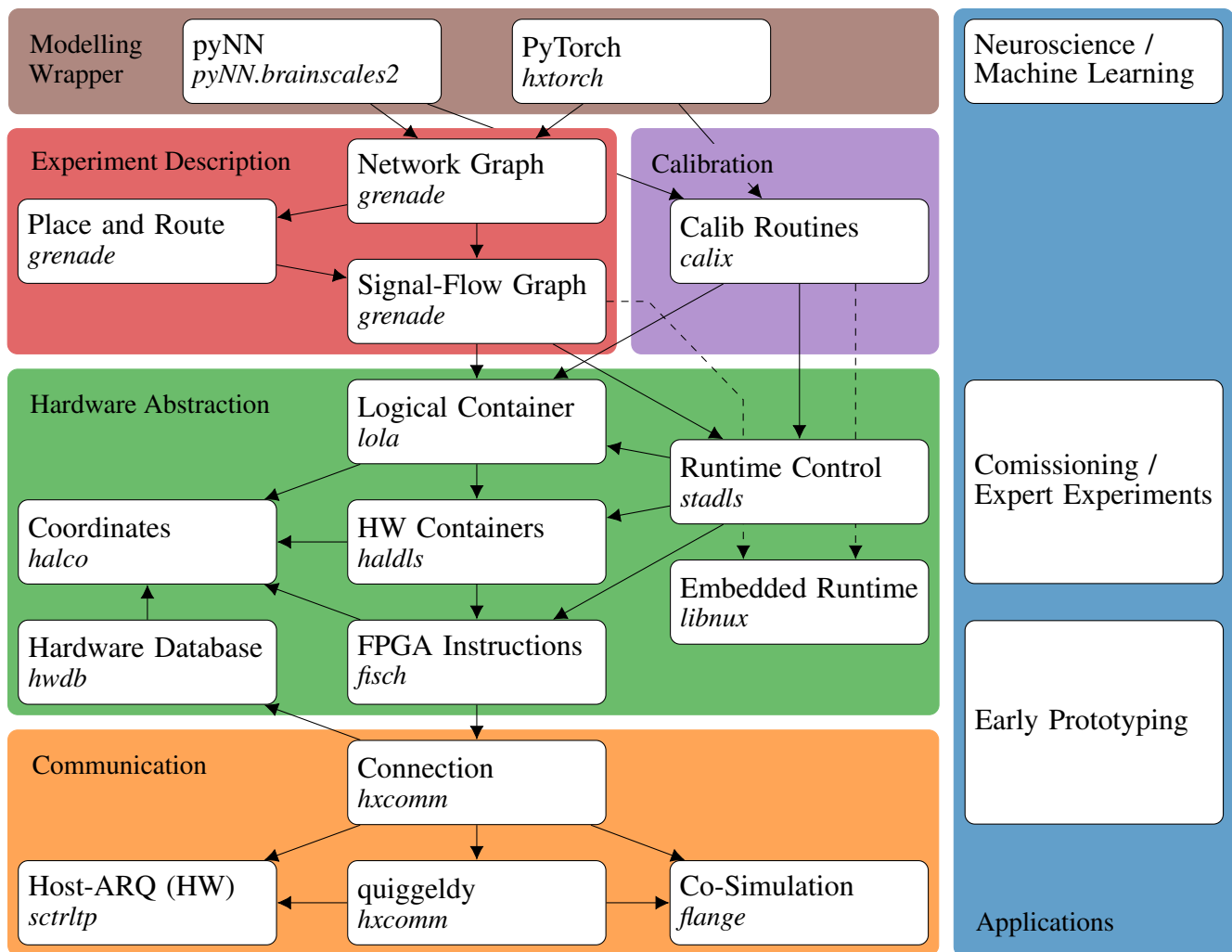


Figure 3. Overview of the BSS-2 software architecture and its applications. Left side: Colored boxes in the background represent the separation of the software into different concerns. White boxes represent individual software APIs or libraries with their specific repositories names and dependencies. Right side: Various applications concerning different system aspects. The arrows represent dependencies in the stack, where the dependent points to its dependencies. For embedded operation additional dependencies on *libnux* are needed (dashed arrows).

195 Within this category the lowest layer is *fisch*⁶ (FPGA Instruction Set arCHitecture), the abstraction
 196 of FPGA instructions. Combined with communication software this is already sufficient to provide an
 197 interface for prototyping in early stages of system development, i.e., the possibility to manually read and
 198 write words at memory locations. With knowledge of the hardware's memory layout this allows specifying
 199 addresses and word values directly, e.g. bit i (and all other bits in this word with possibly unrelated effects)
 200 at address j which then enables the neuron k 's event output.

201 The heterogeneous set of entities on the hardware as well as their memory layout is arranged via
 202 geometric pattern and contain symmetries, e.g. a row of neurons or a matrix of synapses. An intuitive
 203 structure of this fragmented address space is provided by the *coordinate* layer *halco*⁷. It represents hardware
 204 components by custom ranged types that can be converted to other corresponding coordinate types, e.g. a

⁶ *fisch* is available at <https://github.com/electronicvisions/fisch>

⁷ *halco* is available at <https://github.com/electronicvisions/halco>

205 `SynapseOnSynapseRow` as a ranged integer $i \in [0, 256)$, that allows conversion to a neuron column, see
206 (Müller et al., 2020a).

207 A software representation of the configuration space of hardware components is implemented by the
208 *container* layer *haldls*⁸. For example a `NeuronConfig` contains a boolean parameter for enabling the
209 spike output. These configuration containers are translatable (e.g. a neuron container represents one, but
210 not a specific one, of the neurons) and also define methods for de- and encoding between their abstract
211 representation and the on-hardware data format given a location via a supplied *coordinate*. A logical
212 function- instead of a hardware subsystem-centered container collection is implemented by the *lola*⁹ layer.
213 For example the `AtomicNeuron` collects the analog and digital configuration of a single neuron circuit,
214 which is scattered over two digital configurations and a set of elements in the analog parameter array.

215 The *runtime control* layer *stadls*¹⁰ provides an interface to describe timed sequences of read and write
216 instructions of pairs of coordinates and containers, e.g. changing the synaptic weight of synapse i, j at time
217 t , as well as event-like response data, e.g. spikes or ADC samples. These timed sequences, also called
218 playback programs, can then be loaded to and executed on the FPGA which records the response data.
219 Afterwards, the recorded data is transferred-back to the host computer.

220 We track the constitution of all hardware setups in a database, *hwdb*¹¹. It is used for compatibility checks
221 between hardware and software as well as for the automated selection of stored calibration data. We also
222 use it to provide the resource scheduling service with information about all available hardware systems.

223 This set of layers is feature-complete to formulate arbitrary hardware-compatible experiments and was
224 used as basis for experiments in Göltz et al. (2021); Czischek et al. (2022); Klassert et al. (2021); Schemmel
225 et al. (2020); Cramer et al. (2022).

226 2.2.3 Embedded runtime

227 In addition to the controlling host system, the two SIMD CPUs on the BSS-2 ASIC require integration into
228 the BSS-2 OS. To enable users to efficiently formulate their programs, we provide a development environ-
229 ment based on C++. It specifically consists of a cross-compilation toolchain based on `gcc` (GNU Project,
230 2018) that has been adapted to the custom SIMD extensions of the integrated microprocessors (Müller et al.,
231 2020a). More abstract functionality is encapsulated in the support library *libnux*¹², which provides various
232 auxiliary functionality for experiment design. Moreover, the hardware abstraction layer of the BSS-2
233 OS (cf. section 2.2.2) supports the SIMD CPUs as an additional cross-compiled target for configuration
234 containers as well as coordinates.

235 2.2.4 Calibration

236 In order to tune all the analog hardware parameters to the requirements given by an experiment, we
237 provide a calibration framework, *calix*¹³. For example, an experiment might require a certain set of synaptic
238 time constants for which analog parameters are to be configured while counteracting circuit inequalities.
239 In section 2.3.6, this layer’s design is explained in detail. The `Python` module supplies a multitude of
240 algorithms and calibrations for each relevant component of the circuitry: A calibration provides a small

⁸ *haldls* is available at <https://github.com/electronicvisions/haldls>

⁹ *lola* is available at <https://github.com/electronicvisions/haldls>

¹⁰ *stadls* is available at <https://github.com/electronicvisions/haldls>

¹¹ *hwdb* is available at <https://github.com/electronicvisions/hwdb>

¹² *libnux* is available at <https://github.com/electronicvisions/libnux>

¹³ *calix* is available at <https://github.com/electronicvisions/calix>

241 experiment based on the hardware abstraction layer, see section 2.2.2, which is executed on the chip for
242 characterization. An iterative algorithm then decides how configuration parameters should be changed in
243 order to match the measured data with given expectations.

244 The user-interfacing part provides functions that take a set of target parameters and return a serializable
245 calibration result that can be injected to experiment toplevels, cf. section 2.2.6. Additionally, we have the
246 option to calibrate the analog circuits locally on chip, using the embedded processors. Aside of enabling
247 arbitrary user-defined calibrations, we provide default calibrations for spiking operation, cf. for example
248 sections 3.1 and 3.2, and non-spiking matrix-vector multiplication, cf. section 3.4 for convenient entry.
249 They are generated nightly via continuous deployment (CD).

250 2.2.5 Experiment Description

251 With rising experiment and network topology complexity, a coherent description ensuring topology
252 and data-flow correctness becomes beneficial. Therefore, a signal-flow graph is defined representing the
253 hardware configuration and experiment flow. Compilation and subsequent execution via the hardware
254 abstraction layer, cf. section 2.2.2, of this graph in conjunction with supplied data, e.g. spike events, then
255 forms an experiment execution. The applied execution model follows the experiment workflow described
256 in section 2.1. It, therefore, restricts flexibility to enable network-topology-based experiment descriptions
257 and the separation of data-flow description and data.

258 While this aids in construction of complex experiments, detailed knowledge of configuration and its
259 interplay is still required. Solving this, a high-level abstract representation of neural network topology
260 building on top of the signal-flow graph description is developed. An automated translation from this high-
261 level abstraction to a valid hardware configuration is handled by a place-and-route algorithm. This enables
262 hardware usage without detailed knowledge of event routing capabilities and interplay of configuration.
263 While relieving users from providing a valid hardware configuration, this automatism requires tradeoffs to
264 be made between the computational complexity of the algorithms and the size of the explored configuration
265 space to find a matching hardware configuration for a given abstract network representation.

266 This layer is contained in *grenade*¹⁴, short for GRaph-based Experiment Notation And Data-flow
267 Execution. Its design is explained in detail in section 2.3.2.

268 2.2.6 Modeling Wrapper

269 Various back-end-agnostic modeling languages emerged to provide access to various simulators or
270 neuromorphic hardware systems to a wide range of researchers. The BSS-2 software stack comprises
271 wrappers to two of such modeling frameworks: PyNN (Davison et al., 2009) via *pyNN.brainscales2*¹⁵ and
272 PyTorch (Paszke et al., 2019) via *hxtorch*¹⁶ (Spilger et al., 2020). Their goal is to provide a common user
273 interface and to embed different back-ends into an existing software ecosystem. This allows users to benefit
274 from a consistent and prevalent interface and integration into their established work-flow. The design of
275 these layers' integration with BSS-2 is explained in detail in section 2.3.4 for PyNN and in section 2.3.5
276 for PyTorch.

¹⁴ *grenade* is available at <https://github.com/electronicvisions/grenade>

¹⁵ *pyNN.brainscales2* is available at <https://github.com/electronicvisions/pynn-brainscales>

¹⁶ *hxtorch* is available at <https://github.com/electronicvisions/hxtorch>

277 2.3 Software Design

278 We base the full-stack software design on the principles laid out in Müller et al. (2020a). We use C++ as
279 the core language to ensure high performance and make use of its compile-time expression evaluation and
280 template metaprogramming capabilities. Due to the heterogeneous hardware architecture we employ type
281 safety for logical correctness and compile-time error detection. Serialization support of configuration and
282 control flow enables multi-site workflows as well as archiving of experiments.

283 In the following, we show enhancements of the hardware abstraction layer, see section 2.2.2, introduced
284 in Müller et al. (2020a) as well as design decisions for the full software stack with high-level user interfaces.
285 First, support for multiple hardware revisions is shown in section 2.3.1. Then, the signal-flow graph-based
286 experiment notation is derived in section 2.3.2. Following, an abstract network description explained
287 in section 2.3.3 closes the gap to the modelling wrappers in PyNN, cf. section 2.3.4 and PyTorch, cf.
288 section 2.3.5. Closing, the calibration framework is described in section 2.3.6.

289 2.3.1 Multi-revision hardware support

290 As platform development progresses, new hardware revisions require software support. This holds true
291 for both, the ASIC and the surrounding support hardware like the FPGA and system printed circuit boards
292 (PCBs). Additionally, the platform constitution evolves, e.g. by introduction of a mobile system with still
293 one chip but different support hardware or a multi-chip setup.

294 After a potential development of a second revision, a heterogeneous set of hardware setups may co-
295 exist. For one generation of chips, it is typically possible to combine different revisions with different
296 surrounding hardware configurations, leading to a number of combinations given by the Cartesian product
297 $N = M_{\text{ASIC}} \times M_{\text{Platform}_1} \times \dots \times M_{\text{Platform}_P}$, where M_{Platform_i} is the number of configurations for a given
298 part of the platform, e.g. the FPGA and M_{ASIC} is the revision of the BSS-2 ASIC.

299 We provide simultaneous software support by dependency separation and extraction of common code
300 for each affected component across all affected software layers. This way, code duplication is minimized,
301 maintainability of common features is ensured and divergence of software support is prevented. Moreover,
302 phasing-out or retiring hardware revisions is possible without effecting the software infrastructure of
303 other revisions. The to be implemented software reduces to $N' = M_{\text{ASIC}} + M_{\text{Platform}_1} + \dots + M_{\text{Platform}_P}$
304 constituents, the combinations are rolled-out automatically. We use C++ namespaces for separation and
305 C++ templates for common code, which depends on the individual platform's constituents.

306 2.3.2 Signal-flow graph-based experiment notation

307 As stated in section 2.2.2, the hardware abstraction developed in Müller et al. (2020a) is already feature-
308 complete to formulate arbitrary hardware-compatible experiments. However, it lacks a representation
309 of intrinsic relations between different configurable entities. For example, the hard-wired connections
310 between synapse drivers and synapse rows are not represented in their respective configuration but only
311 given implicitly.

312 Neural networks are predominantly described as graphs. For spiking neural networks single neurons or
313 collections thereof and their connectivity form a graph (Davison et al., 2009; Gewaltig and Diesmann,
314 2007; Goddard et al., 2001). In machine-learning, the two major frameworks PyTorch (Paszke et al., 2019)
315 and Tensorflow (Abadi et al., 2015) use a graph-based representation of tensor computation or are moving
316 into this direction (PyTorch's JIT intermediate representation (Facebook, Inc., 2021a) and XLA back
317 end (Facebook, Inc., 2021b; Suhan et al., 2021)).

318 Inspired by this, we implement a signal-flow graph-based experiment abstraction. A signal-flow
 319 graph (Mason, 1953) is a directed graph, where vertices receive signals from their in-neighborhood,
 320 perform some operation, and transmit an output signal to their out-neighborhood. We integrate this
 321 representation at the lowest possible level to fully incorporate all hardware features without premature
 322 abstraction.

323 For BSS-2, the graph-based abstraction is applied at two granularities, see fig. 4. First, the initial
 324 static network configuration as well as virtualized computation using the on-chip embedded processors is
 325 abstracted as a signal-flow graph. Second, data-flow between multiple individual real-time experiments
 326 distributed over chips and time are described as a graph.

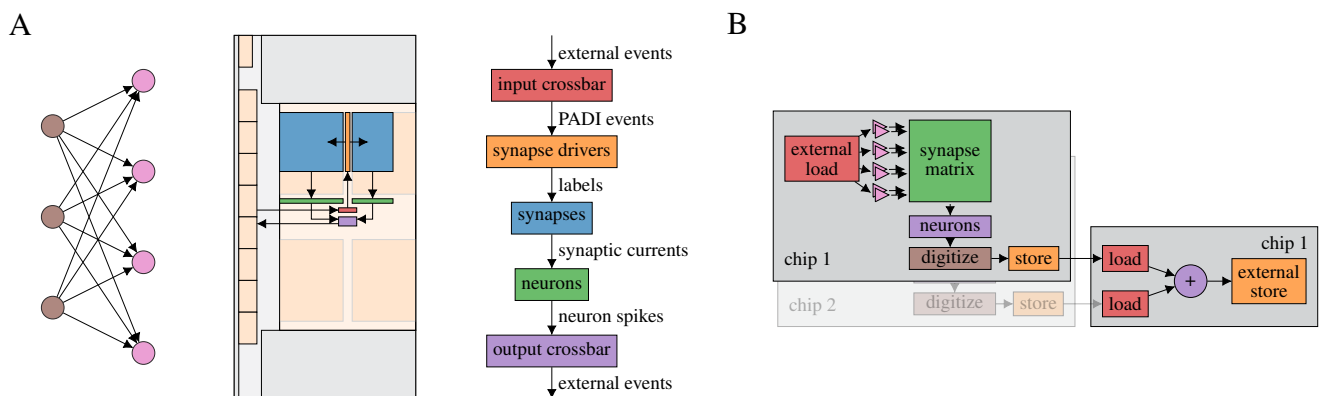


Figure 4. Signal-flow graph-based experiment abstraction on BSS-2. **(A)** Placed feed-forward network represented as signal-flow graph. Left: Abstract network; Middle: Actual layout on the chip, the arrows represent the graph edges; Right: The network graph structure enlarged with signal type annotation on the edges. The color links the same entities in the middle (chip schematic) and right subfigure (vertical data-flow graph). **(B)** Non-spiking network distributed over two physical chips, adapted from Spilger et al. (2020). The result of two matrix multiplications on chips 1 and 2 is added on chip 1. The latter execution instance depends on the output of the two former instances.

327 The signal-flow graph representation yields multiple advantages. Type safety in the graph constituents
 328 facilitates experiment correctness regarding on-chip connectivity and helps to avoid inherently dysfunctional
 329 experiments already during specification. Debugging benefits from visualisation of the graph representation,
 330 which directly contains implicit on-chip connectivity. Finally, the signal-flow graph is the ideal source of
 331 relationship information for on-chip entity allocation optimization or merging of digital operations.

332 However, the actual signals are not part of the signal-flow graph representation. They are either provided
 333 separately (e.g. external events serving as input), will only be present locally upon execution (e.g. synaptic
 334 current pulses) or will be generated by execution (e.g. recorded external events). We implement the
 335 experiment workflow described in section 2.1 consisting of an initial static configuration followed by a
 336 collection (batch) of time evolutions, see fig. 2.

337 The signal-flow graph is a recipe for compilation towards the lower-level hardware abstraction layer,
 338 cf. Müller et al. (2020a), and eventual execution. The specific implementation of the compilation and
 339 execution process is separate from the graph representation in order to allow extensibility and multiple
 340 solutions for different requirement profiles. Here, we present a just-in-time (JIT) execution implementation.
 341 It supports both, spiking and non-spiking experiments. For every execution instance, the local subgraph is
 342 compiled into a sequence of instructions, executed and its results processed in order for them to serve as

343 inputs for the out-neighborhood. While it is feature-complete for the graph representation, it introduces
 344 close coupling between the execution on the neuromorphic hardware and the controlling host computer.
 345 Host-based compilation can be performed concurrently to hardware execution, increasing parallelism.
 346 Figure 5 shows concurrent execution of multiple execution instances (A) and the compilation and execution
 of a single execution instance (B).

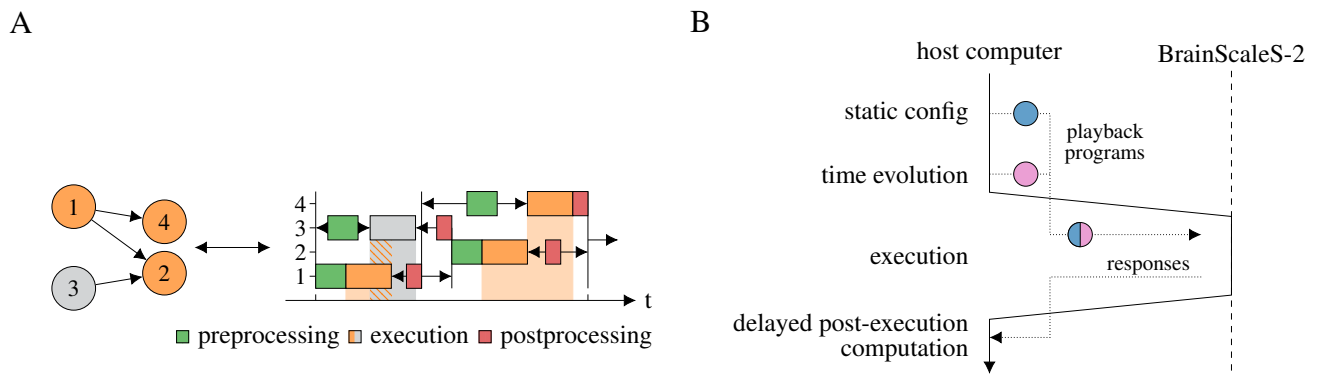


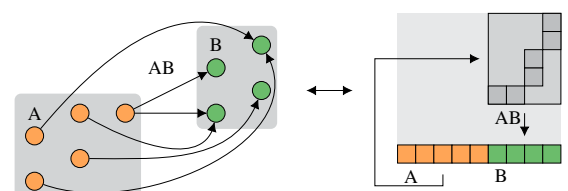
Figure 5. JIT compilation and execution of signal-flow graph of multiple execution instances and within a single execution instance. (A) JIT execution of a graph on two physical chips, adapted from Spilger et al. (2020). Left: Execution instance 3 is to be executed on another physical chip than the other execution instances. Right: The execution of instance 3, depicted in gray, can be performed concurrently to execution instance 1. (B) JIT compilation and execution of a single execution instance subgraph. First, the static configuration is extracted by a vertex visit and transformed to hardware configuration where applicable. Then, the real-time execution is built by a vertex visit. This built program is executed on the neuromorphic hardware and results are transmitted back to the host computer. Finally, delayed digital operations, which require output data from the execution, are performed on the host computer.

347

348 2.3.3 Abstract network description

349 The signal-flow graph-based notation from section 2.3.2 eases creation of correct experiments while
 350 minimizing implicit knowledge. However, knowledge of hardware routing capabilities is still required
 351 to create a graph-based representation of the hardware configuration which performs as expected. This
 352 should not be required to formulate high-level experiments. To close this gap, an abstract representation
 353 similar to PyNN (Davison et al., 2009), consisting of populations as collections of neurons and projections
 354 as collections of synapses, is developed. Given this description, an algorithm finds an event routing
 355 configuration to fulfill the abstract requirements and generates a concrete hardware configuration. This
 356 step is called routing. Figure 6 visualizes an abstract network description and one corresponding hardware
 configuration.

Figure 6. Abstract network notation. Population A consisting of five neurons is connected to population B consisting of four neuron via projection AB. Left: Abstract network; Right: Placed and routed on the hardware, where the projection AB consists of synapses in the two-dimensional synapse matrix and the populations A and B are located in the neuron row, compare fig. 1 (D).



357

358 2.3.4 Integration of PyNN

359 When it comes to modeling spiking neural networks, a widely used API is PyNN (Davison et al., 2009).
 360 It is supported by various neural simulators like NEST (Gewaltig and Diesmann, 2007), NEURON (Hines
 361 and Carnevale, 2003) and Brian (Stimberg et al., 2019), as well as by neuromorphic hardware platforms like
 362 SpiNNaker (Rhodes et al., 2018) or the predecessor hardware of BSS-2: BSS-1 (Müller et al., 2020b) and
 363 Spikey (Brüderle et al., 2009). With the aim of easy access to BSS-2, we expose its hardware configuration
 364 via the PyNN interface. The module `pyNN.brainscales2` implements the PyNN-API for BSS-2. It
 365 offers a custom cell type, `HXNeuron`, which corresponds to a physical neuron circuit on the hardware and
 366 replicates the `lola.AtomicNeuron` from the hardware abstraction layer, see section 2.2.2. This allows
 367 to set parameters directly in the hardware domain and gives expert users the possibility to precisely control
 368 the hardware configuration while at the same time take advantage of high-level features such as neuron
 369 populations and projections. Figure 7 illustrates how these parameters are available in the corresponding
 370 interfaces. An additional neuron type supporting the translation from neuron model parameters in SI units
 371 is currently in the planning. Otherwise, the PyNN program looks the same as for any other back end. Since
 372 the PyNN-API is free from hardware placement specifications, they are algorithmically determined by
 373 mapping and routing in *grenade*, cf. section 2.3.3. This step is performed automatically upon invocation of
 374 `pyNN.run()`, so that the user is not required to have any particular knowledge about event routing on the
 375 hardware. Nevertheless, the interface allows that an experimenter can adjust any low-level configuration
 376 aside from neuron parameters and synaptic weights.

```

neuron = lola.AtomicNeuron()
neuron.leak.v_leak = 650
neuron.leak.i_bias = 420
neuron.leak.enable_division = True

pyNN.Population(1, pyNN.HXNeuron({
    "leak_v_leak": 650,
    "leak_i_bias": 420,
    "leak_enable_division": True}))

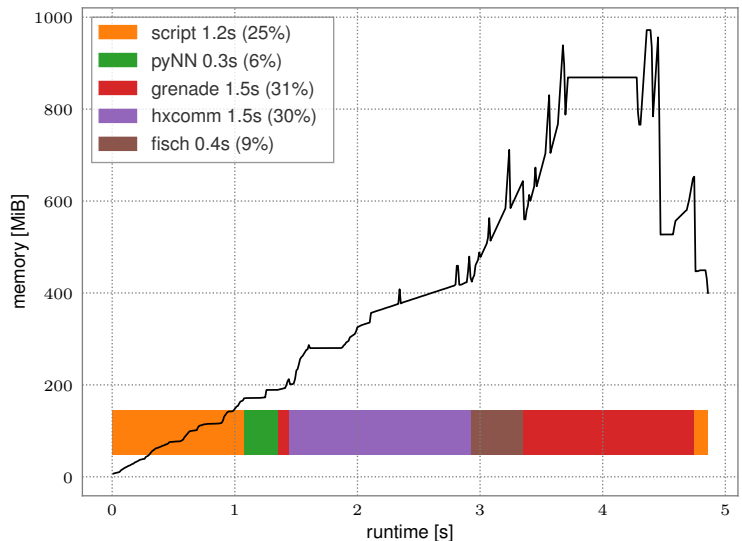
```

Figure 7. Comparison between `lola.AtomicNeuron` and `pyNN.HXNeuron`.

377 To exploit the full potential of the accelerated hardware the software implementation's overhead shall be
 378 minimal. Figure 8 presents runtime and memory consumption analysis of the whole PyNN-based stack
 379 for a high spike count benchmark experiment. 12 neurons are excited by a regular spike train with 1 MHz
 380 frequency and their activity is recorded for one second. These settings are chosen as they roughly equate to
 381 the maximum recording rate without loss.

382 The initial overhead of importing Python libraries and setting up the PyNN environment only needs to
 383 be performed once for every experiment and is independent of the network topology itself. Run time on
 384 hardware is about 1.5 s of which roughly 125 ms are initial configuration and 278 ms are transmission of
 385 the input spike train. Post-processing the 1.2×10^7 received spikes (*fisch* and *grenade*) takes about 1.9 s,
 386 i.e., in the same order of magnitude as the actual hardware run. Peak memory consumption is reached
 387 during post-processing of results obtained after the hardware execution which corresponds to roughly 3
 388 times the minimum memory footprint of the recorded spike train. With this the stack is well suited to also
 389 handle experiments with high spike count without introducing a bottleneck.

Figure 8. Run time analysis of a PyNN-based experiment with large spike count. Population of 12 neurons is excited by a regular spike train with frequency of 1 MHz. The network is emulated for 1 s on hardware resulting in 1.2×10^7 spike events. The black line represents memory consumption during execution. Horizontal bars represent time consumption in software layers. The annotations in the legend present the individual run time of steps and percentage of the overall run time.



390 2.3.5 Integration into PyTorch

391 To enable access to BSS-2 for machine learning applications, we develop a thin wrapper layer to
 392 the PyTorch-API. This extension is called *hxtorch* and was introduced in Spilger et al. (2020) for non-
 393 spiking hardware operation emulating analog multiply-accumulate operations and compositions thereof.
 394 There, we build on top of the same signal-flow graph experiment description as for the spiking mode of
 395 operation, cf. section 2.3.2. Operations are mapped to the hardware size by using temporal serialization and
 396 physical concurrency. The PyTorch extension enhances this by automatic gradient calculation for training.
 397 Same as PyTorch, we implement a functional API in C++ wrapped to Python (e.g. `hxtorch.matmul`
 398 comparable to `torch.matmul`) and add modules/layers on top in Python (e.g. `hxtorch.nn.Linear`
 399 comparable to `torch.nn.Linear`). In contrast, our operations are quantized to the hardware-intrinsic
 400 digital resolution (5 bit unsigned activations, 6 bit weights plus sign bit and 8 bit signed results). Execution
 401 on the hardware is performed individually for each operation using the JIT execution, see section 2.3.2.

402 2.3.6 Calibration framework

403 On BSS-2, there are a multitude of voltages and currents controlling analog circuit behavior. While some
 404 of them can be set to default values, most of them require calibration in order to match experiment-specific
 405 target values and to counteract device-specific mismatch. Fundamentally, the calibration can be executed
 406 on a host computer or locally on chip, using the embedded processors. We provide the Python module
 407 *calix* to handle all aspects of the calibration process.

408 Model parameters are calibrated by iteratively adjusting relevant parts of the hardware configuration. As
 409 an example, the membrane time constant is controlled by a bias current: In order to calibrate the membrane
 410 time constant of all neurons, the neurons' membrane potentials are recorded while they decay back to their
 411 resting potential after an initial perturbation from the resting state. We can perform an exponential fit to the
 412 recorded voltage trace to determine the time constant and iteratively tweak the bias current to reach the
 413 desired target.

414 The calibration routine of each parameter is encapsulated using an object-oriented API providing a
 415 common interface. Mainly, two methods allow the iterative parameter search: one applies a parameter
 416 configuration to the hardware, while the other evaluates an observable to determine circuit behavior. An

417 algorithm calculates parameter updates during the iterative search. In each step, the measurement from the
418 calibration class is compared to the target value and the parameter set is modified accordingly.

419 A functional API is provided for commonly used sets of calibrations, for example for calibration of a
420 spiking leaky-integrate and fire (LIF) neuron. Technical parameters and multidimensional dependencies
421 are handled automatically as required in this case. This yields a simple interface for experimenters for
422 tweaking high-level parameters, while calibration routines for individual parameters remain accessible for
423 expert users.

424 The higher-level calibration functions save their results in a typed data structure, which contains the
425 related analog parameters and digital control bits. Further, success flags indicate whether the calibration
426 targets were reached within the available parameter ranges. These result structures can either directly
427 be applied to a hardware setup or serialized to disk. Application of serialized calibration is beneficial
428 compared to repeating the calibration in experiments due to decreased required time and improved digital
429 reproducibility.

430 Running the calibration on a host computer using `Python` allows for great flexibility in terms of
431 gathering observations from the chip. We can utilize all observables, including a fast ADC, which allows
432 performing fits to measured data – as sketched previously for the calibration of the membrane time constant.
433 While this direct measurement should yield the most accurate results, fitting to a trace for each neuron
434 takes a lot of time. Performing a full LIF neuron calibration takes a few minutes via the `Python` module.
435 And importantly, when scaling this approach to many chips, we need to scale the host computing power
436 accordingly.

437 In order to achieve better scalability, we can control the calibration from the embedded processors, directly
438 on chip, removing the host computer from the loop. However, this approach limits the observables to those
439 easily accessible to the embedded processor, the CADC and spike counters – performing a fit to an MADC
440 trace using the embedded processors would consume lots of runtime and potentially counteract benefits
441 of scaling. As a result, some calibrations have to rely on an indirect measurement of their observable.
442 Again using the neurons' membrane time constant as an example, we can consider the spike rate in a
443 leak-over-threshold setup. However, this introduces a dependency on multiple potentials being calibrated
444 beforehand.

445 Apart from the need for indirect measurements, on-chip and host-based calibration work similarly: An
446 iterative algorithm selects parameters, we configure them on chip and characterize their effects. Using the
447 embedded processors for configuring parameters and acquiring data from the two on-chip readouts is fully
448 supported and naturally faster than fetching them from a host computer. We use the SIMD CPUs' vector
449 units for parallel access to the synapse array and columnar ADCs. This is enabled by cross-compiler-support
450 (cf. section 3.3), by which both the scalar unit and vector unit are integrated and accessible from the C++
451 language.

452 We provide routines for on-chip calibration, which allow all LIF neuron parameters to be calibrated
453 in approximately half a minute, with this number staying constant even when considering large systems
454 comprising many individual chips. Similar to the host-based calibration API, *calix* exposes these on-chip
455 routines as conveniently parametrized functions that can be called within any experiment. Their runtime is
456 mostly limited by waiting for configured analog parameters to stabilize before evaluating the effects on the
457 circuits.

458 2.4 Platform Operation

459 Over the past decade neuromorphic systems evolved from intricate lab setups towards back ends for the
460 more comfortable execution of spiking neural networks (Indiveri et al., 2011; Furber et al., 2012; Benjamin
461 et al., 2014; Davies et al., 2018; Pehle et al., 2022). One major step along this development path is to
462 provide users with seamless access to the systems.

463 Small scale prototype hardware is often connected to a single host machine, e.g., via USB. This is also a
464 common usage mode for different neuromorphic hardware. To access these devices, users have to have
465 (interactive) access to the particular machine the hardware is connected to. This limits the flexibility of the
466 user and is an operational burden as the combination of neuromorphic hardware and host machine has to
467 be maintained. While this tightly coupled mode of operation is sufficient during commissioning and initial
468 experiments, it is not robust enough for higher work-loads and flexible usage.

469 An improvement to the situation sketched above is using a scheduler, e.g., SLURM (Yoo et al., 2003),
470 where users can request a resource, e.g., a specific hardware setup, and the jobs get launched on the
471 matching machine with locally attached hardware. This is the typical mode of access also used for other
472 accelerator-type hardware, e.g., GPU clusters. However, this batch driven way is not always ideal as it often
473 requires accounts on the local compute cluster and does not allow for easy interactive usage. In addition,
474 traditional compute load schedulers optimize for throughput and not latency, therefore the scheduling
475 overhead can be significant especially for hardware that is fast and experiments that are short. In the latter
476 case, job execution rates of the order of Hz and faster are required.

477 Another downside of using a traditional scheduler is that hardware resources are not efficiently utilized
478 when multiple users want to use the same hardware resources at the same time. Therefore, we developed
479 the micro scheduler *quiggeldy* that exposes access to the hardware directly via a network connection,
480 but still manages concurrent access from different users. It decouples the hardware utilization from the
481 user’s surrounding computations such as experiment preparation, updates in iterative workflows or result
482 evaluation. For this to work runtime control, configuration, input stimulus as well as output data must be
483 serializable which is facilitated via cereal (Grant and Voorhies, 2017). The inter-process communication
484 between the user software and the micro scheduler is done with RCF (Delta V Software, 2020). When a
485 user requests multiple hardware runs, it is checked whether certain already performed parts can be omitted,
486 e.g., resets or re-initializations. Experiment interleaving between multiple users is also supported as the
487 initialization state is tracker for each user and is automatically applied when needed.

488 Having the correct software environment for using neuromorphic hardware is also a major challenge.
489 Nowadays, software vendors often provide a container image that includes the appropriate libraries.
490 However, this approach does not necessarily yield well specified and traceable dependencies, but only
491 a “working” black-box solution. We overcome this downside by using the Spack (Gamblin et al., 2015)
492 package manager with a meta-package that explicitly tracks all software dependencies and their version
493 needed to run experiments on and develop for the neuromorphic hardware. An automatically built container
494 embedding the Spack installation enables encapsulation and eased distribution. This Spack meta-package is
495 also used for the EBRAINS’ JupyterLab service and will eventually be deployed to all HPC sites involved
496 in EBRAINS (ebr, 2022). The latter will facilitate multi-site workflows involving neuromorphic hardware
497 and traditional HPC.

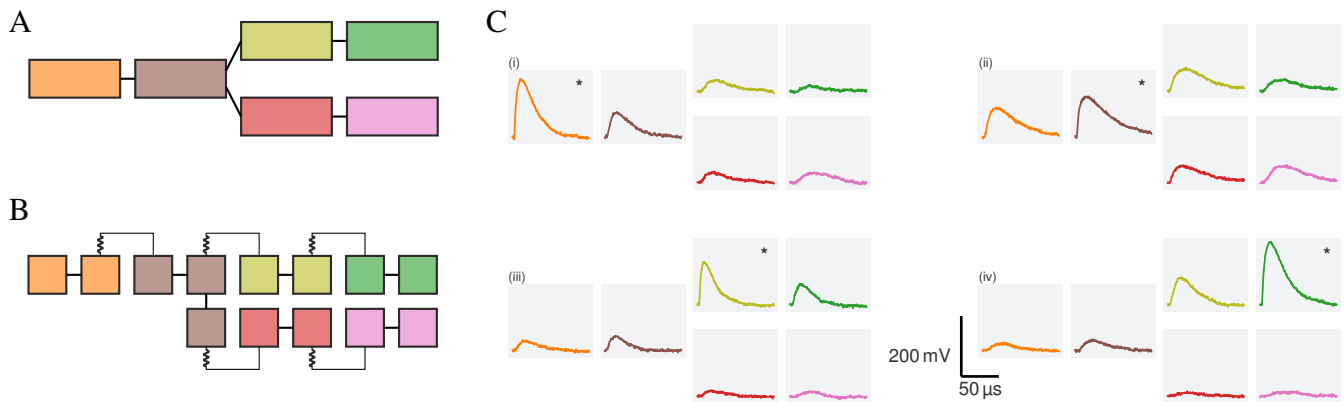


Figure 9. Pulse propagation along a dendrite which branches into two sub-branches. **(A)** Each branch is modeled by two compartments (rectangles). Different compartments are connected via resistors (lines). **(B)** Hardware configuration: neuron circuits (squares) are arranged in two rows on BSS-2, compare fig. 1 (D). Each compartment is represented by at least two neuron circuits. Circuits which form a single compartment are directly connected via switches (straight lines); compartments are connected via resistors. For details see Kaiser et al. (2021). **(C)** Membrane responses to synaptic input: we inject synaptic input at four different compartments; the compartment at which the input is injected is marked by a *. The membrane traces of the different compartments are arranged as in sub-figure (A). For the top left quadrant (i) the input is injected in the first compartment and decreases in amplitude while it travels along the chain. The response in both branches is symmetric. A similar behavior can be observed when the input is injected in the second compartment, (ii). Due to the symmetry of the model, we only display membrane responses for synaptic input to the upper branch. When injecting the input in the first compartment of the upper branch (iii) the input causes a noticeable depolarization within the same branch and the main branch but does not cause a strong response in the lower sister branch. Note: all values are given in the hardware domain.

3 APPLICATIONS

498 In this section, we show-case a range of applications of BSS-2. Each application involves use of unique
 499 hardware features or modes of operation and motivates parts of the software design.

500 First, we describe biological multi-compartmental modelling in section 3.1 concluding in the development
 501 of an API for structural neurons. Continuing, functional modelling with spiking neural network (SNN) is
 502 demonstrated for a pattern-generation task in section 3.2, which leads to embedding of spiking BSS-2 usage
 503 into the machine learning framework PyTorch and involves host-based training as well as local learning
 504 on the SIMD CPUs. Then, embedded operation, where the SIMD CPUs are the experiment orchestrator
 505 of BSS-2, is displayed and their implications detailed in section 3.3. Following, the non-spiking mode of
 506 operation implementing ANNs and its PyTorch interface is characterized in section 3.4. Afterwards, user
 507 adoption and platform access to BSS-2 is shown in section 3.5. Finally, application of the software stack
 508 for hardware co-simulation, co-design and verification is portrayed in section 3.6.

509 3.1 Biological Modeling Example

510 BSS-2 aims to emulate biological inspired neuron models. Most neurons are not simple point-like
 511 structures but possess intricate dendritic structures. In recent years, the research interest in how dendrites
 512 shape the output of neurons has increased (Major et al., 2013; Gidon et al., 2020; Poirazi and Papoutsis,
 513 2020). As a result, BSS-2 incorporates the possibility to emulate multi-compartmental neuron models in
 514 addition to the AdEx point-neuron model (Aamir et al., 2018; Kaiser et al., 2021).

515 In the following, we use a dendritic branch, which splits into two sub-branches, to illustrate how multi-
516 compartmental neuron models are represented in our system, cf. fig. 9. At first, we look at a simplified
517 representation of the model, subfigure (A). The main branch consists of two compartments, connected via
518 a resistance; at the second compartment, the branch splits in two sub-branches, which themselves consist
519 of two compartments each. On hardware this model is replicated by connecting several neuron circuits via
520 switches and tunable resistors, cf. fig. 9 (B). Each compartment consists of at least two neuron circuits,
521 directly connected via switches, compare colors in subfigure (A) and (B). With the help of a dedicated
522 line at the top of the neuron circuits these compartments can then be connected via resistors to form the
523 multi-compartmental neuron model; for more details see Kaiser et al. (2021).

524 In software, the `AtomicNeuron` class stores the configuration of a single neuron circuit and therefore can
525 be used to configure the switches and resistors as desired. As mentioned in section 2.3.4, the `HXNeuron` ex-
526 poses this data structure to the high-level interface `PyNN`, allowing users to construct multi-compartmental
527 neuron models in a known environment. However, it is cumbersome and error-prone to set individual
528 switches. As a consequence, we implement a dictionary-like hierarchy on top of the `AtomicNeuron`,
529 called `LogicalNeuron` in the logical abstraction layer, cf. section 2.2.

530 We use a builder pattern approach to construct these logical neurons: the user creates a neuron morphology
531 by defining which neuron circuits constitute a compartment and how these compartments are connected.
532 Upon finalization of the builder, the correctness of the neuron model configuration of the neuron model
533 is checked; if the provided configuration is valid, a `LogicalNeuron` is created. This `LogicalNeuron`
534 stores the morphology of the neuron as well as the configuration of each compartment.

535 The coordinate system of the BSS-2 software stack, cf. section 2.2.2, allows to place the final logical
536 neuron at different locations on the chip (Müller et al., 2020a). This is achieved by saving the relation
537 between the different neuron circuits defining the morphology in relative coordinates. Once the neuron is
538 placed at a specific location on the chip, the relative coordinates are translated to absolute coordinates.

539 Currently, the logical neuron is only exposed in the logical abstraction layer. In future work, it will
540 be integrated in the `PyNN` API of the BSS-2 system. This will – for instance – allow to easily define
541 populations of multi-compartmental neurons and connections between them.

542 3.2 Functional Modeling Example

543 The BSS-2 system enables energy efficient and fast SNN implementations. Moreover, the system’s
544 embedded SIMD CPU enables highly parallelized on-chip learning with fast access to observables and thus,
545 promises to benefit the computational neuroscience and machine learning community in terms of speed and
546 energy consumption. We demonstrate functional modeling on the BSS-2 system with a pattern-generation
547 task using recurrent spiking neural networks (RSNNs) with an input layer, a recurrent layer and a single
548 readout neuron. The recurrent layer consists of 70 LIF neurons $\{j\}$ with membrane potential v_j^t , receiving
549 spike trains x_i^t from 30 input neurons $\{i\}$. Neurons in the recurrent layer project spike events z_j^t onto the
550 single leaky-integrate readout neuron with potential y^t .

551 RSNNs are commonly trained using backpropagation through time (BPTT) by introducing a variety of
552 surrogate gradients taking account of the discontinuity of spiking neurons (Bellec et al., 2020; Zenke and
553 Ganguli, 2018; Shrestha and Orchard, 2018). However, as BPTT requires knowledge of all network states
554 along the time sequence in order to compute weight updates (backwards locking), it is not just considered
555 implausible from a biological perspective, but also unfavourable for on-chip learning, which effectively
556 enables high scalability due to local learning. Therefore, we utilize e-prop learning rules (Bellec et al.,

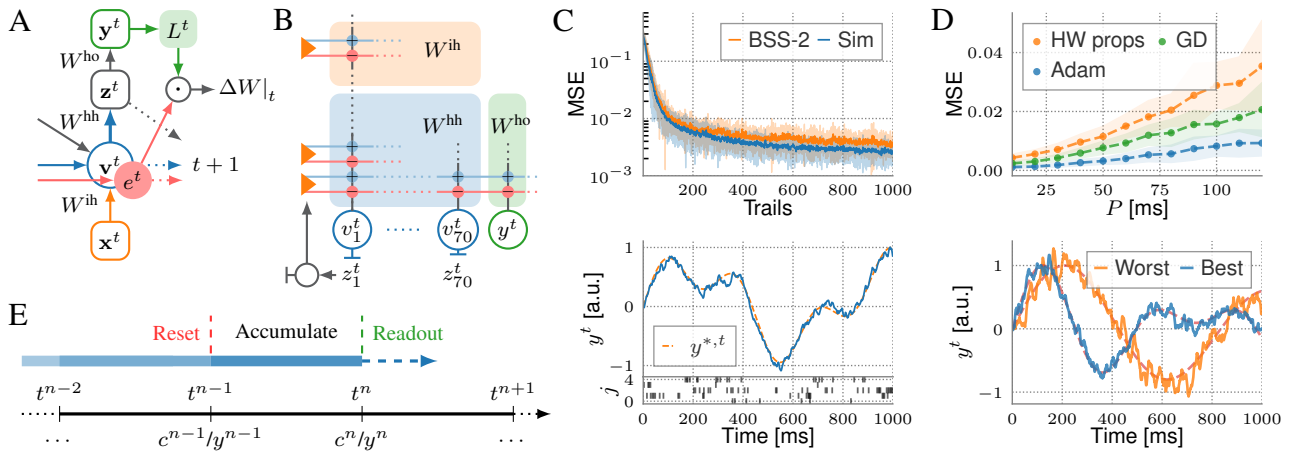


Figure 10. (A) Computational graph of an RSNN for one time step. The contribution to the weight update is computed by merging learning signals L_j^t with eligibility traces e_{ji}^t . (B) Representation of the RSNN on the BSS-2 system using signed synapses. Inputs and recurrent spike trains are routed to the corresponding synapse drivers via the crossbar. (C) s-prop training on hardware. The upper plot depicts the evolution of the mean squared error (MSE) while training the BSS-2 system in-the-loop, where the experiment is executed on BSS-2 and weight updates are computed on the host computer, in comparison to training with the network simulated in software, incorporating basic hardware properties (Sim). In both cases the weights are optimized using the Adam optimizer (Kingma and Ba, 2014). The learned analog membrane trace of the readout neuron after training BSS-2 for 1000 epochs is exemplified in the lower plot, aligned to the spike trains z_j^t of the first five out of 70 recurrent neurons. (D) NASProp simulations. The upper plot depicts the MSE over the update period P after training with Adam in comparison to a training with gradient descent (GD) and a training taking additional hardware properties (noise, weight saturation, etc.) into account (HW props). Optimization with pure GD mimics weight updates computed by the SIMD CPU while on-chip learning. The lower plot shows the worst and best learned readout traces of the target pattern ensemble in simulation. (E) Timing of NASProp weight updates. For each update n at t^n , the correlations c_{ji}^n are merged with the learning signals L_j^n by incorporating the membrane trace y^n .

557 2020), where the gradient for BPTT is factorized into a temporal sum over products of so-called learning
 558 signals L_j^t and synapse-local eligibility traces e_{ji}^t . While the latter accumulates all contributions to the
 559 gradient that can be computed forward in time, the first depends on the network's error and still requires
 560 BPTT. However, Bellec et al. (2020) provide suitable approximations for L_j^t , allowing computing the
 561 weight updates online (fig. 10A). Such learning rules are favorable for the BSS-2 system, as the SIMD
 562 CPU can compute the weight updates locally while the network is emulated in parallel.

E-prop-inspired learning on the BSS-2 system is enabled by adapting Bellec et al. (2020, Eq. (28)). Here we replace the membrane potentials v_j^t in e_{ji}^t with the post-synaptic recurrent spike train z_j^t ,

$$e_{ji}^{t+1} \rightarrow z_j^{t+1} \cdot \mathcal{F}_\alpha(z_i^t) := \hat{e}_{ji}^{t+1}, \quad \Delta W_{ji}^{\text{hh}} = -\eta \sum_t L_j^t \mathcal{F}_\kappa(\hat{e}_{ji}^t), \quad (1)$$

563 where \mathcal{F}_x is an exponential filter with decay constant x . The update rule for input weights, derived in
 564 Bellec et al. (2020), is adapted accordingly. The equation for output weights remains untouched. With the
 565 readout neuron's membrane trace y^t and an MSE loss measuring the error to a target trace $y^{*,t}$, the learning
 566 signals are $L_j^t = W_j^{\text{ho}}(y^t - y^{*,t})$. Since this learning rule propagates only spike-based information over
 567 time we refer to it as *s-prop*.

568 Finally, we approach s-prop learning with BSS-2 in the loop (cf. section 2.1). For this, the network,
 569 represented by PyTorch parameters $W^{\text{ih}, \text{hh}, \text{ho}}$, is mapped to a hardware representation (see fig. 10B) via
 570 *hxtorch* (see section 2.3.5), forwarding a spike tensor on-chip. Inherently, *grenade* (see section 2.3.2)
 571 applies a routing algorithm, finds a graph-based experiment description and executes it on hardware for a
 572 given time interval. The routing algorithm allocates two adjacent hardware synapses for one signed synapse
 573 weight in software, one excitatory and one inhibitory. Further, *grenade* records the MADC-sampled readout
 574 trace y^t and the recurrent spike trains \mathbf{z}^t . Both observables are returned as PyTorch tensors for weight
 575 optimization on the host side. Experiment results are displayed in fig. 10C.

Implementing s-prop on-chip requires the SIMD CPU to know and process explicit spike-times. As this comes with a high computational cost, the correlation sensors are utilized to emulate approximations of the spike-based eligibility traces \hat{e}_{ji}^t in analog circuits, thereby freeing computational resources on the SIMD CPU. The correlation sensors model the eligibility traces under nearest-neighbor approximation (Friedmann et al., 2017) and are accessed by the SIMD CPU as an entity c_{ji}^n , accumulated over a period P . Hence, the time sequence is split into N chunks of size P and weight updates on the SIMD CPU are performed at times $t^n = nP + \tilde{t}$, with $n \in \mathbb{N}_0^{<N}$ (cf. fig. 10E) and $\tilde{t} \in [0, P)$ a random offset,

$$\Delta \bar{W}_{ij}^{\text{ih/hh}} = -\eta \sum_n L_j^n \mathcal{F}_{\hat{\kappa}} \left(c_{ji}^{\text{ih/hh}, n} \right) \quad \text{and} \quad \Delta \bar{W}_{kj}^{\text{ho}} = -\eta \sum_n \left(y_k^n - y_k^{*,n} \right) \mathcal{F}_{\hat{\kappa}} \left(\zeta_j^n \right), \quad (2)$$

576 with $\hat{\kappa} = \exp(-P/\tau_m)$ and ζ_j^n being the recurrent spike count in interval n . Due to the updates rules'
 577 accumulative nature, we refer to them as neuromorphic accumulative spike propagation (NASProp).
 578 Simulations in fig. 10D verify that NASProp endows RSNNs with the ability to solve the pattern-generation
 579 task reasonable well.

580 NASProp's SIMD CPU implementation effectively demonstrates full on-chip learning on the BSS-2
 581 system. In high-level software, on-chip learning is implemented in a PyTorch model, defined in *hxtorch*,
 582 holding parameters for the network's projections. Its `forward` method implicitly executes the experiment
 583 on the BSS-2 system for a batch of input sequences. Currently, this model learning on-chip serves as a
 584 mere black box for the specific network at hand with a static number of layers, as for on-chip spiking
 585 networks the network's topology needs to be known upon execution. Therefore, this approach is considered
 586 a first step from common PyTorch models to spiking on-chip models.

587 As for in-the-loop learning, on forwarding a batch of inputs sequences, *grenade* maps the software
 588 network to a hardware representation with signed synapses and configures the chip accordingly. Moreover,
 589 before executing a batch element, *grenade* starts the plasticity kernel on the SIMD CPU, computing weight
 590 updates in parallel to the network's emulation. The plasticity rule implementation relies on *libnux* (cf.
 591 section 2.2.3) and utilizes the VU extension for accessing hardware observables (e.g. c_{ji}^n and y^n) and
 592 computing weight updates row-wise in parallel, thereby fully exploiting the system's speed up factor.

593 In *hxtorch*, learning parameters are configured in a configuration object exposed to Python, which is
 594 injected to *grenade* and passed to the SIMD CPU before batch execution on hardware begins. As different
 595 projections in the network have different update rules, relying on population-specific observables, the
 596 network's representation on hardware (cf. fig. 10B) is communicated to the SIMD CPU. This allows
 597 for identifying signed hardware synapses and neurons with projections and populations on the SIMD
 598 CPU. Finally, before each batch element is executed, *grenade* has the ability to write trial-specific details
 599 onto the SIMD CPU (e.g. random offset \tilde{t} and the synapse row to perform updates for). Hence, smooth
 600 on-chip learning is granted by reliable communication between little-endian host engine and the embedded

601 big-endian SIMD CPU. For serialization of information from and to the SIMD CPU we deploy the C++
602 header-only library *bitsery* (Vinkelis, 2020), allowing for seamless transmission of objects between systems
603 of differing endianness.

604 Due to changing hardware weights during on-chip training, the adjusted weights are reverse mapped to
605 the software representation and stored in the network's parameter tensors. Therewith we utilize PyTorch's
606 native functionality to load and store network parameters. Reverse network mapping is implemented in the
607 *hxtorch* on-chip-learning model by accessing the hardware routing result and is performed implicitly in the
608 model's `forward` method after experiment execution.

609 Successful implementations of plasticity rules for on-chip learning are facilitated by providing trans-
610 parency of SIMD CPU programs by means for tracing and recording data. To that end, *libnux* (cf.
611 section 2.2.3) facilitates logging of any information into a dedicated SIMD CPU memory region, easily
612 accessed from the host engine. Moreover, logging can be redirected to the FPGA-controlled dynamic
613 random-access memory (DRAM), effectively allowing extensive logging of whole learning processes and
614 hardware observables.

615 3.3 Embedded Operation

616 Apart from operating BSS-2 tightly coupled to a host computer, the integrated microprocessors can
617 act as system controllers. They can orchestrate the control flow of the experiment and undertake tasks
618 within it. These tasks may include calibration routines, virtual environment simulation or optimizer loops.
619 Embedding them in proximity to the neural network core yields latency and data-locality advantages. In the
620 following, we describe three exemplary experiments that make exhaustive use of the embedded processors
621 as system controllers.

622 First, Wunderlich et al. (2019) introduce an embedded environment simulation of a simplified version of
623 the Pong video game on the SIMD CPU, see left panel in fig. 11. One of the two involved agents plays
624 optimally by design, the other one is represented by a SNN on BSS-2. During the experiment, the latter is
625 trained on-chip using a reward-based spike timing dependent plasticity (STDP) rule. This set-up therefore
626 unites the control flow, virtual environment simulation and learning rule within a single program running
627 on the integrated processors.

628 Second, Stradmann et al. (2021) describe the application of the BSS-2 system for inference of ANNs that
629 detect atrial fibrillation in medical electrocardiogram (ECG) data. Targeting applications in energy efficient
630 devices, they aim for as little periphery as possible and therefore let the embedded processors orchestrate
631 all classification routines. The resulting tight loop between the analog inference engine and digital data in-
632 and outputs allows for low classification latencies and high throughput of more than 3600 ECG traces per
633 second.

634 Third, Schreiber et al. (2022) presents the emulation of an insect model with strong biological inspiration
635 on BSS-2. The simplified brain model is embedded into an agent that is fed with stimuli from a simulated
636 environment, see right panel in fig. 11. While the neural network is emulated as a SNN within the analog
637 core, the agent itself as well as its virtual environment are both simulated on the SIMD CPU. The authors
638 specifically challenge the virtual insects with a simple path integration task: As depicted in the right panel
639 of fig. 11, a simulated swarm-out phase is followed by a period of free flight, where the agent is supposed to
640 return to its nest. The complexity of this task and the comparably low number of involved neurons requires
641 precisely controlled dynamics, which they achieve by integrating experiment specific on-chip calibration
642 routines directly on the SIMD CPUs (cf. section 2.3.6).

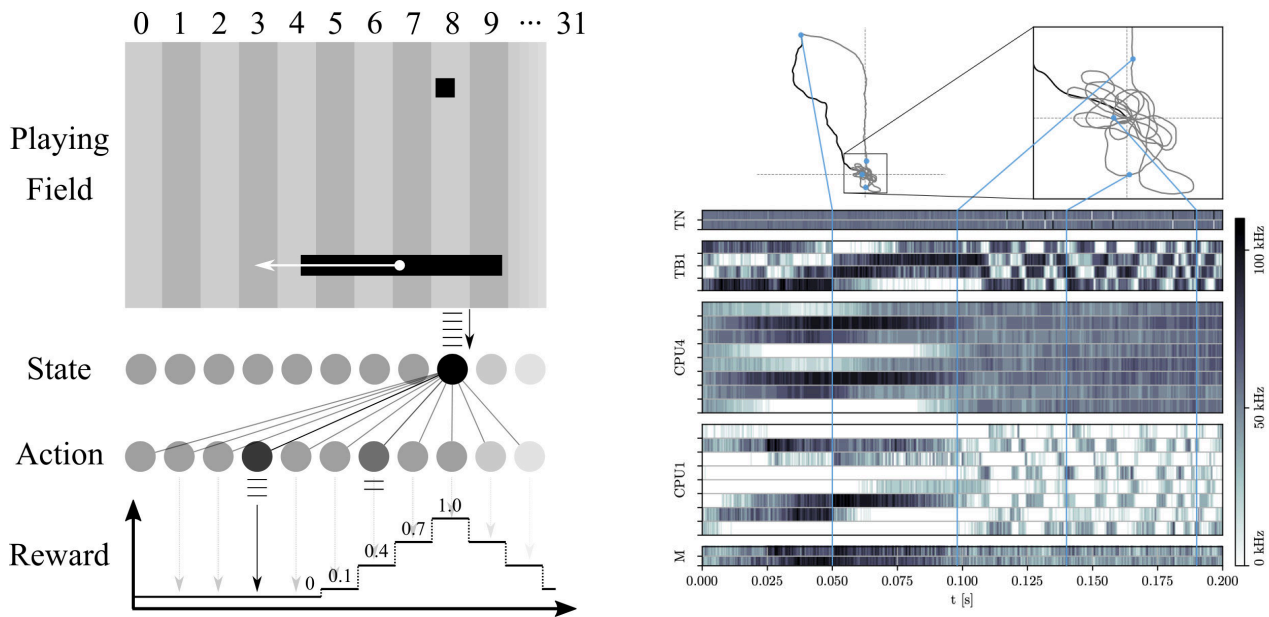


Figure 11. (left) Reinforcement learning: the chip implements a spiking neural network sensing the current ball position and controlling the game paddle. It is trained via a reward-based STDP learning rule to achieve almost optimal performance. The game environment, the motor command and stimulus handling, the reward calculation and the plasticity is performed by a C++ program running on the on-chip processor. Figure taken from Wunderlich et al. (2019). (right) Recording of a virtual insect navigating a simulated environment. The top panels show the forced swarm-out path in black. During this phase, the SNN emulated by the analog neuron and synapse circuits on BSS-2 perform path integration. Afterwards, the insect flies freely and successfully finds its way back to the starting point and circles around it (gray trajectory). The bottom panel shows the neuronal activity during the experiment. The environment simulation as well as the interaction with the insect is performed by a C++ program running on the on-chip processor. Figure taken from Pehle et al. (2022).

643 Supporting these complex experiments on the embedded processors and their interaction with the
 644 controlling host computer poses specific requirements to the BSS-2 OS. Especially, a cross-compilation
 645 toolchain for the SIMD CPU is required.

646 As described in section 2.2.3, we therefore provide a cross-compiler based on `gcc` (GNU Project, 2018),
 647 which in addition to the processor’s scalar unit also integrates its custom vector unit in C++ (Müller et al.,
 648 2020a). Additional hardware specific functionality is encapsulated in the support library *linux*. It abstracts
 649 access to configuration data and observables in the analog neural network core, like synaptic weights or
 650 correlation measurements. The exchange of such data with the host is facilitated by integration of the lean,
 651 cross-platform binary serialization library *bitsery* (Vinkelis, 2020).

652 For execution, the compiled programs need to be placed in system memory — in case of BSS-2, each
 653 SIMD CPU has direct access to 16 kB SRAM. For a complete calibration routine or complex locally
 654 simulated environments, this may not suffice. We therefore utilize the controlling FPGA as memory
 655 controller: It allows the on-chip processors to access externally connected DRAM with significantly larger
 656 capacity at the cost of higher latency. Programs for the embedded processor can place instructions and data
 657 onto both the internal SRAM and the external memory via compiler attributes. This allows fine-grained
 658 decisions about the access-latency requirements of specific instruction and data sections.

659 Similar to experiments designed for operation from the host system, embedded experiments often
 660 require reconfiguration of parts of BSS-2. The hardware abstraction layer introduced in the BSS-2 OS

661 (cf. section 2.2.2) has therefore been prepared for cross-compilation on the embedded processors. As a
 662 result, the described container and coordinate system can be used in experiment programs running on the
 663 on-chip SIMD CPUs.

664 3.4 Artificial Neural Networks

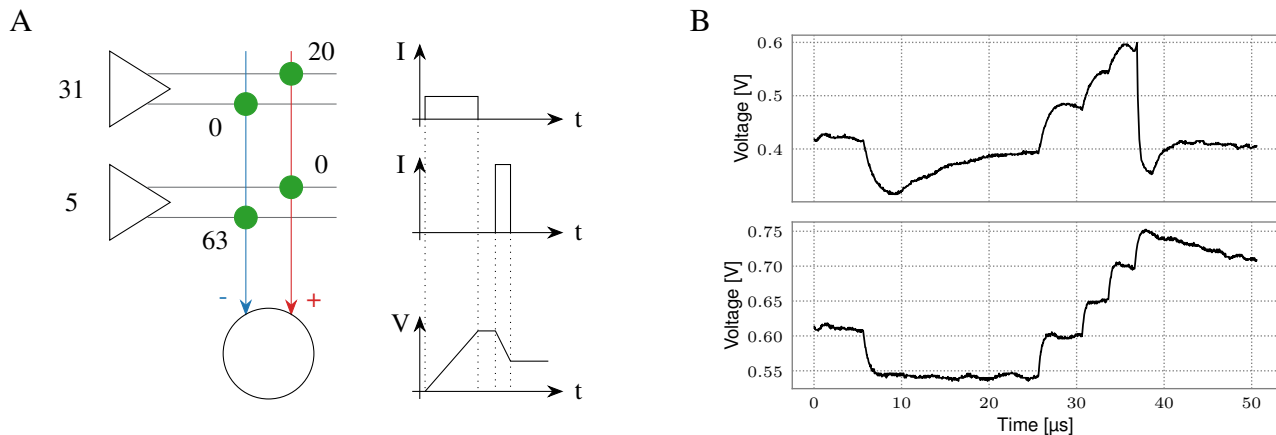


Figure 12. Matrix-vector multiplication for ANN inference. **(A)** Scheme of a multiply-accumulate operation. Vector entries are input via synapse drivers (left) in 5 bit resolution. They are multiplied by the weight of an excitatory or inhibitory synapse, yielding 6 bit plus sign weight resolution. The charge is accumulated on neurons (bottom). Figure taken from Weis et al. (2020). **(B)** Comparison between a spiking (top) and an integrator (bottom) neuron. Both neurons receive identical stimuli, one inhibitory and multiple excitatory inputs. While the top neuron shows a synaptic time constant and a membrane time constant, the lower is configured close to a pure integrator. We use this configuration for ANN inference. Please note that for visualization purposes the input timing (bottom) has been slowed to match the SNN configuration (top). The integration phase typically lasts less than 2 μ s.

665 The BSS-2 hardware supports a non-spiking operation mode which supports artificial neural networks
 666 (ANNs) implementing multiply-accumulate (MAC) operations (Weis et al., 2020). The operation within
 667 the analog core is sketched in fig. 12A. Each entry in the vector operand stimulates one or two rows of
 668 synapses, when using unsigned or signed weights, respectively. The activations have an input resolution of
 669 5 bit, controlling the duration of synapses' activation. Similar to the spiking operation, synapses emit a
 670 current pulse onto the neurons' membranes depending on their weight, which has a resolution of 6 bit. We
 671 implement signed weights by combining an excitatory and an inhibitory synapse into one logical synapse.
 672 Once all entries in the input vector have been sent to the synapses, the membrane potential resembles the
 673 result of the MAC operations. It is digitized for all neurons in parallel using the CADC, yielding an 8 bit
 674 result resolution.

675 As a user interface, we have developed an extension to the PyTorch machine learning framework (Paszke
 676 et al., 2019), *hxtorch* (Spilger et al., 2020). It partitions ANN models into chip-sized MAC operations
 677 that are executed on hardware using *grenade*, see section 2.2.5. Apart from a special MAC program used
 678 for each multiplication, the majority of code is shared between spiking and non-spiking operation. With
 679 the leak term disabled, the neurons' membranes represent the integrated synaptic currents, as shown in
 680 fig. 12B. As the MAC operation lacks any real-time requirements, it is executed as fast as possible to
 681 optimize energy efficiency. In terms of circuit parameterization, this means we choose a small synaptic
 682 time constant in order for the membrane potential to stabilize quickly. Therefore, a subset of the existing

683 spiking calibration routines can be reused here, cf. section 2.3.6. There is only one additional circuit – the
 684 encoding of input activations to activation times in synapse drivers – that needs to be calibrated.

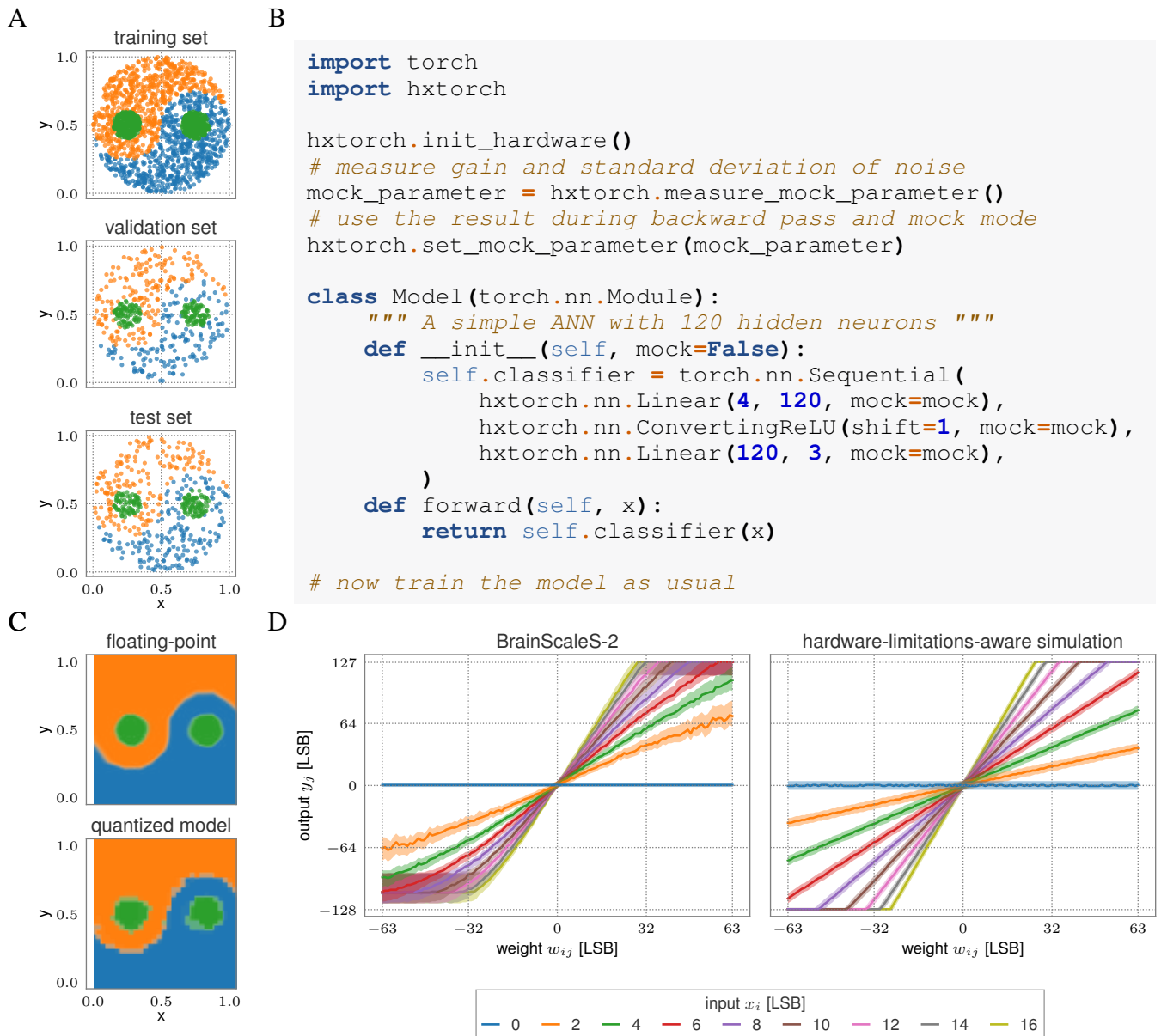


Figure 13. (A) The Yin-Yang dataset (Kriener et al., 2021) used for the experiment. (B) Hardware initialization and model description with *hxtorch*. (C) Network response of the trained model depending on the input. Top: 32 bit floating-point precision; bottom: quantized model on BSS-2 (5 bit activations, 6 bit plus sign weights). (D) Output of the MAC operation on BSS-2 (left) compared to the linear approximation (right). The solid line indicates the median, the colored bands contain 95% of each neuron’s outputs across 100 identical MAC executions.

685 Defining an ANN model in *hxtorch* works similar to PyTorch: The *hxtorch* module provides linear and
 686 convolutional layer classes as a replacement for their PyTorch equivalents. We introduce a few additional
 687 parameters controlling the specifics of hardware execution, e.g. the time interval between sending successive
 688 entries in the input vector to the synapse matrix, or the option to repeat the vector for efficacy scaling. This

689 enables the user to optimize saturation effects when driving the input currents as well as the gain of the
690 MAC operation for the particular experiment. For both we provide default values as a starting point. The
691 activation function `ConvertingReLU` additionally converts signed 8 bit output activations into unsigned
692 5 bit input activations for the following layer by a bitwise right shift.

693 Trained deep neural network models can be transferred to BSS-2 by first quantizing them with PyTorch
694 and subsequently mapping their weights to the hardware domain. For quantization, we need to consider the
695 intrinsic gain factor of the hardware MAC operation.

696 Figure 13 shows an example application of a deep neural network with BSS-2, using the yin-yang dataset
697 from Kriener et al. (2021). One of the three classes – yin, yang, or dot – are to be determined from four
698 input coordinates $(x, y, 1 - x, 1 - y)$. The network is first trained with 32 bit floating point accuracy using
699 PyTorch, achieving 98.9 % accuracy. After quantizing with PyTorch to the hardware resolution of 5 bit
700 activations and 6 bit plus sign weights, this drops to 94.0 %. Porting the model to BSS-2, after running a
701 few epochs of hardware-in-the-loop training, an accuracy of 95.8 % is finally reached.

702 In addition to running the ANN on the BSS-2 hardware, a hardware-limitations-aware simulation is
703 available. It can be enabled per layer via the `mock` parameter (see fig. 13B). For mock mode, we simply
704 assume a linear MAC operation, using a hardware-like gain factor. To investigate possible effects of
705 the analog properties of the BSS-2 hardware on the inference and training, additional Gaussian noise of
706 the accumulators and multiplicative fixed-pattern deviations in the weight matrix can be simulated. The
707 comparison with actual hardware operation shown in fig. 13 D illustrates how this simple model already
708 captures the most dominant non-linearities of the system. More sophisticated software representations
709 that embrace second-order effects across multiple hardware instances have been proposed by Klein et al.
710 (2021). They have shown how pre-training with faithful software models can significantly decrease
711 hardware allocation time while at the same time increasing classification accuracy compared to plain
712 hardware-in-the-loop training.

713 3.5 User Adoption and Platform Access

714 The BSS-2 software stack aims to enable researchers to exploit the capabilities of the novel neuromorphic
715 substrate. Support for common modeling interfaces like PyNN and PyTorch provides a familiar entry point
716 for a wide range of users. However, not all aspects of the hardware can fully be abstracted away, requiring
717 users to familiarize themselves with unique facets of the system. To flatten the learning curve several
718 tutorials —verified in continuous integration (CI) as ‘executable’ documentation— as well as example
719 experiments are provided¹⁷. They range from introducing the hardware via single neuron dynamics to
720 learning schemes like plasticity rate coding. In addition to the scientific community, they also target
721 students, for example exercises accompanying a lecture about Brain Inspired Computing and hands-on
722 tutorials.

723 A convenient entry point to explore novel hardware are interactive web-based user interfaces. That is
724 why we integrated the BSS-2 system into the EBRAINS Collaboratory¹⁸ (ebr, 2022). The Collaboratory
725 provides a dynamic VM hosting on multiple HPC sites for Jupyter notebooks running in a comprehensive
726 software environment. An BSS-2-specific experiment service manages multi-user access to the hardware
727 located in Heidelberg utilizing the *quiggeldy* micro scheduler, see section 2.4. It allows for seamless
728 interactive execution of experiments running on hardware with execution rates of over 10 Hz. This, for

¹⁷ The tutorials and example experiments are available at <https://github.com/electronicvisions/brainscales2-demos>

¹⁸ Platform access is available via <https://ebrains.eu>

729 example, was utilized during hands-on tutorials at the NICE 2021 conference (nic, 2021). The execution
 730 rates of that demonstration are shown in fig. 14.

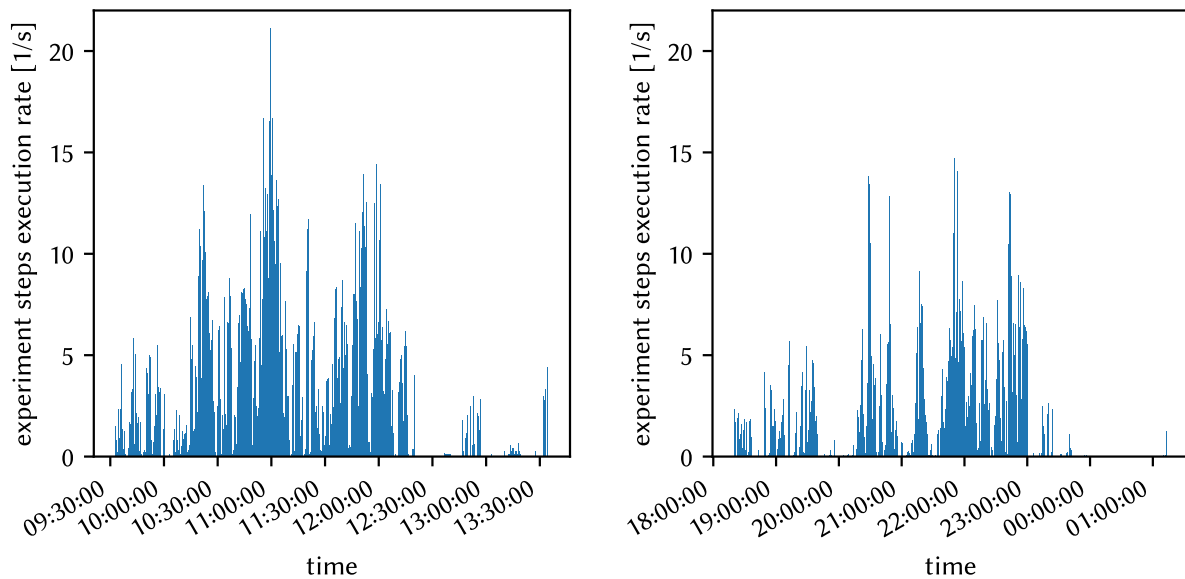


Figure 14. Rate of executed experiment-steps via *quiggeldy* during the two BSS-2 hands-on tutorials at NICE 2021. Experiments were distributed among eight hardware setups. In total there were 86 077 hardware runs executed.

731 Furthermore, EBRAINS has begun to provide a comprehensive software distribution that includes typical
 732 neuroscientific software libraries next to the BSS-2 client software. As of now, this software distribution
 733 has been already deployed at two HPC centers and work is under way to extend this to all sites available
 734 in the EBRAINS community. Leaving interactive demos aside, this automatic software deployment will
 735 simplify multi-site workflows significantly —including BSS-2 systems— as the scientist is not responsible
 736 for software deployment anymore.

737 3.6 Hardware/Software Co-Development

738 The BSS-2 platform consists of two main hardware components: the ASIC implementing an analog
 739 neural network core and digital periphery, as well as an FPGA used for experiment control and digital
 740 communication. Development of these hardware components is primarily driven by simulations of their
 741 analog and digital behavior, where — especially in the case of the ASIC — solid pre-fabrication test
 742 strategies need to be employed. Given the complexity of the system, integration tests involving all
 743 subsystems are required to ensure correct behavior.

744 Replicating the actual hardware systems, the setup for these simulated integration tests pose very similar
 745 requirements on the configuration and control software. The BSS-2 OS therefore provides a unified interface
 746 to both, circuit simulators and hardware systems. For the connection to the simulators, we introduce an
 747 adapter library (*flange*) as an optional substitution for the network transport layer. Implementing an
 748 additional *hxcomm* back-end, *flange* allows for the transparent execution of hardware experiments in
 749 simulation.

750 This architecture enables various synergies between hardware and software development efforts — specif-
 751 ically, co-design of both components already in early design phases. On system level, this methodology

752 helps to preempt interface mismatch between components of various different subsystems. Positive implications for software developers include the possibility of very early design involvement as well as enhanced
753 debug information throughout the full product life cycle: Having simulation models of the hardware
754 components of the system allows for the inspection of internal signals within the FPGA and ASIC during
755 program runtime. In particular, we have made use of this possibility during the development of a compiler
756 toolchain for the embedded custom SIMD microprocessors, where the recording of internal state helps
757 to understand the system's behavior. Hardware development, on the other hand, strongly profits from
758 software-driven verification strategies and test frameworks. BSS-2 OS especially allows to run the very
759 same test suites on current hardware as well as simulations of future revisions. These shared test suites
760 are re-used across all stages of the platform's life cycle for multiple hardware generations, therefore ever
761 accumulating verification coverage.
762

4 DISCUSSION

763 This work describes the software environment for the latest BrainScaleS (BSS) neuromorphic architecture
764 (Pehle et al., 2022): the BrainScaleS-2 (BSS-2) operating system. In Müller et al. (2020b) we introduced
765 the operating system for the BrainScaleS-1 (BSS-1) wafer-scale neuromorphic hardware platform. New
766 basic concepts of the second-generation software architecture were described in Müller et al. (2020a).
767 For example, we introduced a concise representation of “units of configuration” and “experiment runs”
768 supporting asynchronous execution by extensive usage of ‘future’ variables. Key concepts already existing
769 in BSS-1 —e.g., the type-safe coordinate system— were extended for BSS-2. In particular, the systematic
770 use of ‘futures’ now allows higher software levels to transparently support experiment pipelining and
771 asynchronous experiment execution in general. Additionally, dividing experiments into a definition and an
772 execution phase also facilitates experiment correctness, software stack flexibility —by decoupling hardware
773 usage from experiment definition— as well as increased platform performance by enabling a separation of
774 hardware access from other aspects of the experiment.

775 The new software framework is expert-friendly: we designed the software layers to facilitate composition
776 between higher- and lower-level application programming interfaces (APIs). Domain experts can therefore
777 define experiments on a higher abstraction level in certain aspects, and are still able to access low-level
778 functionality. A software package for calibration routines —the process of tuning hardware parameters to
779 the requirements defined by an experiment— provides algorithms and settings for typical parameterizations
780 of the chip, including support for multi-compartmental neurons and non-spiking use cases. An experiment
781 micro scheduler service allows to pipeline experiment runs, and even preempt longer experiment sessions
782 of individual users, to decrease hardware platform latency for other user sessions. Enabling multiple
783 high-level modeling interfaces —such as PyNN and PyTorch— to cover a larger user base was one
784 of the new requirements for BSS-2. To achieve this, we provide a separate high-level representation
785 of user-defined experiments. This signal-graph-based representation is generally suited for high-level
786 configuration validation, optimization, and transformation from higher- to lower-level abstractions. The
787 modeling API wrappers merely provide conversions between data types and call semantics. The embedded
788 microprocessors allow for many new applications: Initially designed to increase flexibility for online
789 learning rules (Friedmann et al., 2017), they have been also used for: environment simulations (Schreiber
790 et al., 2022; Pehle et al., 2022), online calibration (section 3.3), general optimization tasks, as well
791 as experiment control (Wunderlich et al., 2019). We ported our low-level chip configuration interface
792 to the embedded processors and thereby allow for code sharing between host and embedded program
793 parts in addition to a software library for embedded use cases. Apart from features directly concerning

794 platform users, we enhanced the support for multiple hardware revisions in parallel facilitating hardware
795 development, commissioning and platform operation. In combination with a dedicated communication
796 layer, this enables not only support for multiple communication backends between host computer and
797 field-programmable gate array (FPGA), such as gigabit ethernet (GbE) or a memory-mapped interface
798 for hybrid FPGA-CPU systems, but also for co-simulation and therefore co-development of software
799 and hardware. Finally, we operate BSS-2 as a research platform. As a result of our contributions to
800 the design and implementation of the EBRAINS (ebr, 2022) software distribution, interactive usage of
801 BSS-2 is now available to a world-wide research community. To summarize, we motivated key design
802 decisions and demonstrated their implementation based on existing use cases: Support for multiple top-level
803 APIs for ‘biological’ and ‘functional’ modeling; support for the embedded microprocessors including
804 structured data exchange with the host, a multi-platform low-level hardware-abstraction layer, and an
805 embedded execution runtime and helper library; support for artificial neural networks in host-based and
806 standalone applications; focus on the user community by providing an integrated platform; sustainable
807 hardware-software co-development.

808 To build a versatile modeling platform, BSS-2 is a neuromorphic system that improved upon successful
809 properties of predecessors, both, in terms of hardware and software. Simulation speed continues to be
810 an important point in computational neuroscience. The development of new approaches to numerical
811 simulation promising lower execution times and better scalability is an active field of research (Knight and
812 Nowotny, 2018, 2021; Abi Akar et al., 2019), as is improving existing simulation codes (Kunkel et al., 2014;
813 Jordan et al., 2018). Whereas parameter sweeps scale trivially, systematically studying model dynamics
814 over sufficiently long periods as well as iterative approaches to training and plasticity can only benefit
815 from increases in simulation speed. The physical modeling approach of the accelerated neuromorphic
816 architectures allows for a higher emulation speed than state-of-the-art numerical simulations (Zenke and
817 Gerstner, 2014; van Albada et al., 2021). BSS-2 can serve as an accelerator for spiking neural networks
818 and therefore opens up opportunities to work on scientific questions that aren’t accessible by numerical
819 simulation. However, to deliver on this promise in reality, both, hardware and software need to be carefully
820 designed, implemented and applied. The publications building on BSS-2 are evidence of what is possible in
821 terms of modeling on accelerated neuromorphic hardware (Bohnstingl et al., 2019; Billaudelle et al., 2020,
822 2021; Cramer et al., 2019, 2022; Czischek et al., 2022; Göltz et al., 2021; Kaiser et al., 2021; Klassert et al.,
823 2021; Klein et al., 2021; Müller et al., 2020a; Schreiber et al., 2022; Spilger et al., 2020; Stradmann et al.,
824 2021; Weis et al., 2020; Wunderlich et al., 2019).

825 We believe that these publications offer a first glimpse of what will be possible in a scaled-up system.
826 The next step on the roadmap is a multi-chip BSS-2 setup employing EXTOLL (Neuwirth et al., 2015;
827 Resch et al., 2014) for host and inter-chip connectivity. First multi-chip experiments have been performed
828 on a lab setup (Thommes et al., 2022). Additionally, a multi-chip system reusing BSS-1 wafer-scale
829 infrastructure is in the commissioning phase and will provide up to 46 BSS-2 chips. Similar to BSS-1, a
830 true wafer-scale version of BSS-2 will provide an increase in terms of resources by one order of magnitude
831 and thus will enable research that not only looks at dynamics at different temporal scales, but also on
832 larger spatial scales. In terms of software we have been adapting our roadmap continuously to match
833 modelers’ expectations. For example, we work on future software abstractions that will allow for flexible
834 descriptions of spiking network models with arbitrary topology in a machine learning framework. PyTorch
835 libraries such as BindsNET (Hazan et al., 2018) or Norse (Pehle and Pedersen, 2021) enable efficient
836 machine-learning-inspired modeling with spiking neural networks and would benefit from neuromorphic
837 hardware support.

ACKNOWLEDGMENTS

838 The authors wish to thank all present and former members of the Electronic Vision(s) research group
839 contributing to the BrainScaleS-2 neuromorphic platform.

FUNDING

840 This work has received funding from the EC Horizon 2020 Framework Programme under grant agreements
841 785907 (HBP SGA2) and 945539 (HBP SGA3), the Deutsche Forschungsgemeinschaft (DFG, German
842 Research Foundation) under Germany's Excellence Strategy EXC 2181/1-390900948 (the Heidelberg
843 STRUCTURES Excellence Cluster), the German Federal Ministry of Education and Research under grant
844 number 16ES1127 as part of the *Pilotinnovationswettbewerb 'Energieeffizientes KI-System'*, the Helmholtz
845 Association Initiative and Networking Fund [Advanced Computing Architectures (ACA)] under Project
846 SO-092, as well as from the Manfred Stärk Foundation, and the Lautenschläger-Forschungspreis 2018 for
847 Karlheinz Meier.

REFERENCES

- 848 (2021). *NICE '20: Proceedings of the Neuro-Inspired Computational Elements Workshop* (New York, NY,
849 USA: Association for Computing Machinery). doi:10.1145/3381755
- 850 [Webpage] (2022). EBRAINS research infrastructure. <https://ebrains.eu>
- 851 Aamir, S. A., Müller, P., Kiene, G., Kriener, L., Stradmann, Y., Grübl, A., et al. (2018). A mixed-signal
852 structured AdEx neuron for accelerated neuromorphic cores. *IEEE Transactions on Biomedical Circuits
853 and Systems* 12, 1027–1037. doi:10.1109/TBCAS.2018.2848203
- 854 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., et al. (2015). *TensorFlow: Large-Scale
855 Machine Learning on Heterogeneous Distributed Systems*. Whitepaper
- 856 Abi Akar, N., Cumming, B., Karakasis, V., Küsters, A., Klijn, W., Peyser, A., et al. (2019). Arbor—a
857 morphologically-detailed neural network simulation library for contemporary high-performance com-
858 puting architectures. In *2019 27th euromicro international conference on parallel, distributed and
859 network-based processing (PDP)* (IEEE), 274–282
- 860 Bellec, G., Scherr, F., Subramoney, A., Hajek, E., Salaj, D., Legenstein, R., et al. (2020). A solution to
861 the learning dilemma for recurrent networks of spiking neurons. *Nature Communications* 11, 3625.
862 doi:10.1038/s41467-020-17236-y
- 863 Benjamin, B. V., Gao, P., McQuinn, E., Choudhary, S., Chandrasekaran, A. R., Bussat, J.-M., et al. (2014).
864 Neurogrid: A mixed-analog-digital multichip system for large-scale neural simulations. *Proceedings of
865 the IEEE* 102, 699–716
- 866 Billaudelle, S., Cramer, B., Petrovici, M. A., Schreiber, K., Kappel, D., Schemmel, J., et al. (2021).
867 Structural plasticity on an accelerated analog neuromorphic hardware system. *Neural networks : the
868 official journal of the International Neural Network Society* 133, 11–20. doi:10.1016/j.neunet.2020.09.
869 024
- 870 Billaudelle, S., Stradmann, Y., Schreiber, K., Cramer, B., Baumbach, A., Dold, D., et al. (2020). Ver-
871 satile emulation of spiking neural networks on an accelerated neuromorphic substrate. In *2020 IEEE
872 International Symposium on Circuits and Systems (ISCAS)* (IEEE). doi:10.1109/iscas45731.2020.
873 9180741
- 874 Bohnstingl, T., Scherr, F., Pehle, C., Meier, K., and Maass, W. (2019). Neuromorphic hardware learns to
875 learn. *Frontiers in Neuroscience* 2019, 1–14. doi:10.3389/fnins.2019.00483

- 876 Brette, R. and Gerstner, W. (2005). Adaptive exponential integrate-and-fire model as an effective description
877 of neuronal activity. *J. Neurophysiol.* 94, 3637 – 3642. doi:10.1152/jn.00686.2005
- 878 Brüderle, D., Müller, E., Davison, A., Müller, E., Schemmel, J., and Meier, K. (2009). Establishing
879 a novel modeling tool: a python-based interface for a neuromorphic hardware system. *Frontiers in*
880 *Neuroinformatics* 3, 17. doi:10.3389/neuro.11.017.2009
- 881 Cramer, B., Billaudelle, S., Kanya, S., Leibfried, A., Grübl, A., Karasenko, V., et al. (2022). Surrogate
882 gradients for analog neuromorphic computing. *Proceedings of the National Academy of Sciences* 119
- 883 Cramer, B., Stöckel, D., Krefft, M., Wibrals, M., Schemmel, J., Meier, K., et al. (2019). Control of criticality
884 and computation in spiking neuromorphic networks with plasticity
- 885 Czischek, S., Baumbach, A., Billaudelle, S., Cramer, B., Kades, L., Pawlowski, J. M., et al. (2022). Spiking
886 neuromorphic chip learns entangled quantum states. *SciPost Phys.* 12, 39. doi:10.21468/SciPostPhys.12.
887 1.039
- 888 Dally, W. J., Turakhia, Y., and Han, S. (2020). Domain-specific hardware accelerators. *Commun. ACM* 63,
889 48–57. doi:10.1145/3361682
- 890 Davies, M., Srinivasa, N., Lin, T.-H., China, G., Cao, Y., Choday, S. H., et al. (2018). Loihi: A
891 neuromorphic manycore processor with on-chip learning. *IEEE Micro* 38, 82–99. doi:10.1109/MM.
892 2018.112130359
- 893 Davison, A. P., Brüderle, D., Eppler, J., Kremkow, J., Müller, E., Pecevski, D., et al. (2009). PyNN: a
894 common interface for neuronal network simulators. *Front. Neuroinform.* 2. doi:10.3389/neuro.11.011.
895 2008
- 896 [Software] Delta V Software (2020). Remote call framework
- 897 Facebook, Inc. (2021a). *PyTorch JIT Overview*
- 898 Facebook, Inc. (2021b). *PyTorch on XLA Devices*
- 899 Friedmann, S., Schemmel, J., Grübl, A., Hartel, A., Hock, M., and Meier, K. (2017). Demonstrating hybrid
900 learning in a flexible neuromorphic hardware system. *IEEE Transactions on Biomedical Circuits and*
901 *Systems* 11, 128–142. doi:10.1109/TBCAS.2016.2579164
- 902 Furber, S. B., Lester, D. R., Plana, L. A., Garside, J. D., Painkras, E., Temple, S., et al. (2012).
903 Overview of the SpiNNaker system architecture. *IEEE Transactions on Computers* 99. doi:http:
904 //doi.ieeecomputersociety.org/10.1109/TC.2012.142
- 905 Gamblin, T., LeGendre, M., Collette, M. R., Lee, G. L., Moody, A., de Supinski, B. R., et al. (2015). The
906 spack package manager: Bringing order to hpc software chaos. In *Proceedings of the International*
907 *Conference for High Performance Computing, Networking, Storage and Analysis* (New York, NY, USA:
908 ACM), SC '15, 40:1–40:12. doi:10.1145/2807591.2807623
- 909 Gewaltig, M.-O. and Diesmann, M. (2007). NEST (NEural Simulation Tool). *Scholarpedia* 2, 1430.
910 doi:10.4249/scholarpedia.1430
- 911 Gidon, A., Zolnik, T. A., Fidzinski, P., Bolduan, F., Papoutsis, A., Poirazi, P., et al. (2020). Dendritic action
912 potentials and computation in human layer 2/3 cortical neurons. *Science* 367, 83–87. doi:10.1126/
913 science.aax6239
- 914 [Software] GNU Project (2018). The GNU Compiler Collection 8.1. Website. Free Software Foundation
915 Inc.
- 916 Goddard, N. H., Hucka, M., Howell, F., Cornelis, H., Shankar, K., and Beeman, D. (2001). Towards
917 neuroml: model description methods for collaborative modelling in neuroscience. *Philos Trans R Soc*
918 *Lond B Biol Sci* 356, 1209–28

- 919 Göltz, J., Kriener, L., Baumbach, A., Billaudelle, S., Breitwieser, O., Cramer, B., et al. (2021). Fast
920 and energy-efficient neuromorphic deep learning with first-spike times. *Nature Machine Intelligence* 3,
921 823–835. doi:10.1038/s42256-021-00388-x
- 922 [Software] Grant, W. S. and Voorhies, R. (2017). cereal - a C++11 library for serialization
- 923 Hazan, H., Saunders, D. J., Khan, H., Patel, D., Sanghavi, D. T., Siegelmann, H. T., et al. (2018). BindsNET:
924 A machine learning-oriented spiking neural networks library in python. *Frontiers in Neuroinformatics*
925 12, 89. doi:10.3389/fninf.2018.00089
- 926 Hines, M. and Carnevale, N. (2003). *The NEURON simulation environment*. (M.A. Arbib). 769–773
- 927 Indiveri, G., Linares-Barranco, B., Hamilton, T. J., van Schaik, A., Etienne-Cummings, R., Delbruck, T.,
928 et al. (2011). Neuromorphic silicon neuron circuits. *Frontiers in Neuroscience* 5. doi:10.3389/fnins.
929 2011.00073
- 930 Jordan, J., Ippen, T., Helias, M., Kitayama, I., Sato, M., Igarashi, J., et al. (2018). Extremely scal-
931 able spiking neuronal network simulation code: From laptops to exascale computers. *Frontiers in*
932 *Neuroinformatics* 12, 2. doi:10.3389/fninf.2018.00002
- 933 Kaiser, J., Billaudelle, S., Müller, E., Tetzlaff, C., Schemmel, J., and Schmitt, S. (2021). Emulating dendritic
934 computing paradigms on analog neuromorphic hardware. *Neuroscience* doi:10.1016/j.neuroscience.
935 2021.08.013
- 936 Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint*
- 937 Klassert, R., Baumbach, A., Petrovici, M. A., and Gärttner, M. (2021). Variational learning of quantum
938 ground states on spiking neuromorphic hardware
- 939 Klein, B., Kuhn, L., Weis, J., Emmel, A., Stradmann, Y., Schemmel, J., et al. (2021). Towards addressing
940 noise and static variations of analog computations using efficient retraining. In *Machine Learning*
941 *and Principles and Practice of Knowledge Discovery in Databases* (Cham: Springer International
942 Publishing), 409–420. doi:10.1007/978-3-030-93736-2_32
- 943 Knight, J. C. and Nowotny, T. (2018). GPUs outperform current HPC and neuromorphic solutions in terms
944 of speed and energy when simulating a highly-connected cortical model. *Frontiers in neuroscience* 12,
945 941
- 946 Knight, J. C. and Nowotny, T. (2021). Larger GPU-accelerated brain simulations with procedural
947 connectivity. *Nature Computational Science* 1, 136–142
- 948 Kriener, L., Göltz, J., and Petrovici, M. A. (2021). The yin-yang dataset. *arXiv preprint*
- 949 Kungl, A. F., Schmitt, S., Klähn, J., Müller, P., Baumbach, A., Dold, D., et al. (2019). Accelerated
950 physical emulation of bayesian inference in spiking neural networks. *Frontiers in Neuroscience* 13, 1201.
951 doi:10.3389/fnins.2019.01201
- 952 Kunkel, S., Schmidt, M., Eppler, J. M., Plesser, H. E., Masumoto, G., Igarashi, J., et al. (2014). Spiking
953 network simulation code for petascale computers. *Frontiers in neuroinformatics* 8, 78
- 954 Major, G., Larkum, M. E., and Schiller, J. (2013). Active properties of neocortical pyramidal neuron
955 dendrites. *Annual Review of Neuroscience* 36, 1–24. doi:10.1146/annurev-neuro-062111-150343
- 956 Mason, S. J. (1953). Feedback theory-some properties of signal flow graphs. *Proceedings of the IRE* 41,
957 1144–1156. doi:10.1109/JRPROC.1953.274449
- 958 Müller, E., Mauch, C., Spilger, P., Breitwieser, O. J., Klähn, J., Stöckel, D., et al. (2020a). Extending
959 BrainScaleS OS for BrainScaleS-2. *arXiv preprint*
- 960 Müller, E., Schmitt, S., Mauch, C., Schmidt, H., Montes, J., Ilmberger, J., et al. (2020b). The operating
961 system of the neuromorphic BrainScaleS-1 system. *arXiv preprint* Submitted to Neurocomputing OSP

- 962 Neuwirth, S., Frey, D., Nuessle, M., and Bruening, U. (2015). Scalable communication architecture
963 for network-attached accelerators. In *2015 IEEE 21st International Symposium on High Performance*
964 *Computer Architecture (HPCA)* (IEEE), 627–638. doi:10.1109/HPCA.2015.7056068
- 965 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., et al. (2019). Pytorch: An imperative
966 style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*
967 32, eds. H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (Curran
968 Associates, Inc.). 8024–8035
- 969 Pehle, C., Billaudelle, S., Cramer, B., Kaiser, J., Schreiber, K., Stradmann, Y., et al. (2022). The
970 BrainScaleS-2 accelerated neuromorphic system with hybrid plasticity. *Frontiers in Neuroscience* 16.
971 doi:10.3389/fnins.2022.795876
- 972 [Software] Pehle, C. and Pedersen, J. E. (2021). Norse — A deep learning library for spiking neural
973 networks. doi:10.5281/zenodo.4422025. Documentation: <https://norse.ai/docs/>
- 974 Poirazi, P. and Papoutsi, A. (2020). Illuminating dendritic function with computational models. *Nature*
975 *reviews. Neuroscience* 21, 303–321. doi:10.1038/s41583-020-0301-7
- 976 PowerISA (2010). *PowerISA Version 2.06 Revision B*. Specification, Power.org
- 977 Resch, M. M., Bez, W., Focht, E., Kobayashi, H., and Patel, N. (eds.) (2014). *Sustained Simulation*
978 *Performance 2014: Proceedings of the Joint Workshop on Sustained Simulation Performance*. University
979 of Stuttgart (HLRS) and Tohoku University (Springer)
- 980 Rhodes, O., Bogdan, P. A., Brenninkmeijer, C., Davidson, S., Fellows, D., Gait, A., et al. (2018). spynnaker:
981 A software package for running pynn simulations on spinnaker. *Frontiers in Neuroscience* 12, 816.
982 doi:10.3389/fnins.2018.00816
- 983 Schemmel, J., Billaudelle, S., Dauer, P., and Weis, J. (2020). Accelerated analog neuromorphic computing.
984 *arXiv preprint*
- 985 Schmitt, S., Klähn, J., Bellec, G., Grübl, A., Güttler, M., Hartel, A., et al. (2017). Neuromorphic hardware
986 in the loop: Training a deep spiking network on the brainscales wafer-scale system. *Proceedings of the*
987 *2017 IEEE International Joint Conference on Neural Networks* doi:10.1109/IJCNN.2017.7966125
- 988 Schreiber, K., Wunderlich, T., Spilger, P., Billaudelle, S., Cramer, B., Stradmann, Y., et al. (2022). Insectoid
989 path integration on accelerated neuromorphic hardware. In preparation
- 990 Shrestha, S. B. and Orchard, G. (2018). SLAYER: Spike layer error reassignment in time. In *Advances*
991 *in Neural Information Processing Systems*, eds. S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
992 N. Cesa-Bianchi, and R. Garnett (Curran Associates, Inc.), vol. 31
- 993 Spilger, P., Müller, E., Emmel, A., Leibfried, A., Mauch, C., Pehle, C., et al. (2020). hxtorch: PyTorch
994 for BrainScaleS-2 — perceptrons on analog neuromorphic hardware. In *IoT Streams for Data-Driven*
995 *Predictive Maintenance and IoT, Edge, and Mobile for Embedded Machine Learning* (Cham: Springer
996 International Publishing), 189–200. doi:10.1007/978-3-030-66770-2_14
- 997 Stimberg, M., Brette, R., and Goodman, D. F. (2019). Brian 2, an intuitive and efficient neural simulator.
998 *eLife* 8. doi:10.7554/eLife.47314
- 999 Stradmann, Y., Billaudelle, S., Breitwieser, O., Ebert, F. L., Emmel, A., Husmann, D., et al. (2021).
1000 Demonstrating analog inference on the BrainScaleS-2 mobile system. *arXiv preprint*
- 1001 Suhan, A., Libenzi, D., Zhang, A., Schuh, P., Saeta, B., Sohn, J. Y., et al. (2021). Lazytensor: combining
1002 eager execution with domain-specific compilers
- 1003 Thommes, T., Bordukat, S., Grübl, A., Karasenko, V., Müller, E., and Schemmel, J. (2022). Demonstrating
1004 BrainScaleS-2 inter-chip pulse communication using EXTOLL. In *Neuro-inspired Computational*
1005 *Elements Workshop (NICE ’22), March 29 – April 1, 2022*. Accepted

- 1006 van Albada, S. J., Pronold, J., van Meegen, A., and Diesmann, M. (2021). Usage and scaling of an
1007 open-source spiking multi-area model of monkey cortex. In *Brain-Inspired Computing*, eds. K. Amunts,
1008 L. Grandinetti, T. Lippert, and N. Petkov (Cham: Springer International Publishing), 47–59
1009 [Software] Vinkelis, M. (2020). Bitsery. <https://github.com/frailt/bitsery>
1010 Weis, J., Spilger, P., Billaudelle, S., Stradmann, Y., Emmel, A., Müller, E., et al. (2020). Inference with
1011 artificial neural networks on analog neuromorphic hardware. In *IoT Streams for Data-Driven Predictive*
1012 *Maintenance and IoT, Edge, and Mobile for Embedded Machine Learning* (Cham: Springer International
1013 Publishing), 201–212. doi:10.1007/978-3-030-66770-2_15
1014 Wunderlich, T., Kungl, A. F., Müller, E., Hartel, A., Stradmann, Y., Aamir, S. A., et al. (2019). Demon-
1015 strating advantages of neuromorphic computation: A pilot study. *Frontiers in Neuroscience* 13, 260.
1016 doi:10.3389/fnins.2019.00260
1017 Yavuz, E., Turner, J., and Nowotny, T. (2016). GeNN: a code generation framework for accelerated brain
1018 simulations. *Scientific reports* 6, 1–14
1019 Yoo, A. B., Jette, M. A., and Grondona, M. (2003). Slurm: Simple linux utility for resource management.
1020 In *Workshop on Job Scheduling Strategies for Parallel Processing* (Springer), 44–60
1021 Zenke, F. and Ganguli, S. (2018). SuperSpike: Supervised learning in multilayer spiking neural networks.
1022 *Neural Computation* 30, 1514–1541. doi:10.1162/neco_a_01086
1023 Zenke, F. and Gerstner, W. (2014). Limits to high-speed simulations of spiking neural networks using
1024 general-purpose computers. *Frontiers in Neuroinformatics* 8. doi:10.3389/fninf.2014.00076

