

# Training and Deploying Spiking NN Applications to the Mixed-Signal Neuromorphic Chip Dynap<sup>TM</sup>-SE2 with Rockpool

Ugurcan Cakal MSc<sup>1</sup>, Prof Ilkay Ulusoy PhD<sup>2</sup> and Dr Dylan Muir PhD<sup>1</sup>

<sup>1</sup>SynSense AG, Thurgauerstrasse 60, Zürich, 8050, Switzerland.

<sup>2</sup>Electrical and Electronics Engineering, METU, Ankara, 06800, Turkey.

Contributing authors: [ugurcan.cakal@synsense.ai](mailto:ugurcan.cakal@synsense.ai); [ilkay@metu.edu.tr](mailto:ilkay@metu.edu.tr); [dylan.muir@synsense.ai](mailto:dylan.muir@synsense.ai);

## Abstract

Mixed-signal neuromorphic processors provide extremely low-power operation for edge inference workloads, taking advantage of sparse asynchronous computation within Spiking Neural Networks (SNNs). However, deploying robust applications to these devices is complicated by limited controllability over analog hardware parameters, unintended parameter and dynamics variations of analog circuits due to fabrication non-idealities. Here we demonstrate a novel methodology for offline training and deployment of spiking neural networks (SNNs) to the mixed-signal neuromorphic processor Dynap-SE2. The methodology utilizes an unsupervised weight quantization method to optimize the network’s parameters, coupled with adversarial parameter noise injection during training. The optimized network is shown to be robust to the effects of quantization and device mismatch, making the method a promising candidate for real-world applications with hardware constraints. This work extends Rockpool, an open-source deep-learning library for SNNs, with support accurate simulation of mixed-signal SNN dynamics. Our approach simplifies the development and deployment process for the neuromorphic community, making mixed-signal neuromorphic processors more accessible to researchers and developers.

**Keywords:** mixed-signal, neuromorphic, spiking neural networks, Dynap-SE2

## Introduction

Neuromorphic processors use analog and mixed-signal circuits to emulate the dynamics and computational abilities of biological neurons and synapses. One of the most advanced architectures is Dynap-SE2, which has an asynchronous mixed-signal structure whose analog components operate in the subthreshold operation range, making it an ultra-low power and ultra-low latency application delivery candidate.

Devices such as Dynap-SE2 offer a high degree of realism and configurability, but have been historically difficult to configure, for several reasons. Dynap-SE2 and other similar devices individually instantiate arrays

of synapses and neurons in analog circuits. Sub-threshold operation of the these emulation circuits exposes them to variability of the individual device components, due to variations introduced in the fabrication process. This variability, known as “mismatch”, causes the behaviour of ostensibly identical neurons and synapses to differ across a chip, and between chips [1, 2]. Mixed-signal devices such as Dynap-SE2 do not usually permit individual control over each parameter on the chip, instead grouping parameters such as time-constants, thresholds, and even weight values across a number of neurons and synapses. This grouping reduces the parameter space of the device,

but makes it difficult to calibrate in the face of mismatch. In addition, grouping parameters means that the parameter configuration space of an SNN deployed to the chip is itself heavily constrained.

As a result, training and deploying applications to these devices usually requires many months of manual effort by experienced researchers.

To bridge this gap, this work introduces a fast and efficient software toolchain that provides the potential for commercial application development for Dynap-SE2. Extending Rockpool, an open-source deep-learning library for SNNs [3], the toolchain provides the necessary mechanisms to perform offline optimization of SNNs that can be robustly deployed to a number of Dynap-SE2 chips, while preserving behavior in the face of mismatch. This work provides a Dynap-SE2 simulator, “DynapSim”, which operates in the same parameter space as Dynap-SE2 family processors [4]. DynapSim executes an efficient and accurate simulation of the Dynap-SE2 design dynamics to solve the characteristic circuit transfer functions over time.

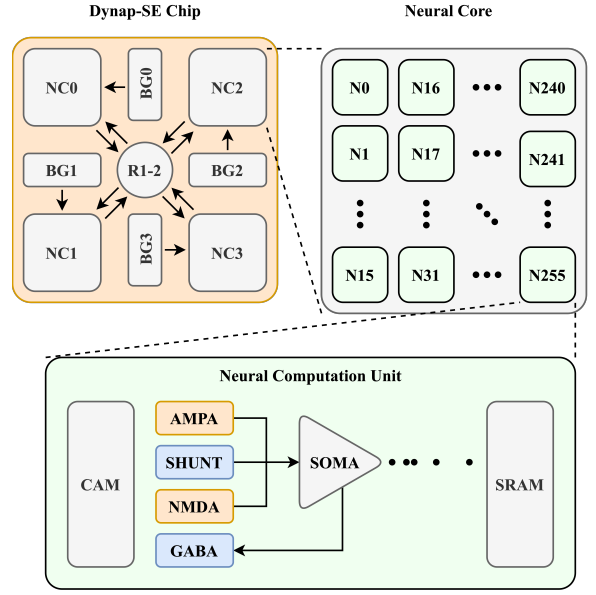
DynapSim is used as a computational neuron model in offline SNN simulations and during training. Rockpool translates optimized networks to equivalent hardware configurations, and manages deployment of these networks to Dynap-SE2 chips.

This manuscript provides an overview of the Dynap-SE2 hardware. We then describe the implementation details of DynapSim, and demonstrate training, deployment and quantitative evaluation of a toy model to Dynap-SE2 hardware. We present the steps to achieve training and deployment using code examples for Rockpool.

In the following weeks, we will be increasing the simulation coverage, enhancing the user experience, proposing more use cases, and updating this manuscript. The most up-to-date implementation can be found here: <https://github.com/synsense/rockpool>.

## Overview of the Hardware

The DYNAMIC Neuromorphic Asynchronous Processor — Scalable 2 (Dynap-SE2) is a mixed-signal chip that inherits the event-driven nature of the DYNAP family [5]. It directly emulates biological behavior using analog spiking neurons and analog synapses as the computing units. The neural cores’ transistors operate in the subthreshold region, resulting in power consumption below 1 mW. Each Dynap-SE2 chip is equipped with 1024 adaptive exponential integrate-and-fire (AdExpIF) analog ultra-low-power spiking neurons



**Fig. 1 Dynap-SE2 Architecture.** NC: Neural Core; R: Router; BG: Bias Generator. Other acronyms: see text.

and 64 synapses per neuron. Fig. 1 displays an abstract overview of the architecture of the chip.

The neural computation unit serves as the primary building block for creating the dynamics in Dynap-SE2. Each neural core (NC) consists of 256 analog neurons that share the same parameter set. The digital memory blocks, Content-Addressable Memories (CAMs) and Static-RAMs (SRAMs), store the transmitting and receiving event configurations, respectively. The synapses and neuron soma carry out analog computations, with four different types of synapses — AMPA, GABA, NMDA, and SHUNT — integrating the incoming events and injecting current into the membrane. AMPA and NMDA activation increase the firing probability, while GABA and SHUNT activation decrease it.

The CAM stores the listening event setting for each of the 64 connections of a neuron, which specifies its synaptic processing unit. The neuron soma integrates the injection currents and holds a temporal state, with configurable paths of charging and discharging capacitors designating the temporal behavior. The membrane current, which is a secondary reading on the membrane capacitance, functions as the temporal state variable. When the membrane current reaches the firing threshold, the neuron’s reset mechanism triggers the event sensing units, which package the event in Address Event Representation (AER) format and broadcast it

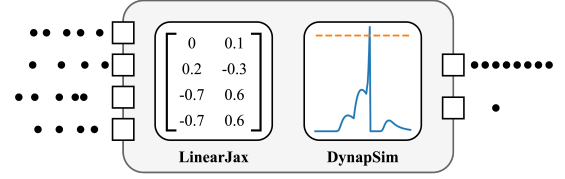
to the indicated locations. In this way, the neuron uses analog sub-threshold circuits to compute the dynamics but conveys the resulting outputs using a digital routing mechanism.

Each neural core holds a parameter group to set the neuronal and synaptic parameters for its 256 neurons and their pre-synaptic synapses. The neurons in the same core share the same parameter values, including time constants, refractory periods, synaptic connection strengths, and other attributes. Special digital-to-analog converters, bias generators (BG), set these parameter current values. In total, there are 70 parameters that can be set to adjust the behavior of the neurons and synapses, including time constants, pulse widths, amplifier gain ratios, and synaptic weight strengths, among others.

For simulation purposes, a custom computational spiking neural model relates the behavioral dynamics of a computational neural setting to the VLSI parameters of the respective circuits. It uses forward Euler updates to predict the time-dependent dynamics and solves the characteristic circuit transfer functions in time. Specifically, a “DynapSim” neuron solves the silicon neuron [6] and silicon synapse [7] circuit equations, making use of assumptions and simplifications from [8]. Further details of the application and implementation can be found in [4].

## Training Pipeline

DynapSim is an extension of the contemporary spiking neural network library, Rockpool [3], and serves as a simulation solution for the device. The solution it offers involves solving characteristic equations of the analog circuits and does not provide a circuit-level accurate simulation. Instead, DynapSim provides an approximate simulation that can be fine-tuned and translated into a device configuration. The simulator is powered by the state-of-the-art high-performance machine learning library JAX [9], which facilitates fast execution. The toolchain we provide performs off-chip gradient-based optimization of an SNN and deploys the trained network to the chip while preserving the optimized behavior. The upcoming sections elaborate on the technique of gradient-based optimization employed to train a Spiking Neural Network (SNN) before its deployment to a Dynap-SE2 chip.



**Fig. 2 Frozen noise classification task.** 60 input channels provide input spiking patterns to the network (left; 4 shown here). The network provides two output channels (right), which should emit high spiking activity when presented with one of two target frozen noise inputs. “LinearJax” and “DynapSim” modules, provided by Rockpool, are used to simulate the weights, synapse and neuron dynamics of the network (see text for further details).

## Toy task: Frozen Noise Classification

The purpose of the frozen noise classification experiment is to evaluate the learning abilities of the implemented simulator. The experiment focuses on training a DynapSim network to accurately classify two distinct random frozen noise patterns. The network comprises two analog neurons with recurrent connections, along with 60 external input connections. The desired outcome is for the first neuron to exhibit a significantly higher firing rate when presented with the first frozen noise, and for the second neuron to exhibit a significantly higher firing rate when presented with the second frozen noise. Fig. 2 illustrates the task at hand.

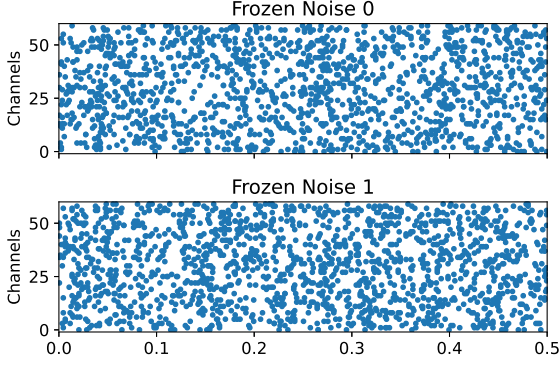
## Data

To run the experiment with the spiking neuron model, a spiking input pattern is necessary. For this specific task, randomly generated discrete Poisson time series with a mean frequency of 50 Hz in a 500 ms duration are used as frozen noise recordings. Each sample comprises 60 channels, and the time-step duration is 1 ms.

For training purposes, two samples are utilized to enable the network to overfit, while 1000 different random samples are reserved for testing the trained network. The optimized network should recognize the 2 critical training samples, by generating high activity on the corresponding output neuron, and low activity on the other neuron. If any other sample is provided, the network should generate random output activity.

## Network

The network architecture used in this study consists of two modules. The first module, **LinearJax**, applies a linear transformation to the input spikes to simulate spike weighting. The second module, **DynapSim**, simulates the time-dependent analog silicon neuron and synapse dynamics. Each neuron in this layer generates a spike



**Fig. 3** Frozen noise recordings used in training.

train as output. Listing 1 shows how to instantiate this network in Rockpool.

---

```
net = Sequential(
    LinearJax((Nin, Nrec)),
    DynapSim((Nrec, Nrec), dt=dt),
)
```

---

**Listing 1** Constructing the SNN in Rockpool.

---

This network architecture can be compared to using a ReLU activation layer following a fully connected layer in classical NNs. The difference, however, lies in the fact that the DynapSim layer computes and maintains a time-dependent state instead of a stateless activation. The output of the DynapSim neurons depends not only on the instantaneous inputs but also on past inputs via internal state variables. The state continues to evolve continuously over time, regardless of when the neuron receives spikes on its input. Additionally, the DynapSim layer encapsulates a recurrent connection matrix that is one of the targets of the optimizer. Lastly, the layer corresponds to a custom analog hardware configuration, and solving the characteristic equations of the analog circuits is a key part of its function.

To limit the complexity of the task, in this case only the weight parameters are trained, while the rest of the neuron and synapse parameters are fixed to their default simulation values. This means that mathematically, only two 2D weight matrices are subject to optimization: the  $60 \times 2$  input weight matrix stored inside LinearJax and the  $2 \times 2$  recurrent weight matrix stored inside DynapSim.

## Response Analysis

Each task requires an appropriate readout strategy. In this one, we first calculate the mean firing rates of the neurons over the duration of the simulation as in Eq. 1.

$$r = \frac{1}{dt \cdot N} \cdot \sum_{i=0}^N S[i] \quad (1)$$

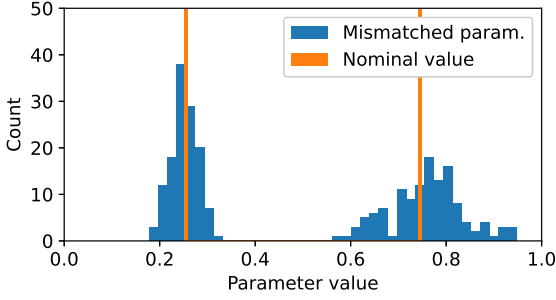
As a performance metric, the ratio between the output neurons' mean firing rate quantify the network's ability to distinguish the two target frozen noise patterns. The firing rate ratio (FRR) is calculated by dividing the higher mean firing rate by the lower mean firing rate read from the decision neurons, as shown in Eq. 2.

$$FRR = \frac{\max(r_0, r_1)}{\min(r_0, r_1)} \quad (2)$$

## Mismatch Simulation

Each optimisation step during training includes a forward and a backward pass. The forward pass simulates the neural dynamics over time and produces the neural activity. The backward pass backpropagates the loss over time, calculating the gradients of parameter values with respect to the loss. The forward computation simulates the circuit behavior in ideal conditions, but the circuits in real Dynap-SE2 devices are subject to parameter mismatch. In order to make trained networks robust against parameter deviations, Büchel et al. proposed to modify values during training by injecting parameter noise as well as by an adversarial attack on parameter values [1].

DynapSim includes an empirically-verified model of parameter mismatch, which we apply in the forward pass, slightly perturbing the parameter values in the network that are subject to change. Mismatch simulation deviates the parameters using a Gaussian distribution, with the mean taken as the nominal parameter values, and variation determined by empirical measurement [1, 2]. New mismatched parameters are set every  $n$  epochs. During optimisation, the network reaches parameter values that obtain a low loss value, in spite of the parameter variation. As a result, SNNs trained in this way are less sensitive to mismatch-induced loss of performance when deployed to mixed-signal neuromorphic hardware. Fig. 4 illustrates the effect of mismatch on parameter values.



**Fig. 4 Mismatch Simulation.** Nominal values for two parameters (orange) are applied to the network, following which DynapSim is used to simulate the parameter mismatch that would be experienced when deployed to a Dynap-SE2 device (blue distributions).

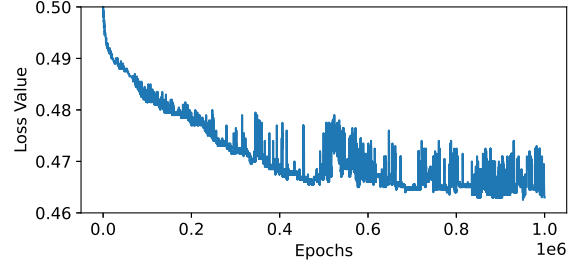
## Optimization

The optimization objective is highly dependent on the task and can be customized according to the requirements. In this particular task, the goal is to increase the firing rate of a specific neuron upon receiving a known frozen noise record. To achieve this, the mean square error (MSE) loss function is utilized. The mathematical formulation of the MSE loss function is given by Eq. 3.

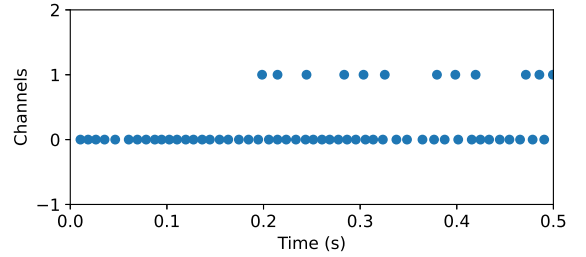
$$f_{\text{MSE}}(y_{\text{out}}, y_{\text{target}}) = \frac{1}{N} \sum_{i=0}^N (y_{\text{out}}[i] - y_{\text{target}}[i])^2 \quad (3)$$

The target train is a uniform spike train that generates an event at every time step from one channel and that generates no events from the other channel. The mean value of the differences in time gives a scalar loss value to be used in error backpropagation. However, the actual neuron model is incapable of producing the ideal spike train aimed due to refractory periods and spike frequency adaptation mechanisms, making it impossible to achieve zero error. The optimizer pushes the neurons to converge towards the ideal spiking regime. The expected behavior during training is that the error will initially be high and gradually decrease to a level that is above zero.

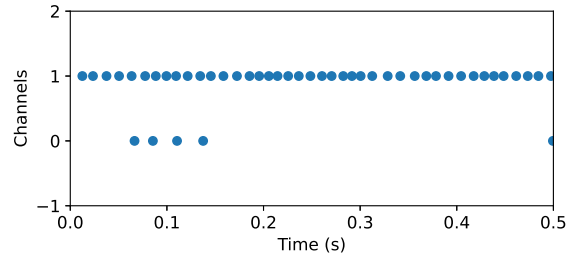
In this experiment, the Adaptive Moment Estimation or Adam algorithm is utilized for optimization, which provides a first-order gradient-based optimization of stochastic objective functions based on adaptive estimates of lower-order moments [10]. So, the training pipeline is similar to a conventional machine learning task, with a few differences. Since the forward computation involves non-differentiable functions, a surrogate function approximation is used in the backward pass [11–14].



**Fig. 5 Mean square error (MSE) loss over the course of training.**



**Fig. 6 Response of the trained network to input class 0.**



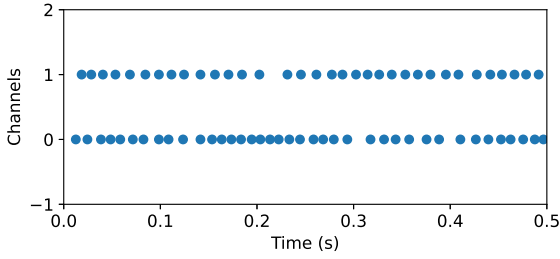
**Fig. 7 Response of the trained network to input class 1.**

Fig. 5 illustrates the decrease in MSE loss over the training process.

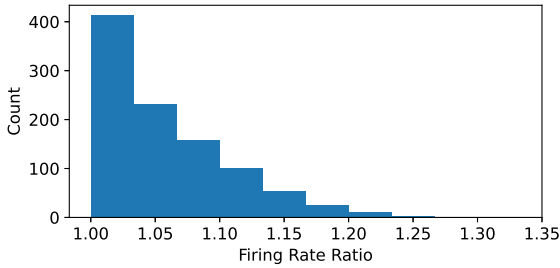
During training, the MSE loss decreased from 0.5 to 0.46 over one million epochs. Even this seemingly small drop results in a significant difference in behavior, allowing the network to classify two similar frozen noise samples. Figs. 6 and 7 show the response of the network to input class 0 and input class 1 respectively.

When the first noise pattern is presented to the network, neuron 0 (channel 0) exhibits high firing activity (164 Hz), while neuron 1 (channel 1) shows significantly lower activity (24 Hz), resulting in an FRR of 6.83. On the other hand, when the second noise pattern is presented, neuron 1 (channel 1) fires almost constantly (122 Hz), while neuron 0 (channel 0) remains quiet as intended (10 Hz), resulting in an FRR of 12.2.





**Fig. 8 Response of the trained network to a random test noise sample.**



**Fig. 9 Distribution of FRR values under 1000 random noise test samples.** Most FRRs are close to one, indicating that the network has learned to reject the non-trained noise inputs. FRRs for the trained classes were above 5.

The clear distinction between the higher and lower firing rates demonstrates that the network is capable of distinguishing between the two input patterns.

To evaluate the network’s recognition capabilities on unseen data, we used a test set of random noise samples to demonstrate that the network only recognizes the training patterns. We generated 1000 frozen noise input patterns with the same mean frequency and length as the target patterns used in training. The FRRs between the decision neurons were recorded to quantify the ability of the network to reject un-trained input patterns. We expect the network to respond with FRRs close to one, indicating the input patterns are not similar to either class 0 or class 1. Fig. 8 shows the network’s response to a random sample.

The FRR in Fig. 8 is close to 1. The distribution of FRR values under 1000 random noise samples is shown in Fig. 9.

During the 500 ms test runs, both the first and second neurons remain active, firing at similar rates. Based on these observations, it can be concluded that the network responds strongly only to the trained target inputs, and rejects the non-trained noise inputs.

---

```
# Define
net = Sequential(
    LinearJax((Nin, Nrec)),
    DynapSim((Nrec, Nrec), dt=dt),
)

# Map
spec = mapper(net.as_graph())
spec.update(autoencoder_quantization(**spec))
config = config_from_specification(**spec)

# Connect & Interface
se2_devices = find_dynapse_boards()
se2 = DynapseSamna(se2_devices[0], **config)
out, state, rec = se2(raster, record=True)
```

**Listing 2 Deploying an SNN to Dynap-SE2.**

---



---

```
net = dynapsim_net_from_config(**config)
out, state, rec = net(raster, record=True)
```

**Listing 3 Extracting an SNN from a hardware configuration.**

---

## Deployment to Dynap<sup>TM</sup>-SE2

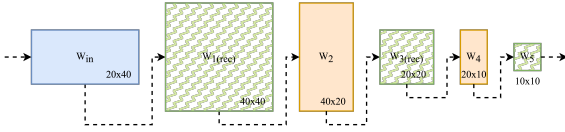
To successfully deploy an SNN using Rockpool, one needs to first build the network in simulation and optimize it using gradient-based or non-gradient-based methods. Rockpool then extracts a computational graph from the optimized network, containing all the necessary parameters for specifying the chip configuration. However, the computational graph does not include information about hardware resource allocation, so a mapping procedure is required to cluster the parameters and find a suitable hardware allocation. The parameters also need to be quantized since the Dynap-SE2 hardware cannot support floating point precision. Finally, the user needs to connect and interface with the chip to deploy the SNN.

In short, successful deployment requires:

1. Building and training an SNN
2. Extracting the computational graph of the SNN
3. Mapping the computational graph to the hardware
4. Quantizing the parameters
5. Connecting and interfacing with the chip

Despite the apparent complexity of these steps, Rockpool simplifies the process down to just a few lines of code, as demonstrated in Listing 2.

Our pipeline also supports “reverse mapping”, whereby a simulation SNN can be extracted from an existing hardware configuration. It’s as easy as in Listing 3.



**Fig. 10** Weight matrix parameters and information flow through an example SNN. Input weights  $W_{in}$  are shown in blue. Blocks of recurrent weights  $W_{n(rec)}$  are hatched. Feed-forward hidden weights  $W_n$  are indicated in orange.

The sections below explain the details of these steps.

## Computational Graph

A computational graph in Rockpool represents the flow of data through a neural network. It includes the computational neuron parameters, hardware parameters, and the order of operation execution.

In Rockpool, the `as_graph()` method extracts a computational graph from a full SNN model. This graph captures all the computationally significant parameters of the network, as well as the network structure. The graph representation enables manipulation of the SNN architecture, facilitating mapping the network parameters to various hardware architectures, and permitting conversion between neuron models.

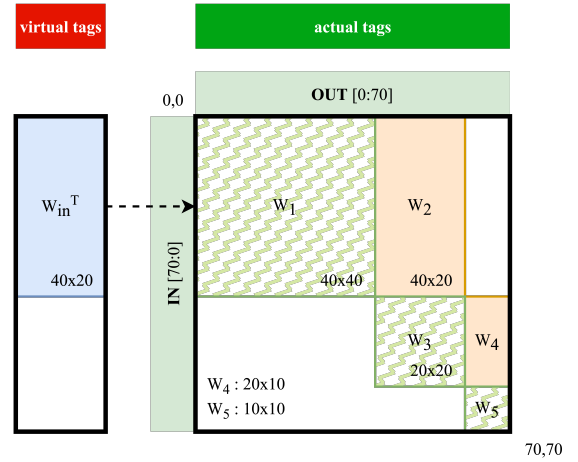
For instance, a trained LIF network can be transformed into an equivalent DynapSim network, with similar behaviour.

## Mapping

The `mapper()` function in the mapping package projecting a computational graph from an arbitrary SNN onto a Dynap-SE2 HDK hardware specification, regardless of whether the network was originally a DynapSim network. `mapper()` clusters the parameters into groups and determines the hardware IDs of neurons. It achieves this by processing a multi-layer network definition to obtain an equivalent single-layer definition, building virtual grids on the weight matrix, and allocating logical regions for intermediate layers in a large recurrent weight matrix onto which the merged network is projected.

Fig. 10 shows the weight parameters and information flow in an example SNN with both feed-forward and recurrent components.

The input weight matrix,  $W_{in}$ , applies a linear transformation to the external input and delivers it to the hardware neurons. Recurrent weight matrices,  $W_{1(rec)}$ ,  $W_{3(rec)}$ , and  $W_{5(rec)}$ , establish the connection weights between hardware neurons, while feed-forward weight



**Fig. 11** Merged weight matrices given the network in Fig. 10. The blocks of weights in Fig. 10 are indicated here with their corresponding colours and labels.

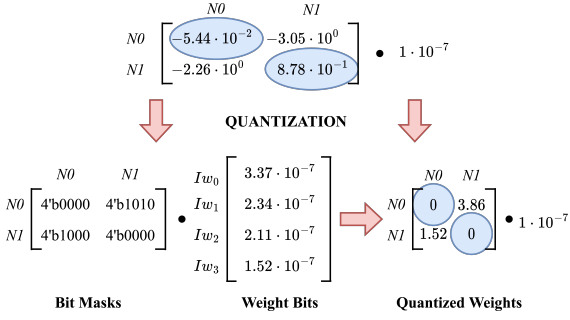
matrices,  $W_2$  and  $W_4$ , connect different groups of neurons to each other. The mapper produces a single equivalent recurrent weight matrix that collects all feed-forward neurons in a single pseudo-recurrent representation.

Fig. 11 shows the merged input and recurrent weight matrices corresponding to the network in Fig. 10.

The input weights and recurrent weights are stored separately, because input neurons on Dynap-SE2 are virtual, and the input side of  $W_{in}$  does not correspond to hardware neurons. Instead,  $W_{in}$  determines how to transmit the external activity to the analog neurons on Dynap-SE2.

All other weight matrices in Fig. 10 are merged into one large recurrent weight matrix. The input neurons are assigned tags (virtual IDs) from the set of virtual input tags (“virtual tags”), while the hardware neurons are assigned tags (hardware IDs) from the list of available hardware neurons (“actual tags”).

Once this is complete, the `mapper()` reduces an SNN down to three connected graph modules: one `DynapseNeurons` object holding the current parameter values of the hardware neurons, one `LinearWeights` object holding the input weights from the external connections to the hardware neurons, and one `LinearWeights` object holding the recurrent weights between the hardware neurons. It is worth noting that while mapping implements the functionality to process the SNN layers and obtain a hardware specification, it cannot guarantee that the network is deployable. Issues such as design rule violations may still arise during the



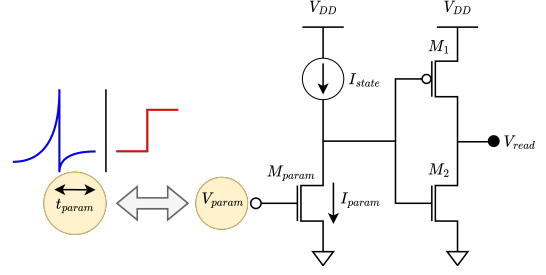
**Fig. 12 Process of weight quantization.** The floating-point weight matrix (top) is converted into a quantized representation, consisting of a bit-mask containing the number of quantal connections made between two neurons as well as connection type (bottom; left) and a set of “base weight” parameters (bottom; middle) indicating the strength of each quantal connection. The reconstructed weight matrix (bottom; right) will be similar to, but not identical to the original matrix. For example, weak connections may be pruned (blue highlights).

later stages of quantization and configuration object generation that the `mapper()` is unable to detect.

## Quantization

During simulation, weight matrices in the layers can have any floating-point value, but when deploying to the hardware of Dynap-SE2, weight settings are limited to a 4-bit restricted connection-specific assignment. To convert weight matrices to device configuration, a quantization phase is necessary. The quantization process for Dynap-SE2 involves two steps: defining 4 base weight parameters for the inner product space and storing connection-specific 4-bit binary weight masks in digital memory cells. The goal of quantization is to find a set of common “base weight” parameter values and binary bit-mask matrix that can reconstruct a floating-point weight matrix with minimal deviation. To accomplish this, an auto-encoder structure, a popular unsupervised machine learning method, is used. The intermediate code representation represents the base weight currents, and the decoder weight matrix provides binary bit-masks.

Moreover, in the simulated network, weight values can be positive or negative, representing the synapse’s excitatory or inhibitory behavior. The sign of each weight is used to determine the synapse type for that connection: inhibitory GABA synapses for negative values and excitatory AMPA synapses for positive values. In this way the negative and positive weights produce the desired effect on the post-synaptic membrane. Fig. 12 shows the weight quantization procedure.



**Fig. 13 Translation between behavior and parameters.** The ideal behaviour of a neuron (e.g. a time constant; left) corresponds to setting an accurate analog bias value to control the behaviour of a circuit (right).

The unsupervised auto-encoder training aims to discover a hardware configuration that can replicate the target weight matrix with minimal deviation. The approach used to calculate the matrix reconstruction loss is the mean square error (MSE). This method computes the difference between the absolute values of the original weight matrix and the reconstructed version (Eq. 4).

$$f_{MSE}(W_Q, W) = \frac{1}{N \cdot M} \sum_{i=0}^N \sum_{j=0}^M (W_Q[i, j] - W[i, j])^2 \quad (4)$$

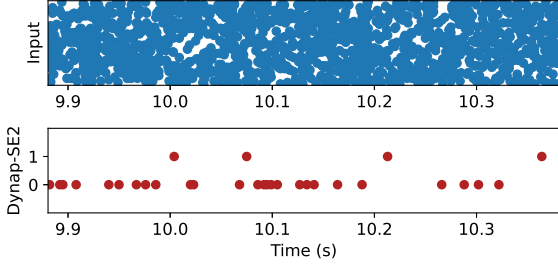
## Deployment

The behavior of a neuron and synapse is characterized by several parameters, including time constants that determine leakage speed and gain ratios that control the amplitude of spike-dependent jumps. While in simulation these parameters can be adjusted mathematically to modify the behaviour of neurons, implementing this parameterisation in VLSI circuits is more complicated. Silicon-based implementations of neurons and synapses rely on adjusting bias voltages and currents. In Rockpool, parameter translation implies translating the behavioral dynamics of a computational neuron model into the VLSI parameters of the corresponding neuron circuits. This involves finding a digital bias generator setting that accurately expresses bias current values in Amperes using empirical lookup tables. Fig. 13 exemplifies what parameter translation means in this reference frame.

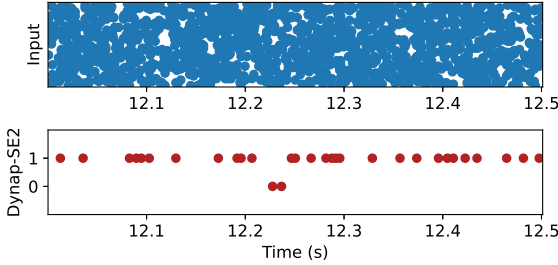
## Results

We used Rockpool and DynapSim to train an SNN as described above, then used the mapping quantization





**Fig. 14** Result of presenting input class 0 to the deployed network on Dynap-SE2. Top: Input events for class 0. Bottom: Output events from Dynap-SE2 evaluating the deployed trained SNN.



**Fig. 15** Result of presenting input class 1 to the deployed network on Dynap-SE2. Top: Input events for class 1. Bottom: Output events from Dynap-SE2 evaluating the deployed trained SNN.

and deployment facilities of Rockpool to configure an SNN on the Dynap-SE2 chip. We converted the frozen noise patterns to real-time AER sequences for injecting into a Dynap-SE2 device. Each event in the noise pattern is encapsulated with its time and address and sent to the chip, where an on-board FPGA circuit converts the AER events to digital pulses that stimulate the synaptic input gates of the neurons. The analog neurons on the Dynap-SE2 then process the inputs and produce output events. Whenever a neuron fires, digital circuits on the FPGA capture the event’s timestamp and source address and encapsulate it as an AER event, which is temporarily stored in buffers implemented inside the FPGA. The output of the hardware evaluation of an input is recorded as AER event sequences.

Figs. 14 and 15 demonstrate that the combined effects of quantization and device mismatch do not result in loss of performance. The hidden temporal information that the network has learned is still present, and the network is able to differentiate between the training samples.

The firing rate ratio (FRR) described in Sec. 2 is utilized to assess the neurons’ ability to distinguish frozen noise patterns. If the FRR is noticeably higher, it indicates that the neuron has made a decision because

one of the neurons fired significantly less and the other one fired considerably more. Tab. 1 displays the firing responses of the simulated, quantized, and hardware networks when presented with frozen noise 0 (FN0), frozen noise 1 (FN1), and the test samples.

1000 random test samples were presented for inference in simulation, and 10 samples on Dynap-SE2 hardware. In all cases, mismatch was either simulated (for “simulated” and “quantized” results), or was present on the Dynap-SE2 hardware device. Test samples produced high output firing rates in general ( $>100$  Hz simulation; TEST in Tab. 1), but with low FRR close to 1. The worst-case FRR for test samples in simulation was 1.5; for the quantized network was 1.6; for inference on the HW was 3.2. Target samples produced similar maximum firing rates, but with higher FRR.

The experimental results indicate that quantization and mismatch cause information loss when converting an optimized network to a hardware configuration. However, despite this loss of information and device mismatch, the decision mechanism remains functional. The hardware deployed model also successfully distinguishes trained input samples, with high FRR.

## Training speed

Gradient-based spiking neural networks training is known to take a long time due to the complexity of the dynamical equations that spiking neurons solve in time, which involve many floating point operations. Additionally, backpropagation through time implies additional memory overhead compared with backpropagation in classical ANN optimisation problems. However, recent advancements in machine learning tools have presented opportunities for performance improvements.

DynapSim utilizes JAX, one of the latest high-performance machine learning tools, to solve dynamical equations and compute surrogate gradients. JAX’s just-in-time (JIT) compilation support for functions written in a purely functional way significantly reduces execution time in the optimization loop. The simulator is designed to comply with JIT constraints, and as a result, JIT reduces the execution time by multiple orders of magnitude.

To illustrate the benefits of JIT, we executed the training script on two different machines, and the performance results are presented in Tab. 2.

Using JAX-JIT reduced the computation time by up more than 3000 $\times$ , making it possible to optimize SNN structures employing complex neuron models

Frozen Noise	Simulated			Quantized			Hardware		
	N0 (Hz)	N1 (Hz)	FRR	N0 (Hz)	N1 (Hz)	FRR	N0 (Hz)	N1 (Hz)	FRR
FN class 0	164	14	11.7	136	58	2.3	18	2	9.0
FN class 1	24	122	5.1	58	144	2.5	0	36	$\infty$
TEST (mean)	138	123	1.1	143	123	1.2	17	14	1.9
TEST (max FRR)	150	100	1.5	154	94	1.6	32	10	3.2

**Table 1 Output firing rates and FRRs for target and test input samples.** “Simulated” results are from PC-based simulation of the trained DynapSim network in Rockpool. “Quantized” results are from PC-based simulation of the quantized model in Rockpool. “Hardware” results are obtained from running inference of the model deployed to Dynap-SE2 hardware. 1000 test samples were used for the Simulated and Quantized results. Only 10 test samples were used for inference on hardware. FN: Frozen Noise (target input sample); TEST: Random poisson test samples; N0: Output neuron for class 0; N1: Output neuron for class 1.

Attribute	Machine 1	Machine 2
CPU	8 Core Apple M1 Pro	Intel Core i7-7500
RAM	32 GB	16 GB
OS	macOS 12.4	Ubuntu 20.04
Epoch/s (JAX)	0.7	0.4
Epoch/s (JAX-JIT)	2600	1350
Duration (JAX)	15 days	28 days
Duration (JAX-JIT)	6.5 minutes	12.5 minutes
Speedup	3714×	3375×

**Table 2 Training time comparison.** The training process was executed identically using the same training code on two machines. We compared the training speed when using non-accelerated JAX, and when using JAX-JIT compilation to the CPU on each machine.

without the need for giant computer clusters or waiting for weeks to see the results. Rockpool / DynapSim is able to use the JIT facilities of JAX to target GPUs and TPUs, enabling scalable use of large computational resources when available, for efficient training of SNNs.

## Conclusion

We demonstrated a new approach and toolchain for gradient-based training and automated deployment of SNN applications to Dynap-SE2. The training pipeline is shown to run 1 million epochs in minutes instead of weeks, by exploiting the just-in-time compilation features of JAX. DynapSim is a huge step towards building commercial SNN applications for mixed-signal neuromorphic processors.

The deployment strategy offers an unsupervised method for weight quantization, removing the need for manual calibration and tuning of hardware bias parameters. The resulting quantized network’s parameters are automatically translated to a hardware configuration. This application is the first of its kind in terms

that a network is optimized offline using backpropagation and automatically deployed to the Dynap-SE2 chip. Although the task introduced here is relatively simple, it proposes a novel methodology for application development targeting Dynap-SE2. The approach, metrics, and evaluation strategies can easily be applied to more complex tasks.

Results indicate that the optimized network is robust to the effects of quantization and device mismatch, which are common challenges in hardware implementation. This is a promising finding that suggests the potential for the use of spiking neural networks in real-world applications where hardware constraints and variability are significant factors. Still, further studies are needed to investigate the network’s performance under different quantization and device mismatch scenarios and to generalize the findings to different spiking neural network architectures and applications.

Rockpool simplifies the modeling and deployment process for the neuromorphic community, addressing a key obstacle that has limited the accessibility of mixed-signal neuromorphic processors to only high-end academic and industrial research. DynapSim paves the way for building commercial applications using mixed-signal neuromorphic technologies.

## References

- [1] Büchel, J., Zendrikov, D., Solinas, S., Indiveri, G., Muir, D.R.: Supervised training of spiking neural networks for robust deployment on mixed-signal neuromorphic processors. *Scientific Reports* **11** (2021). <https://doi.org/10.1038/s41598-021-02779-x>
- [2] Zendrikov, D., Solinas, S., Indiveri, G.: Brain-inspired methods for achieving robust

- computation in heterogeneous mixed-signal neuromorphic processing systems. *bioRxiv* (2022) <https://arxiv.org/abs/https://www.biorxiv.org/content/early/2022/10/27/2022.10.26.513846.full.pdf>. <https://doi.org/10.1101/2022.10.26.513846>
- [3] Muir, D.R., Bauer, F., Weidel, P.: Rockpool Documentation. Zenodo (2019). <https://doi.org/10.5281/zenodo.3773845>
- [4] Çakal, U.: DynapSIM: A fast, optimizable, and mismatch aware mixed-signal neuromorphic chip simulator. Master’s thesis, Middle East Technical University (2022). <https://hdl.handle.net/11511/98616>
- [5] Moradi, S., Qiao, N., Stefanini, F., Indiveri, G.: A scalable multicore architecture with heterogeneous memory structures for dynamic neuromorphic asynchronous processors (DYNAPs). *IEEE Transactions on Biomedical Circuits and Systems* **12**(1), 106–122 (2018). <https://doi.org/10.1109/TBCAS.2017.2759700>
- [6] Livi, P., Indiveri, G.: A current-mode conductance-based silicon neuron for address-event neuromorphic systems. In: 2009 IEEE International Symposium on Circuits and Systems, pp. 2898–2901 (2009). <https://doi.org/10.1109/ISCAS.2009.5118408>
- [7] Bartolozzi, C., Indiveri, G.: Synaptic dynamics in analog VLSI. *Neural Computation* **19**(10), 2581–2603 (2007). <https://doi.org/10.1162/neco.2007.19.10.2581>
- [8] Chicca, E., Stefanini, F., Bartolozzi, C., Indiveri, G.: Neuromorphic electronic circuits for building autonomous cognitive systems. *Proceedings of the IEEE* **102**(9), 1367–1388 (2014). <https://doi.org/10.1109/JPROC.2014.2313954>
- [9] Bradbury, J., Frostig, R., Hawkins, P., Johnson, M.J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., Zhang, Q.: JAX: composable transformations of Python+NumPy programs (2018). <http://github.com/google/jax>
- [10] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings (2015). <http://arxiv.org/abs/1412.6980>
- [11] Lee, J.H., Delbruck, T., Pfeiffer, M.: Training deep spiking neural networks using backpropagation. *Frontiers in Neuroscience* **10** (2016). <https://doi.org/10.3389/fnins.2016.00508>
- [12] Zenke, F., Ganguli, S.: SuperSpike: Supervised Learning in Multilayer Spiking Neural Networks. *Neural Computation* **30**(6), 1514–1541 (2018) [https://arxiv.org/abs/https://direct.mit.edu/neco/article-pdf/30/6/1514/1039264/neco\\_a\\_01086.pdf](https://arxiv.org/abs/https://direct.mit.edu/neco/article-pdf/30/6/1514/1039264/neco_a_01086.pdf). [https://doi.org/10.1162/neco\\_a\\_01086](https://doi.org/10.1162/neco_a_01086)
- [13] Neftci, E.O., Mostafa, H., Zenke, F.: Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks. *IEEE Signal Processing Magazine* **36**(6), 51–63 (2019). <https://doi.org/10.1109/MSP.2019.2931595>
- [14] Kaiser, J., Mostafa, H., Neftci, E.: Synaptic plasticity dynamics for deep continuous local learning (decolle). *Frontiers in Neuroscience* **14** (2020). <https://doi.org/10.3389/fnins.2020.00424>