

EXPLORING THE LIMITATIONS OF LAYER SYNCHRONIZATION IN SPIKING NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural-network processing in machine learning applications relies on layer synchronization. This is practiced even in artificial Spiking Neural Networks (SNNs), which are touted as consistent with neurobiology, in spite of processing in the brain being in fact asynchronous. A truly asynchronous system however would allow all neurons to evaluate concurrently their threshold and emit spikes upon receiving any presynaptic current. Omitting layer synchronization is potentially beneficial, for latency and energy efficiency, but asynchronous execution of models previously trained with layer synchronization may entail a mismatch in network dynamics and performance. We present and quantify this problem, and show that models trained with layer synchronization either perform poorly in absence of the synchronization, or fail to benefit from any energy and latency reduction, when such a mechanism is in place. We then explore a potential solution direction, based on a generalization of backpropagation-based training that integrates knowledge about an asynchronous execution scheduling strategy, for learning models suitable for asynchronous processing. We experiment with 2 asynchronous neuron execution scheduling strategies in datasets that encode spatial and temporal information, and we show the potential of asynchronous processing to use less spikes (up to 50%), complete inference faster (up to 2x), and achieve competitive or even better accuracy (up to $\sim 10\%$ higher). Our exploration affirms that asynchronous event-based AI processing can be indeed more efficient, but we need to rethink how we train our SNN models to benefit from it.

1 INTRODUCTION

Artificial Neural Networks (ANNs) are the foundation behind many of the recently successful developments in AI, such as in computer vision Szegedy et al. (2017); Voulodimos et al. (2018) and natural language processing Vaswani et al. (2017); Brown et al. (2020). To match the complexity of the ever more demanding tasks, networks have grown in size, with advanced large language models having billions of parameters Zhao et al. (2024). With this the power consumption exploded Luccioni et al. (2023), limiting the deployment to large data centers. In an effort to learn from our brain’s superior power efficiency, and motivated in neuroscience research, SNNs Maass (1997) bolster as an alternative. They use discrete binary or graded spikes (events) for communication, are suited for processing sparse features He et al. (2020), and when combined with asynchronous event-based processing are assumed to enhance latency and energy efficiency. Sparsity leads to fewer synaptic operations, resulting in low energy consumption, and asynchronous operation potentiates concurrent evaluation of all neurons in the network purely event-driven, leading to low latency.

Conventional highly parallel ANN compute accelerators, such as Graphics Processing Units (GPUs) and Tensor Processing Units (TPUs), which are primarily optimized for dense vectorized matrix operations, face inherent challenges in exploiting unstructured and temporal sparsity for improving their energy efficiency. Targeting the commonplace practiced way of executing ANNs layer-after-layer has left them with poor support for asynchronous processing too (for improving latency). At best, they parallelize processing within a layer and/or pipeline processing across layers. This leaves an exploration space for neuromorphic processors that try to excel in handling the event-driven nature of SNNs and leverage asynchronous concurrent processing, offering efficiency advantages in various tasks Ivanov et al. (2022); Kang et al. (2020a); Müller-Cleve et al. (2022).

However, despite this advancement, the training of SNNs today very often conveniently relies on conventional end-to-end ANN training methods for performance Dampfhoer et al. (2023), which organize/synchronize computations per-layer rather than event-driven per neuron Guo et al. (2023). Specifically, at any given discrete timestep within a neuron layer, first, the total of all presynaptic currents (from the preceding layer) must be computed and integrated, before postsynaptic neurons in the current layer update their state and evaluate their activation function (i.e. emitting new spikes). For consistency, neuron evaluation in one layer must thus complete and synchronize before proceeding to evaluate neurons of a next layer. This *breadth-first* processing approach (with per layer synchronization), while it facilitates use of vectorized computing hardware (such as GPUs) during training, it introduces dependencies on per-layer synchronization that could impact model accuracy if altered during inference. To avoid this situation, even asynchronous neuromorphic processors, such as Loihi Davies et al. (2018), have integrated mechanisms to ensure (and enforce) layer synchronization.

This leaves a crucial (efficiency) aspect of SNNs relatively unexplored: the ability to allow spiking completely asynchronously across the network without having separate phases for integration and firing, just like in our brain Zeki (2015). In that neurophysical *modus operandi*, neurons can fire and receive currents anywhere in the network at any time, completely independent, a concept we term "*network asynchrony*". Allowing network asynchrony can be advantageous Pfeiffer & Pfeil (2018), as we confirm in our results. Spike activity can quickly propagate deep into the network without being bound by synchronization barriers, thus reducing latency. Furthermore, adhering to layer synchronization could lead to increased computational overhead as the network scales, as suggested by Amdahl's law Rodgers (1985). This implies that the overhead grows non-linearly by adding more computational units to a group that needs to be synchronized at some point in time Yavits et al. (2014). With network asynchrony, such groups can be kept smaller, and the number of synchronization moments can be minimized, potentially reducing the waiting time.

In this paper, we demonstrate this problem and explore solutions. Using a simulation environment that implements the concept of network asynchrony, we provide quantitative results on benchmark datasets (with different spatio-temporal information content) and network topologies of two or more hidden layers, which show the performance degradation and latency/energy inefficiency resulting from changes in model dynamics when trained with layer synchronization and later deployed for asynchronous inference. Next, we explore a potentially promising solution by proposing a generalisation of gradient (backpropagation-based) training, that can be parameterized with various neuron execution scheduling strategies for asynchronous processing, and vectorization abilities present in various neuromorphic processors. We show that using this training method, it is possible not only to recover the compromised accuracy, but also to fulfil the expectation of saving energy and improving latency under asynchronous processing (when compared to the conventional breadth-first processing in GPUs). This work opens a path for design space explorations aimed to bridge the efficiency gap between neural network model training and asynchronous processor design.

2 RELATED WORK

The term "asynchronous processing" is often used whenever neurons can be active in parallel and communicate asynchronously, even if some synchronization protocol is enforced to control the order of spikes and synaptic current processing Rathi et al. (2023); Kang et al. (2020b); Shahsavari et al. (2023); Yousefzadeh et al. (2019). Another kind of asynchrony is related to input coding Guo et al. (2021). In this context, synchronizing means grouping spikes (events) into frames, a topic that has extensively been researched He et al. (2020); Messikommer et al. (2020); Qiao et al. (2024); Ren et al. (2024); Taylor et al. (2024). Neither of the two, however, is the focus of this work. Here, asynchrony refers solely to activity within the network not being artificially bound by any order restriction. This has been researched for simulating biologically accurate neural networks Lytton & Hines (2005); Magalhães et al. (2019; 2020), but remains under-explored in the context of SNNs for machine learning. In this context, processing in layers is just one convenient way Mo & Tao (2022).

Few event-driven neuromorphic processors, such as μ Brain Stuijt et al. (2021) and Speck Caccavella et al. (2023), are fully event-driven and lack any layer synchronization mechanism. This makes them notorious for training out-of-the-box with mainstream model training tools that rely on per-layer synchronization (e.g., PyTorch). Speck developers propose to train their models with hardware in-the-loop Liu et al. (2024) to reduce the mismatch between the training algorithms and the inference

hardware. However this method does not provide a general solution for training asynchronous neural networks. Others, like SpiNNaker Furber et al. (2014) and TrueNorth Akopyan et al. (2015), use timer-based synchronization, or Loihi Davies et al. (2018) and Seneca Tang et al. (2023) use barrier synchronization between layers, in order to warrant that the asynchronous processing dynamics on hardware are aligned at layer boundaries with models trained in software (with layer synchronization). This, however, entails an efficiency penalty, as we will show.

Functionally, the most suitable type of model for end-to-end asynchronous processing is probably rate-coded SNN models, whereby neurons can integrate state and communicate independently from any other neuron. These models are trained like ANNs Zenke & Vogels (2021); Lin et al. (2018); Diehl et al. (2015) or converted from pre-trained ANNs Rueckauer et al. (2017); Kim et al. (2020). Therefore, to be accurate under asynchronous processing, it is required to run the inference for a long time, reducing the latency and energy benefit of using SNNs Sengupta et al. (2019).

Alternative to rate-coding models are temporal-coding models, with time-to-first spike (TTFS) Park et al. (2020); Srivatsa et al. (2020); Kheradpisheh & Masquelier (2020); Comşa et al. (2022) or order encoding Bonilla et al. (2022). They are very sparse (hence energy efficient) but very cumbersome and elaborate to convert from ANNs Painkras et al. (2013); Srivatsa et al. (2020) or train directly Kheradpisheh & Masquelier (2020); Comşa et al. (2022), less tolerant to noise Comşa et al. (2022), and their execution so far, while event-based, requires some form of synchronization or a reference time. Efficiency is thus only attributable to the reduced number of spikes, all of which need to be evaluated in order before a decision is reached. Other alternative encodings for event-based processing of SNNs include phase-coding Kim et al. (2018) and burst coding Park et al. (2019), which are, however, no more economical than rate coding and have not been shown, to our knowledge, to attain competitive performance.

The approach presented in this paper is, in fact, unique in enabling the trainability of models for event-based asynchronous execution, providing efficiency from processing only a subset of spikes, and delivering consistent performance and tolerance to noise. Also relevant to the works in this paper are SparseProp Engelken (2024) and EventProp Wunderlich & Pehle (2021) on efficient event-based simulation and training, respectively. EventProp Wunderlich & Pehle (2021) is potentially more economical than discrete-time backpropagation for training event-based models for asynchronous processing (and fully compatible with the work in this paper), but it has not been shown how its complexity scales beyond 1 hidden layer. SparseProp Engelken (2024) proposes an efficient neuron execution scheduling strategy for asynchronous processing in software simulation. An effective hardware implementation of this scheduling strategy can be found in Monti et al. (2017), which has inspired the herein proposed “*momentum schedule*”. Moreover, we advance this research by demonstrating how to train SNN models using this event scheduler.

3 METHODS FOR SIMULATING AND TRAINING OF ASYNCHRONOUS SNNs

In this section, we provide an overview of our methodology. We start by showing how one can simulate asynchronous SNN processing (implementing network asynchrony) and then we explain how this can be integrated in back-propagation training, to prepare good performing and efficient asynchronous models.

3.1 SIMULATING ASYNCHRONOUS SNNs

The SNNs used in this work consist of L layers of Leaky Integrate and Fire (LIF) neurons He et al. (2020), with each layer l for $1 \leq l \leq L$ having $N^{(l)}$ neurons and being fully-connected to the next layer $l + 1$, except for the output layer. Each input feature is connected to all neurons in the first layer. Every connection has a synaptic weight. Using these weights, we can compute the incoming current x per neuron resulting from spikes, as per equation 3.

Each LIF neuron has a membrane potential exponentially decaying over time based on some membrane time constant τ_m . We use the analytical solution (see equation 4) to compute the decay. By keeping track of the membrane potential $u[t]$ and the elapsed time Δt since time t , it is possible to precisely calculate the membrane potential for $t + \Delta t$. Therefore, computations are required only when $x[t + \Delta t] > 0$.

To determine if a neuron spikes, a threshold function Θ checks if the membrane potential exceeds a threshold U_{thr} , following equation 5. If $\Theta(u) = 1$, the membrane potential is hard reset to 0. Soft reset is another option Guo et al. (2022) as well as refractoriness Sanaullah et al. (2023). However for simplicity here we only experiment with hard reset.

3.1.1 EVENT-DRIVEN STATE UPDATES

We can vectorize the computations introduced in the previous section to perform them on a per-layer basis. This is herein referred to as the "layered" inference approach, which implies layer synchronization, since all neurons of one layer need to evaluate their state, before any neuron in a down-stream layer does the same. An alternative to this situation is an event-driven approach, where the computations for state updates are applied in response to new spike arrivals. If additionally spikes can be emitted (threshold evaluation) at any point in time (e.g in response to any individual synaptic current), by any neuron in the network, affecting the states of postsynaptic neurons independently of others, then (we assume) this approach can achieve a true representation of network asynchrony.

When using network asynchrony, the network dynamics evolve with each spike. Any single spike can generate multiple currents downstream, each linked to a specific stimuli at timestep t , determined by the initiating input activity. Let K_t represent the atomic computation that is executed at one neuron in the network, in response to stimulus presented at the network at timestep t , and which will update the state of the neuron and evaluate its threshold function (activation). The exact computations for a LIF neuron can be found in section A.2. The order in which K_t is executed across neurons inside the entire network (and not just in one layer) affects the output of the network due to the non-linearity in K_t (see figure 1 for an example of the effect this can have). Here for simplicity, it is assumed that the fan-out currents resulting from the same spike are processed as a group at once (e.g. no heterogeneous propagation delays exist), and processing is done when an emitted spike is scheduled for processing, i.e., when the spike is propagated. This may not necessarily be the same moment as when the spike was emitted due to network asynchrony. Because of these assumptions, analyzing the spike propagation order gives the same insights as analyzing the execution order of K_t (across neurons in the network).

Unlike the layered approach, in this modus operandi because of the per-input current evaluation of the firing threshold neurons could theoretically fire more than once at the same timestep t (in response to the same input stimulus). We prevent this by inducing a neuron to enter a refractory state for the remainder of time t . In this state, the neuron keeps its membrane potential at 0. This makes the model causally simpler, more tractable, and more "economic" when energy consumption is coupled to spike communication. These aspects also may also give an explicit retrospectively justification for the neurophysically observed refractoriness Berry & Meister (1997).

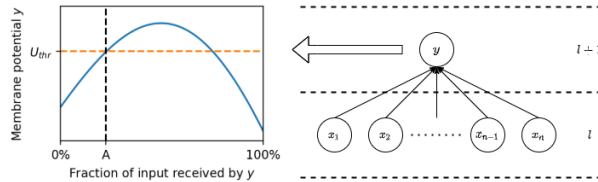


Figure 1: With network asynchrony, currents from neurons x_i in layer l may come to neuron y in layer $l + 1$ at any time in any order. It may spike before all inputs have been received in a discrete timestep. Illustrated is the exceedance of the firing threshold at point A . All inputs, and thus information after A is missed, i.e., the spiking decision is made based on partial information. In this example, if layer synchronization is enforced, neuron y would wait for 100% of the input and will not spike.

Two resolutions of time We apply input framing as typically done for inputs from a DVS sensor, accumulating timestamped events in discrete time-bins, and providing these frames as inputs to the network in discrete *timesteps*. The *timestep size* refers to the time-interval between the timestamps of the first (or last) events in the time-bins of two consecutive timesteps. Timesteps correspond to one concept of timing (resolution) with same semantics as in RNNs, and which we refer to here as *macro-time*. At each timestep, input event frames initiate network activity, in a so-called *forward pass*, making all new spikes contingent on the timestep.

The order of spikes propagated through the network during a timestep establishes another notion of “time” in a forward pass (we call it “micro-time”), which is discretized in so-called *forward steps*. A forward step is associated with the processing of one or a group of events (in case of vectorization), and the number of forward steps measures time to complete inference. One can realize that in case of layered inference the number of forward steps is fixed and equals the number of layers, but this is not the case for asynchronous processing.

3.1.2 VECTORIZED NETWORK ASYNCHRONY

The event-driven (neuron state) update rules for network asynchrony as introduced in the previous section can be vectorized by selecting a number of spikes for processing them simultaneously. This allows us to consider the entire spectrum of possibilities between per-layer synchronization at one extreme (by assuming a vector size equal to a layer size), and “complete” asynchrony at the other extreme, where each spike event is processed entirely independently of all others. Additionally, vectorization makes acceleration possible by exploiting the parallelization features and vector pipelines of accelerators, where these models execute, leading to pragmatic simulation of network asynchrony.

During the simulation, the states of all the $N = \sum_{l=1}^L N^{(l)}$ neurons in the network are stored in vectors. Vector $\mathbf{x} \in \mathbb{R}^N$ tracks the computed input currents for the neurons, $\mathbf{u} \in \mathbb{R}^N$ the membrane potential of the neurons, $\mathbf{s} \in \mathbb{N}^N$ the emitted spikes awaiting processing, and $\mathbf{c} \in \mathbb{N}^N$ which neurons have spiked in the current forward pass. For any of those vectors, the indices from $\sum_{k=1}^{l-1} N^{(k)}$ to $\sum_{k=1}^l N^{(k)}$ represent the values for the neurons in layer l , for $1 \leq l \leq L$ where $N^{(0)} = 0$.

Algorithm 1 outlines the processing during a forward pass (propagation of the spikes of an input frame through the network). Input spikes are available in $\mathbf{s}_{\text{in}} \in \mathbb{N}^{N_{\text{in}}}$ (not to be confused with \mathbf{s}) where N_{in} is the number of input features. Each *forward pass* consists of *forward steps* (the code within the while loop), which update the state of a set of neurons based on the spikes selected for propagation. A complete list of parameters can be found in section A.5.

Algorithm 1 Vectorized network asynchrony forward pass

Input: input spikes $\mathbf{s}_{\text{in}} \in \mathbb{N}^{N_{\text{in}}}$, previous forward pass time $t_0 \in \mathbb{R}$, current forward pass time $t_1 \in \mathbb{R}$, neuron state $\mathbf{u} \in \mathbb{R}^N$ at time t_0 , forward group size $F \in \mathbb{N}_{>0}$

Output: spike count $\mathbf{c} \in \mathbb{N}^N$

$\Delta t \leftarrow t_1 - t_0$

$\mathbf{u} \leftarrow \text{NeuronDecay}(\mathbf{u}, \Delta t)$

$\mathbf{x} \leftarrow \text{InputLayerForward}(\mathbf{s}_{\text{in}})$

$\mathbf{s} \leftarrow \mathbf{0}^N$

▷ Vector with zeros of length N

$\mathbf{c} \leftarrow \mathbf{0}^N$

while ($\text{Sum}(\mathbf{x}) \neq 0$ or $\text{Sum}(\mathbf{s}) \neq 0$) and $\neg \text{EarlyStop}(\mathbf{s})$ **do**

$\mathbf{s}_{\text{new}}, \mathbf{u} \leftarrow \text{NeuronForward}(\mathbf{x}, \mathbf{u})$

$\mathbf{s} \leftarrow \mathbf{s} + \mathbf{s}_{\text{new}}$

▷ Enqueue new spikes

$\mathbf{c} \leftarrow \mathbf{c} + \mathbf{s}_{\text{new}}$

▷ Update spike count

$\mathbf{s}_{\text{selected}} \leftarrow \text{SelectSpikes}(\mathbf{s}, F)$

$\mathbf{s} \leftarrow \mathbf{s} - \mathbf{s}_{\text{selected}}$

▷ Dequeue selected spikes

$\mathbf{x} \leftarrow \text{NetworkLayersForward}(\mathbf{s}_{\text{selected}})$

end while

The forward pass ends either when all spike activations have been processed (“On spiking done” stop condition), or when any neuron in the output layer has spiked one or more (default is one) forward steps ago (“On output” stop condition). The latter is checked in the `EarlyStop` function.

The `SelectSpikes` function defines how to select a subset of the emitted spikes for propagation. The function selects F spikes at a time, or less if there are less than F remaining spikes ready to be propagated. The method of spike selection is determined by a *scheduling policy*.

For the experiments here two policies are exercised. The intend is that different scheduling strategies can be tested for their effectiveness in capturing tractably asynchronous processing dynamics. The first, Random Scheduling (RS), randomly picks spikes from the entire network. The second, Momentum Scheduling (MS), prioritizes spikes from neurons based on their membrane potential upon exceeding their threshold.

The neuron model-specific (LIF) behavior is expressed in the `NeuronDecay` and `NeuronForward` functions. This entails the computations for state updates (see section 3.1) for all neurons in the network (`NeuronDecay`) or for only those neurons receiving input currents in the forward step (`NeuronForward`), with added restriction that spiking is only allowed once per forward pass.

The network architecture is defined by the `InputLayerForward` and `NetworkLayersForward` functions. These functions compute the values of synaptic currents from spikes. This follows from equation 3.

Imitating neuromorphic accelerator hardware Neuromorphic processors come in an number of variations, but most of them have a template architecture that interconnects many tiny processing cores with each other. This design enables a scalable architecture that supports fully distributed memory and compute systems. In this template, each processing core has its own small synchronous execution domain, but the cores operate asynchronously among each other. For instance, the architectures of Loihi Davies et al. (2018), SENECA Yousefzadeh et al. (2022), SpiNNaker Painkras et al. (2013), Speck Caccavella et al. (2023), TrueNorth Akopyan et al. (2015), POETS Shahsavari et al. (2021) and μ Brian Stuijt et al. (2021) all follow this template. Our asynchronous processing simulation environment imitates the behavior of a generalized neuromorphic processor by simulating the asynchronous interactions of several tiny synchronous domains and can be also configured to reproduce more specialized intrinsic behaviors of many of those processors (see table 5).

3.2 TRAINING ASYNCHRONOUS SNNs

We use backpropagation to train the model weights, and specifically when stimulus is presented across multiple timesteps (forward passes), such as for sequential or temporal data, then this is Backpropagation Through Time (BPTT). Following common practise to address the non-differentiability of the threshold function, the surrogate gradient method is used Zenke & Vogels (2021). Specifically the arctan function Fang et al. (2021), (see section A.3 for details) provides a continuous and smooth approximation of the threshold function.

Class prediction is based on the softmaxed membrane potentials (for CIFAR-10) or spike counts over time (for the other datasets) of the neurons in the output layer, as described in section A.4. The loss is minimized using the Adam optimizer, with $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

3.2.1 UNLAYERED BACKPROPAGATION

We refer to “conventional” SNN training with per-layer synchronization using backpropagation Dampfhofer et al. (2023), as “layered backpropagation”.

The vectorized network asynchronous processing approach is differentiable as well, and can be used with backpropagation. We refer to this as “unlayered backpropagation”. Combined with BPTT unrolling, this method implies a two level unrolling. At the outer level, unrolling is based on discretization of time and thus the input across (as usually fixed number of) timesteps. At the inner level, unrolling is based on the F -grouping (and vectorized processing) of spikes in forward steps, subject to the scheduling policy, applied to the emitted spikes by neurons; from the beginning of the current timestep up-until the output is read.

Note, that the dependence on the scheduling policy, applied on a variable number of activations (across timesteps and data samples), in F -sized groups is the fundamental difference from layered back-propagation, and leads to different gradient state being built up in the computation graph. This state now captures the dynamics of asynchronous processing. For a single *backward pass* in a timestep of BPTT, which is applied in a similar way as the layered backpropagation equivalent, it is given that:

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial c_t} \sum_{i=1}^{N_t} \left(\frac{\partial c_t}{\partial s_i} \frac{\partial s_i}{\partial W} \right) \quad (1)$$

where t is the time(step) of the *forward pass*, W refers to the trainable weights, L_t is the loss, c_t is the spike count at the end of the forward pass, and s_i is the emitted spikes vector at the end of the forward step i . The overall gradient (state) depends on the total number of *forward steps* N_t in the forward pass and the spikes processed in each forward step. The number of steps scales linearly with the number of spikes processed in the forward pass, and it is important to understand that due to asynchronous processing the number of spikes processed until the evaluation of the loss function may

be (well) less than the total number of spikes emitted throughout the forward pass. Since for every forward step, the computations are repeated, the time complexity scales linearly with the number of forward steps, $O(N_t)$. The same applies to the space complexity.

During the backward pass, the spikes in s_i which are not selected for processing can skip the computation f for that step:

$$\frac{\partial s_i}{\partial W} = \frac{\partial f(s_{i-1;\text{selected}}, u_{i-1}, x_t, W)}{\partial W} + \frac{\partial s_{i-1;\text{not selected}}}{\partial W} \quad (2)$$

Skip connections have been researched in deep ANNs and identified as a contributor to the stability of the training process Orhan & Pitkow (2018). This may apply to the skip in unlayered backpropagation as well. To what extent this is the case is not explored in this work.

Note that under this generalization, layered backpropagation corresponds to F -group the size of a layer, populated with a static round-Robin execution schedule following neuron index order within a layer, and re-reset across consecutive layers).

3.2.2 REGULARIZATION TECHNIQUES

During training, we use regularization to prevent overfitting and/or enhance model generalization. These techniques are not used during inference.

Input spike dropout. Randomly omits input spikes with a given probability. The decision to drop each spike is independent according to a Bernoulli distribution.

Weight regularization. Adds weight decay to the loss function: $L_{\lambda_W}(\mathbf{W}) = L(\mathbf{W}) + \lambda_W \|\mathbf{W}\|_2^2$ where λ_W is the regularization coefficient, L is the loss given the weights, and \mathbf{W} are all the weights.

Refractory dropout. With some probability, do not apply the refractory effect, allowing a neuron to fire again within the same forward pass.

Momentum noise. When using momentum scheduling, noise sampled from $U(0, 1)$ and multiplied by some constant λ_{MS} is added to the recorded membrane potential while selecting spikes.

4 RESULTS

4.1 EXPERIMENTAL SETUP

We carried out our experiments on SNN models trained primarily in three common benchmarking datasets, each of them has different structure: N-MNIST Orchard et al. (2015), SHD Cramer et al. (2020), and DVS gestures Amir et al. (2017). N-MNIST has purely spatial structure, SHD purely temporal, and DVS gestures combines both spatial and temporal (input framing in DVS gesture is done such that an entire gesture motion and contour is not revealed within a single frame). We also repeat some of the experiments with a fourth dataset, CIFAR-10 Krizhevsky (2012), in a more complex VGG-style network, which we will discuss in section 5. More details on the datasets and network topologies, can be found in section A.6.1. Table 1, summarizes the parameterization of the experiments and the reported results. The network architecture and hyperparameters are given in table 6. State-of-the art performance for these tasks can be achieved with reasonably shallow and wide models. We chose however to train narrow, but deeper network architectures so that the effects of absence of layer synchronization can be revealed in the comparison (because of this bias in network topology the accuracy results shown can vary from the top state of the art). In section A.7, we also provide an additional ablation with regard to how various training hyperparameters affect accuracy (forward group size, refractory dropout, and momentum noise during training).

4.2 NETWORK ASYNCHRONY INCREASES NEURON REACTIVITY

As observed in figure 2 (top row) during asynchronous inference, neurons are more reactive, i.e. a neuron can spike after integrating only a small number of incoming currents. With layer synchronization, this effect is averaged out as neurons are always required to consider all presynaptic currents (which are more likely to cancel each other out). For inference with $F = 8$, artifacts in the number of currents to spike can be observed, because each forward step propagates currents resulting from 8

Table 1: Parameterization of experiments and results.

Parameterization	Description
Training method	Training with layered backpropagation is marked as "Layered" and with unlayered backpropagation as "Unlayered [scheduling strategy]".
Inference method	Inference with layer synchronization is marked as "Layered" and with network asynchrony as "Async [scheduling strategy]".
Scheduling strategy	How to select spikes from the queue. Can be random ("RS"), or based on the membrane potential just before spiking ("MS").
Forward group size F	Number of spikes to select for processing at the same time. Default is 8, both during training and inference.
Stop condition	Forward pass terminates: If all network activity drains ("On spiking done") or one forward step after the first spike is emitted by the output layer ("On output"). The default is "On spiking done" during training and "On output" during inference.

spikes, causing neurons to integrate more currents than necessary for firing. Similar artifacts might also occur in neuromorphic chips equipped with fixed-width vector processing pipelines; making these observations insightful into the behavior of such hardware.

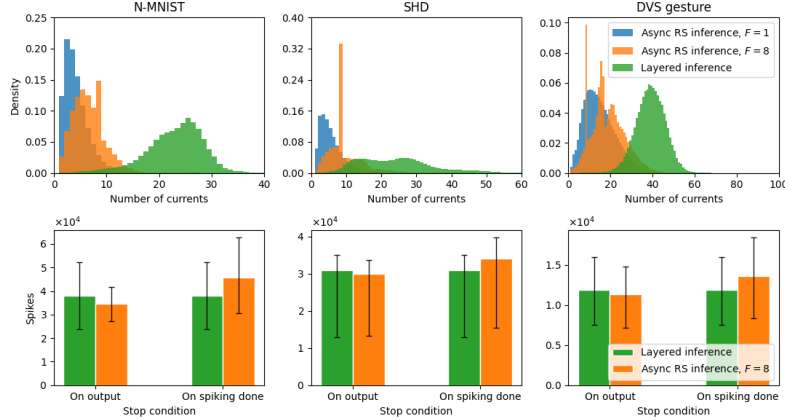


Figure 2: **(Top row)** Number of currents integrated by a neuron before spiking, recorded per neuron and per forward pass for all samples and neurons (excluding the neurons in the input layer). The Y-axis shows the relative frequency of the number of currents integrated before spiking. **(Bottom row)** Mean number of spikes per neuron during inference of all samples in the test. Error bars show the 25th and 75th percentiles. (Models in this figure were trained with layered backpropagation.)

One expects that more reactive neurons imply higher activation density. Interestingly, this is not necessarily the case for the models trained for asynchronous inference! In figure 2 (bottom row) we see that if we wait for the network to “drain” of spike activity during a forward pass the total number of spikes will indeed be higher. But if the forward pass terminates as soon as a decision is made, asynchronous models are consistently sparser. This is because asynchronous processing allows spike activity to freely flow through to the output and not be blocked at every layer for synchronization.

4.3 UNLAYERED BACKPROPAGATION RECOVERS ACCURACY AND INCREASES SPARSITY

Network asynchrony negatively affects the performance of the models trained with layered backpropagation in all three datasets (table 2). This is likely the “Achilles’ heel” of neuromorphic processing today. However, we observe that the accuracy loss is remediated when training takes into account asynchronous processing dynamics, which also significantly increases sparsity (by about 2x). We witnessed that the accuracy of models trained and executed asynchronously is consistently superior under the two scheduling policies we considered. This result conjectures that neuromorphic AI is competitive and more computationally efficient.

Figure 3 reveals another interesting result. It depicts how accuracy evolves as we allow more forward steps in the forward pass after the initial output during asynchronous inference. We see that because of the free flow of key information *depth-first*, models trained with unlayered backpropagation obtain the correct predictions as soon as the output layer gets stimulated. Activity that is likely triggered by “noise” in the input is integrated later on. Momentum scheduling is particularly good at exploiting this to boost accuracy.

Table 2: Accuracy and activation density results. More details about these metrics in section A.6.2.

Training	Inference	N-MNIST		SHD		DVS gesture	
		Acc.	Density	Acc.	Density	Acc.	Density
Layered	Layered	0.949	3.987	0.783	14.386	0.739	46.473
Layered	Async RS	0.625	3.652	0.750	13.905	0.701	44.140
Unlayered RS	Async RS	0.956	1.504	0.796	5.100	0.777	25.140
Unlayered MS	Async MS	0.963	1.476	0.816	6.224	0.856	26.686

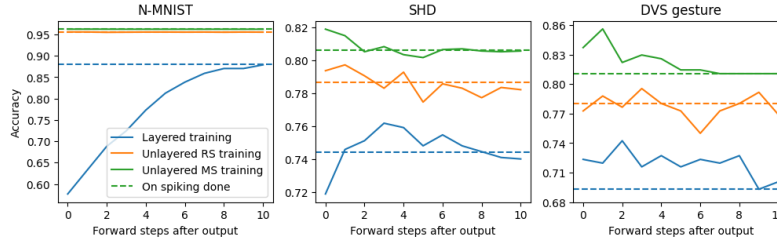


Figure 3: Accuracy as function of forward steps after the first spike in the output layer. Given that $F = 8$, each extra forward step processes another 8 spikes, assuming enough spikes are available. Dashed lines show the accuracy after all spike activity has been “drained” out of the network.

4.4 NETWORK ASYNCHRONY AND ENERGY EFFICIENCY

To provide a quantification of energy savings from asynchronous processing in the trained models, table 3 reports the mean number of synaptic operations and energy consumption on the μ Brain neuromorphic chip Stuijt et al. (2021) per sample.

Training	Inference	N-MNIST		SHD		DVS gesture	
		SOs	Energy	SOs	Energy	SOs	Energy
Layered	Layered	3.521	0.9155	50.82	13.21	158.8	41.29
Layered	Async RS	3.225	0.8386	49.12	12.77	150.9	39.22
Unlayered RS	Async RS	1.328	0.3454	18.02	4.684	85.92	22.34
Unlayered MS	Async MS	1.304	0.3389	21.99	5.717	91.2	23.71

Table 3: Mean num of Synaptic Operations (SOs $\times 10^4$) and energy consumption (μ J) for classification of one sample, using 26 pJ/SO as measured for μ Brain Stuijt et al. (2021); ignoring static power.

4.5 NETWORK ASYNCHRONY REDUCES LATENCY

An equally important result concerns the inference latency reduction under asynchronous network processing. The models trained with unlayered backpropagation have significantly lower latency than those trained with layer synchronization. Assuming a unit latency for processing one spike, the comparatively worst case latency will be given under sequential processing, as a function of the number of spikes processed until an inference decision is made. Figure 4 shows a distribution of the inference latencies across the entire test-set. It should be clear by now that this is because “the important spikes” in these models quickly reach the output layer, uninhibited by layer synchronization.

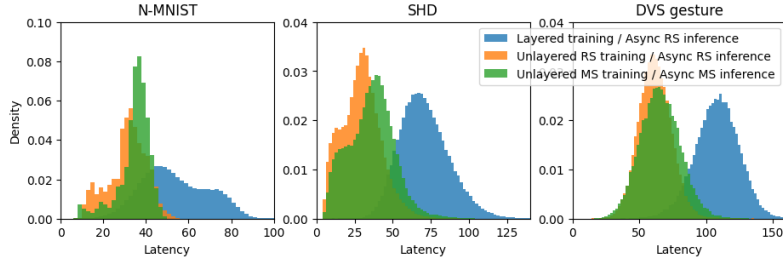


Figure 4: Latency per forward pass for all samples, in number of spikes until a decision in the output layer is reached. The Y-axis shows the relative frequency of recorded latencies.

4.6 UNLAYERED BACKPROPAGATION IS RESOURCE-INTENSIVE

We also tried to confirm the results in a scaled up setup, namely with a deeper VGG-7 like network (details in A.9), trained on the CIFAR-10 dataset. Figure 5a confirms the problem when training with layered backpropagation and the running asynchronous inference. When we try to train asynchronous models with unlayered backpropagation and the simpler RS scheduling policy the memory cost however becomes prohibitive unless we substantially increase the forward group size F . Nevertheless, even with as large as $F = 512$ (during training, the smallest possible with an NVIDIA RTX 5000 GPU) and tested with $F = 64$, accuracy steeply recovers to 71%, while activation density drops to about 1/4 (confirming our previous observations). We anticipate that with much smaller F during training (higher degree of network asynchrony), the accuracy can be completely recovered.

Unfortunately, the computational cost of unlayered backpropagation, in the current framework (and implementation) is rather high, especially when training with a small F . Each time F is halved, the time and memory requirements approximately double as discussed in sections 3.2.1 and A.8.

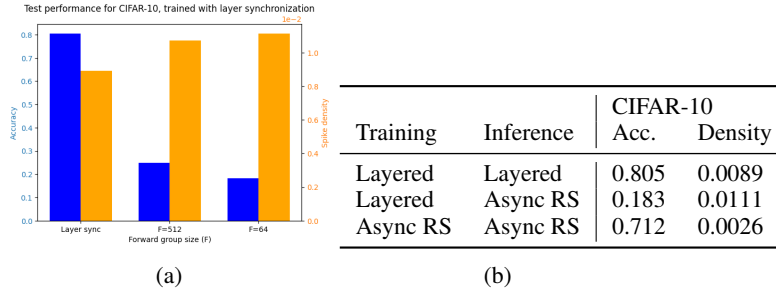


Figure 5: (a) Impact of network asynchrony on a CIFAR-10 trained VGG-style model. A larger decrease in accuracy occurs with smaller forward group sizes accompanied by an increase in activations. (b) Accuracy and activation density results on CIFAR-10

5 DISCUSSION

We pinpoint a crux for neuromorphic AI processing in training SNN models with conventional backpropagation and then naively assuming that they will execute consistently, and energy and latency efficiently in asynchronous event-driven neuromorphic processors. We identified the source of this problem in neglecting the dynamics of asynchronous processing, and we found a way to factor them in the gradient training process. The resulting models not only tend to recover the affected performance (even exceed it) but also deliver the anticipated energy and latency benefits touted by neuromorphic computing. Execution scheduling strategies may be a missing bridge between neural algorithms/models and neuromorphic hardware architecting. This study has merely scratched the surface of a research exploration in this direction, but our observations is that unless we rethink our training methods (take into account asynchronous processing dynamics), and design choices of neuromorphic accelerators (layer synchronization primitives considered harmful), SNNs and brain-inspired computing could fail to deliver in the long-term either in performance or efficiency.

REPRODUCIBILITY STATEMENT

The methods and appendices should provide enough information to reproduce the results. To help with the reproducibility, we also provided code.

REFERENCES

- Filipp Akopyan, Jun Sawada, Andrew Cassidy, Rodrigo Alvarez-Icaza, John Arthur, Paul Merolla, Nabil Imam, Yutaka Nakamura, Pallab Datta, Gi-Joon Nam, et al. Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neuromorphic chip. *IEEE transactions on computer-aided design of integrated circuits and systems*, 34(10):1537–1557, 2015.
- Arnon Amir, Brian Taba, David Berg, Timothy Melano, Jeffrey McKinstry, Carmelo Di Nolfo, Tapan Nayak, Alexander Andreopoulos, Guillaume Garreau, Marcela Mendoza, et al. A low power, fully event-based gesture recognition system. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7243–7252, 2017.
- Michael Berry and Markus Meister. Refractoriness and neural precision. *Advances in neural information processing systems*, 10, 1997.
- Lina Bonilla, Jacques Gautrais, Simon Thorpe, and Timothée Masquelier. Analyzing time-to-first-spike coding schemes: A theoretical approach. *Frontiers in Neuroscience*, 16:971937, 09 2022. doi: 10.3389/fnins.2022.971937.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Caterina Caccavella, Federico Paredes-Vallés, Marco Cannici, and Lyes Khacef. Low-power event-based face detection with asynchronous neuromorphic hardware. *arXiv preprint arXiv:2312.14261*, 2023.
- Iulia-Maria Comşa, Krzysztof Potempa, Luca Versari, Thomas Fischbacher, Andrea Gesmundo, and Jyrki Alakuijala. Temporal coding in spiking neural networks with alpha synaptic function: Learning with backpropagation. *IEEE Transactions on Neural Networks and Learning Systems*, 33(10):5939–5952, 2022. doi: 10.1109/TNNLS.2021.3071976.
- Benjamin Cramer, Yannik Stradmann, Johannes Schemmel, and Friedemann Zenke. The heidelberg spiking data sets for the systematic evaluation of spiking neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 33(7):2744–2757, 2020.
- Manon Dampfhofer, Thomas Mesquida, Alexandre Valentian, and Lorena Anghel. Backpropagation-based learning techniques for deep spiking neural networks: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A neuromorphic manycore processor with on-chip learning. *Ieee Micro*, 38(1):82–99, 2018.
- Peter U. Diehl, Daniel Neil, Jonathan Binas, Matthew Cook, Shih-Chii Liu, and Michael Pfeiffer. Fast-classifying, high-accuracy spiking deep networks through weight and threshold balancing. In *2015 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2015. doi: 10.1109/IJCNN.2015.7280696.
- Rainer Engelken. Sparseprop: Efficient event-based simulation and training of sparse recurrent spiking neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Wei Fang, Zhaofei Yu, Yanqi Chen, Timothée Masquelier, Tiejun Huang, and Yonghong Tian. Incorporating learnable membrane time constant to enhance learning of spiking neural networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2661–2671, 2021.

- S. B. Furber, F. Galluppi, S. Temple, and L. A. Plana. The spinnaker project. *Proceedings of the IEEE*, 102(5):652–665, 2014. doi: 10.1109/JPROC.2014.2304638.
- Wenzhe Guo, Mohammed E. Fouda, Ahmed M. Eltawil, and Khaled Nabil Salama. Neural coding in spiking neural networks: A comparative study for robust neuromorphic systems. *Frontiers in Neuroscience*, 15, 2021. ISSN 1662-453X. doi: 10.3389/fnins.2021.638474. URL <https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2021.638474>.
- Yufei Guo, Yuanpei Chen, Liwen Zhang, YingLei Wang, Xiaode Liu, Xinyi Tong, Yuanyuan Ou, Xuhui Huang, and Zhe Ma. Reducing information loss for spiking neural networks. In *European Conference on Computer Vision*, pp. 36–52. Springer, 2022.
- Yufei Guo, Xuhui Huang, and Zhe Ma. Direct learning-based deep spiking neural networks: a review. *Frontiers in Neuroscience*, 17:1209795, 2023.
- Weihua He, YuJie Wu, Lei Deng, Guoqi Li, Haoyu Wang, Yang Tian, Wei Ding, Wenhui Wang, and Yuan Xie. Comparing snns and rnns on neuromorphic vision datasets: Similarities and differences. *Neural Networks*, 132:108–120, 2020.
- Dmitry Ivanov, Aleksandr Chezhegov, Mikhail Kiselev, Andrey Grunin, and Denis Larionov. Neuro-morphic artificial intelligence systems. *Frontiers in Neuroscience*, 16:1513, 2022.
- Laxmi R Iyer, Yansong Chua, and Haizhou Li. Is neuromorphic mnist neuromorphic? analyzing the discriminative power of neuromorphic datasets in the time domain. *Frontiers in neuroscience*, 15: 608567, 2021.
- Minseon Kang, Yongseok Lee, and Moonju Park. Energy efficiency of machine learning in embedded systems using neuromorphic hardware. *Electronics*, 9(7):1069, 2020a.
- Ziyang Kang, Lei Wang, Shasha Guo, Rui Gong, Shiming Li, Yu Deng, and Weixia Xu. Asie: An asynchronous snn inference engine for aer events processing. *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, 16(4):1–22, 2020b.
- Saeed Reza Kheradpisheh and Timothée Masquelier. Temporal backpropagation for spiking neural networks with one spike per neuron. *International Journal of Neural Systems*, 30, 03 2020. doi: 10.1142/S0129065720500276.
- Jaehyun Kim, Heesu Kim, Subin Huh, Jinho Lee, and Kiyoun Choi. Deep neural networks with weighted spikes. *Neurocomputing*, 311:373–386, 2018. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2018.05.087>. URL <https://www.sciencedirect.com/science/article/pii/S0925231218306726>.
- Seijoon Kim, Seongsik Park, Byunggook Na, and Sungroh Yoon. Spiking-yolo: Spiking neural network for energy-efficient object detection. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:11270–11277, 04 2020. doi: 10.1609/aaai.v34i07.6787.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- Chit-Kwan Lin, Andreas Wild, Gautham N. Chinya, Yongqiang Cao, Mike Davies, Daniel M. Lavery, and Hong Wang. Programming spiking neural networks on intel’s loihi. *Computer*, 51(3):52–61, 2018. doi: 10.1109/MC.2018.157113521.
- Yuhang Liu, Tingyu Liu, Yalun Hu, Wei Liao, Yannan Xing, Sadique Sheik, and Ning Qiao. Chip-in-loop snn proxy learning: a new method for efficient training of spiking neural networks. *Frontiers in Neuroscience*, 17:1323121, 2024.
- Alexandra Sasha Luccioni, Sylvain Viguier, and Anne-Laure Ligozat. Estimating the carbon footprint of bloom, a 176b parameter language model. *Journal of Machine Learning Research*, 24(253): 1–15, 2023.
- William W Lytton and Michael L Hines. Independent variable time-step integration of individual neurons for network simulations. *Neural computation*, 17(4):903–921, 2005.

- Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 10(9):1659–1671, 1997.
- Bruno Magalhães, Michael Hines, Thomas Sterling, and Felix Schürmann. Fully-asynchronous fully-implicit variable-order variable-timestep simulation of neural networks. In *Computational Science–ICCS 2020: 20th International Conference, Amsterdam, The Netherlands, June 3–5, 2020, Proceedings, Part V 20*, pp. 94–108. Springer, 2020.
- Bruno RC Magalhães, Thomas Sterling, Michael Hines, and Felix Schürmann. Fully-asynchronous cache-efficient simulation of detailed neural networks. In *Computational Science–ICCS 2019: 19th International Conference, Faro, Portugal, June 12–14, 2019, Proceedings, Part III 19*, pp. 421–434. Springer, 2019.
- Nico Messikommer, Daniel Gehrig, Antonio Loquercio, and Davide Scaramuzza. Event-based asynchronous sparse convolutional networks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*, pp. 415–431. Springer, 2020.
- Lingfei Mo and Zhihan Tao. Evtstnn: Event-driven snn simulator optimized by population and pre-filtering. *Frontiers in Neuroscience*, 16:944262, 2022.
- Massimo Monti, Manolis Sifalakis, Christian F. Tschudin, and Marco Luise. On hardware programmable network dynamics with a chemistry-inspired abstraction. *IEEE/ACM Transactions in Networking*, 25(4):2054–2067, aug 2017. ISSN 1063-6692. doi: 10.1109/TNET.2017.2674690. URL <https://doi.org/10.1109/TNET.2017.2674690>.
- Simon F Müller-Cleve, Vittorio Fra, Lyes Khacef, Alejandro Pequeño-Zurro, Daniel Klepatsch, Evelina Forno, Diego G Ivanovich, Shavika Rastogi, Gianvito Urgese, Friedemann Zenke, et al. Braille letter reading: A benchmark for spatio-temporal pattern recognition on neuromorphic hardware. *Frontiers in Neuroscience*, 16:951164, 2022.
- Garrick Orchard, Ajinkya Jayawant, Gregory K Cohen, and Nitish Thakor. Converting static image datasets to spiking neuromorphic datasets using saccades. *Frontiers in neuroscience*, 9:437, 2015.
- Emin Orhan and Xaq Pitkow. Skip connections eliminate singularities. In *International Conference on Learning Representations*, 2018.
- Eustace Painkras, Luis A Plana, Jim Garside, Steve Temple, Francesco Galluppi, Cameron Patterson, David R Lester, Andrew D Brown, and Steve B Furber. Spinnaker: A 1-w 18-core system-on-chip for massively-parallel neural network simulation. *IEEE Journal of Solid-State Circuits*, 48(8): 1943–1953, 2013.
- Seongsik Park, Seijoon Kim, Hyeokjun Choe, and Sungroh Yoon. Fast and efficient information transmission with burst spikes in deep spiking neural networks. In *2019 56th ACM/IEEE Design Automation Conference (DAC)*, pp. 1–6, 2019.
- Seongsik Park, Seijoon Kim, Byunggook Na, and Sungroh Yoon. T2fsnn: Deep spiking neural networks with time-to-first-spike coding. *57th ACM/IEEE Design Automation Conference (DAC)*, pp. 1–6, 2020. doi: 10.1109/DAC18072.2020.9218689.
- Michael Pfeiffer and Thomas Pfeil. Deep learning with spiking neurons: Opportunities and challenges. *Frontiers in neuroscience*, 12:409662, 2018.
- Guanchao Qiao, Ning Ning, Yue Zuo, Pujun Zhou, Mingliang Sun, Shaogang Hu, Qi Yu, and Yang Liu. Spatio-temporal fusion spiking neural network for frame-based and event-based camera sensor fusion. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2024.
- Nitin Rathi, Indranil Chakraborty, Adarsh Kosta, Abhronil Sengupta, Aayush Ankit, Priyadarshini Panda, and Kaushik Roy. Exploring neuromorphic computing based on spiking neural networks: Algorithms to hardware. *ACM Computing Surveys*, 55(12):1–49, 2023.
- Hongwei Ren, Yue Zhou, Haotian FU, Yulong Huang, Xiaopeng LIN, Jie Song, and Bojun Cheng. Spikepoint: An efficient point-based spiking neural network for event cameras action recognition. In *The Twelfth International Conference on Learning Representations*, 2024.

- David P Rodgers. Improvements in multiprocessor system design. *ACM SIGARCH Computer Architecture News*, 13(3):225–231, 1985.
- Bodo Rueckauer, Iulia-Alexandra Lungu, Yuhuang Hu, Michael Pfeiffer, and Shih-Chii Liu. Conversion of continuous-valued deep networks to efficient event-driven networks for image classification. *Frontiers in Neuroscience*, 11, 2017. doi: 10.3389/fnins.2017.00682. URL <https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2017.00682>.
- Sanaullah, Shamini Koravuna, Ulrich Rückert, and Thorsten Jungeblut. Exploring spiking neural networks: a comprehensive analysis of mathematical models and applications. *Frontiers in Computational Neuroscience*, 17:1215824, 2023.
- Abhronil Sengupta, Yuting Ye, Robert Wang, Chiao Liu, and Kaushik Roy. Going deeper in spiking neural networks: Vgg and residual architectures. *Frontiers in neuroscience*, 13:95, 2019.
- Mahyar Shahsavari, Jonathan Beaumont, David Thomas, and Andrew D. Brown. POETS: A Parallel Cluster Architecture for Spiking Neural Network. *International Journal of Machine Learning and Computing*, 11(4):281–285, August 2021. ISSN 20103700. doi: 10.18178/ijmlc.2021.11.4.1048.
- Mahyar Shahsavari, David Thomas, Marcel van Gerven, Andrew Brown, and Wayne Luk. Advancements in spiking neural network communication and synchronization techniques for event-driven neuromorphic systems. *Array*, 20:100323, 2023.
- P Srivatsa, Kyle Timothy Ng Chu, Yaswanth Tavva, Jibin Wu, Malu Zhang, Haizhou Li, and Trevor E. Carlson. You only spike once: Improving energy-efficient neuromorphic inference to ann-level accuracy. *ArXiv*, abs/2006.09982, 2020. doi: 10.48550/arXiv.2006.09982.
- Jan Stuijt, Manolis Sifalakis, Amirreza Yousefzadeh, and Federico Corradi. μ brain: An event-driven and fully synthesizable architecture for spiking neural networks. *Frontiers in neuroscience*, 15: 664208, 2021.
- Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- Guangzhi Tang, Kanishkan Vadivel, Yingfu Xu, Refik Bilgic, Kevin Shidqi, Paul Detterer, Stefano Traferro, Mario Konijnenburg, Manolis Sifalakis, Gert-Jan van Schaik, and Amirreza Yousefzadeh. Seneca: building a fully digital neuromorphic processor, design trade-offs and challenges. *Frontiers in Neuroscience*, 17, 2023. ISSN 1662-453X. doi: 10.3389/fnins.2023.1187252.
- Luke Taylor, Andrew King, and Nicol S Harper. Addressing the speed-accuracy simulation trade-off for adaptive spiking neurons. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Athanasios Voulodimos, Nikolaos Doulamis, Anastasios Doulamis, and Eftychios Protopapadakis. Deep learning for computer vision: A brief review. *Computational intelligence and neuroscience*, 2018, 2018.
- Timo C Wunderlich and Christian Pehle. Event-based backpropagation can compute exact gradients for spiking neural networks. *Scientific Reports*, 11(1):12829, 2021.
- Leonid Yavits, Amir Morad, and Ran Ginosar. The effect of communication and synchronization on amdahl’s law in multicore systems. *Parallel Computing*, 40(1):1–16, 2014.
- Amirreza Yousefzadeh, Mina A Khoei, Sahar Hosseini, Priscila Holanda, Sam Leroux, Orlando Moreira, Jonathan Tapson, Bart Dhoedt, Pieter Simoens, Teresa Serrano-Gotarredona, et al. Asynchronous spiking neurons, the natural key to exploit temporal sparsity. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 9(4):668–678, 2019.

Amirreza Yousefzadeh, Gert-Jan Van Schaik, Mohammad Tahghighi, Paul Detterer, Stefano Traferro, Martijn Hijdra, Jan Stuijt, Federico Corradi, Manolis Sifalakis, and Mario Konijnenburg. Seneca: Scalable energy-efficient neuromorphic computer architecture. In *2022 IEEE 4th International Conference on Artificial Intelligence Circuits and Systems (AICAS)*, pp. 371–374. IEEE, 2022.

Semir Zeki. A massively asynchronous, parallel brain. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 370(1668):20140174, 2015.

Friedemann Zenke and Tim P Vogels. The remarkable robustness of surrogate gradient learning for instilling complex function in spiking neural networks. *Neural computation*, 33(4):899–925, 2021.

Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. Explainability for large language models: A survey. *ACM Transactions on Intelligent Systems and Technology*, 15(2):1–38, 2024.

A APPENDIX / SUPPLEMENTAL MATERIAL

A.1 SPIKING NEURAL NETWORKS

A.1.1 INCOMING CURRENT

For every neuron layer, all the synaptic weights on its inbound connections are kept in a weight matrix $\mathbf{W}^{(l)} \in \mathbb{R}^{N^{(l)} \times N^{(l-1)}}$, where $N^{(0)}$ = number of input features N_{in} . Using these weight matrices, the total incoming current x for a neuron i in the next layer can be computed using:

$$x_i^{(l+1)}[t] = \sum_{j=1}^{N^{(l)}} W_{ij}^{(l+1)} s_j^{(l)}[t] \quad (3)$$

where $s_j^{(l)}[t]$ is 1 if neuron j in layer l has emitted a spike at time t , otherwise 0.

A.1.2 MEMBRANE POTENTIAL DECAY

The decay of the membrane potential is governed by a linear Ordinary Differential Equation (ODE). The analytical solution can be used to compute the decay:

$$u[t] = u[t - \Delta t] \cdot e^{-\frac{\Delta t}{\tau_m}} + x[t] \quad (4)$$

where τ_m is the membrane time constant and Δt the elapsed time.

A.1.3 THRESHOLD FUNCTION

$$\Theta(u) = \begin{cases} 1 & \text{if } u > U_{\text{thr}} \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

where U_{thr} is the membrane potential threshold required for spiking.

A.2 EVENT-DRIVEN STATE UPDATE RULE FOR LIF NEURON

When an input current is received at time t by a neuron with a previous state $u[t_0]$ at some time $t_0 \leq t$, an atomic set of computations is executed. Start with computing the decayed membrane potential $u[t_-] = u[t_0] \cdot e^{-\frac{(t-t_0)}{\tau_m}}$, then update the membrane potential $u[t_+] = u[t_-] + x[t]$ with the input current $x[t]$, and finally set $u[t] = 0$ and emit a spike if $\Theta(u[t_+]) = 1$; otherwise $u[t] = u[t_+]$ without emitting a spike.

A.3 ARCTAN SURROGATE GRADIENT

$$\Theta(u) = \frac{1}{\pi} \arctan(\pi u \frac{\alpha}{2}) \quad (6)$$

where α is a hyperparameter modifying the steepness of the function.

A.4 CLASS PREDICTION

Class prediction involves first calculating the output rates as follows:

$$c_i = \sum_{t \in T} s[i, t] \quad (7)$$

where c_i is the spike count for class i , T are all the timesteps for which a forward pass occurred, and $s[i, t]$ is the output of the neuron representing class i in the output layer at the end of the forward pass at time t . For CIFAR-10, instead the value of c_i is equal to the membrane potential of the output neuron for class i at the end of having processed all the spikes. These values are subsequently used as logits within a softmax function:

$$p_i = \frac{e^{c_i}}{\sum_{j=1}^{N_C} e^{c_j}} \quad (8)$$

where N_C is the total number of classes. The resulting probabilities are then used to compute the cross-entropy loss:

$$L = \sum_{i=1}^{N_C} y_i \log(p_i) \quad (9)$$

where $y \in \{0, 1\}^{N_C}$ is the target class in a one-hot encoded format.

A.5 PARAMETERS FOR THE SIMULATOR

Table 4: Overview of all current simulator parameters. If the text is *italic*, then the parameter was not used for the experiments in this paper.

Name	Value range	Description
Forward group size	$\mathbb{N}_{>0}$	3.1.2
Scheduling policy	RS/MS	3.1.2
Prioritize input	True/False	If input spikes are propagated before any spikes from inside the network.
Stop condition	On spiking done / On output	3.1.2
Forward steps after output	$\mathbb{N}_{\geq 0}$	3.1.2 (only for "On output" stop condition)
Refractory dropout	$[0.0, 1.0]$	3.2.2
Momentum noise	$\mathbb{R}_{\geq 0}$	3.1.2 (only for "MS" scheduling policy)
Membrane time constant	$\mathbb{R}_{>0}$	3.1
Input spike dropout	$[0.0, 1.0]$	3.2.2
<i>Network spike dropout</i>	$[0.0, 1.0]$	The same as input spike dropout, but for spikes from inside the network, applied per forward step.
Membrane potential threshold	$\mathbb{R}_{>0}$	3.1
Timestep size	$\mathbb{N}_{>0}$	3.1.1
<i>Synchronization threshold</i>	$\mathbb{N}_{\geq 0}$	Have neurons wait for a number of input currents (including barrier messages) before being allowed to fire. Can be set per neuron.
<i>Emit barrier messages</i>	True/False	Have neurons emit a barrier message if they have retrieved exactly the number of input currents to exceed the synchronization threshold, but not enough to exceed the membrane potential threshold.

Table 5: Simulator parameters for simulating neuromorphic processors

Synchronization type	Neuromorphic processors	Forward group size	Synchronization threshold [emit barrier messages]
Barrier-based / layer	Loihi Davies et al. (2018), SENECA Yousefzadeh et al. (2022), POETS Shahsavari et al. (2021)	# neu/layer	# neu/layer [True]
Timer-based / core	SENECA Yousefzadeh et al. (2022), SpiNNaker Painkras et al. (2013), TrueNorth Akopyan et al. (2015), POETS Shahsavari et al. (2021)	# neu/core	# neu/core [True]
Asynchronous	Speck Caccavella et al. (2023), μ Brian Stuijt et al. (2021)	# neu/layer	1 [False]

A.6 DETAILS ON THE EXPERIMENTAL SETUP

A.6.1 DATASETS AND NETWORK ARCHITECTURES

The Neuromorphic MNIST (N-MNIST) dataset captures the MNIST digits using a Dynamic Vision Sensor (DVS) camera. It presents minimal temporal structure Iyer et al. (2021). It consists of 60000 training samples and 10000 test samples. Each sample spans approximately 300 ms, divided into three 100 ms camera sweeps over the same digit. Only the initial 100 ms segment of each sample is used in this study.

The Spiking Heidelberg Digits (SHD) dataset Cramer et al. (2020) is composed of auditory recordings with significant temporal structure. It consists of 8156 training samples and 2264 test samples. Each sample includes recordings of 20 spoken digits transformed into spike sequences using a cochlear model, capturing the rich dynamics of auditory processing. The 700 cochlear model output channels are downsampled to 350 channels.

The DVS gesture dataset Amir et al. (2017) focuses on different hand and arm gestures recorded by a DVS camera. Like the SHD dataset, it has significant temporal structure. It focuses on 11 different hand and arm gestures recorded by a DVS camera. It consists of 1176 training samples and 288 test samples. The 128×128 input frame is downsampled to a 32×32 frame.

CIFAR-10 Krizhevsky (2012) is a widely-used image classification dataset consisting of 60000 32×32 color images across 10 classes, including animals and vehicles. It has no temporal structure. Unlike the other datasets, all input is provided in one single timestep.

For all four datasets, input events are assigned a timestamp and an index. In the case of N-MNIST, the index corresponds to a position within a 34×34 pixel frame, with each pixel having a binary polarity value (either 1 or 0), leading to a total of $34 \times 34 \times 2 = 2312$ distinct input indices. For SHD, the index denotes one of the 350 output channels. For DVS gesture, the 128×128 input frame is downsampled to a 32×32 frame. Like N-MNIST, each pixel has a binary polarity value, so in total this gives $32 \times 32 \times 2 = 2048$ input indices. Unlike the other datasets, for CIFAR-10, the input events present a continuous value (so they are currents instead of spikes), while also still being assigned an index and timestamp (although the timestamp is irrelevant).

Table 6: Network architecture and hyperparameters. The architecture is given as [neurons in hidden layers \times number of hidden layers] - [neurons in output layer].

	N-MNIST	SHD	DVS gesture	CIFAR-10
Architecture	[64 \times 3]-10	[128 \times 3]-20	[128 \times 3]-11	see A.9
Timestep size	10 ms	10 ms	20 ms	N/A
Batch size	256	32	32	64
Epochs	50	100	70	150
Learning rate	5e-4	7e-4	1e-4	2e-4
Membrane threshold U_{thr}	0.3	0.3	0.3	0.2
Weight decay constant λ_W	1e-5	1e-4	1e-5	0
Membrane time constant τ_m	1 ms	100 ms	100 ms	N/A
Surrogate steepness α	2	10	10	2
Input spike dropout	0.25	0.2	0.2	0
Forward group size F	8	8	8	512
Refractory dropout	0.8	0.8	0.8	0.7
Momentum noise λ_{MS}	1e-6	0.1	0.1	see A.9

A.6.2 PERFORMANCE METRICS

To evaluate accuracy, output rates c_i for each class i are first calculated as outlined in equation 7. The predicted class corresponds to the one with the highest output rate. Accuracy is then quantified as the ratio of correctly predicted outputs to the total number of samples.

Spike density is computed using:

$$\text{density} = \frac{1}{N_{\text{samples}} \cdot N_{\text{neurons}}} \sum_{i=1}^{N_{\text{samples}}} N_{\text{spikes}}[i] \quad (10)$$

where N_{samples} is the number of samples, N_{neurons} is the number of neurons in the hidden layers, and $N_{\text{spikes}}[i]$ is the total number of spikes during inference of the sample i .

A.7 RESULTS ON HYPERPARAMETERS

Choosing a smaller F (i.e., with a more asynchronous system), may improve accuracy, particularly benefiting models with mechanisms that rely on network asynchrony such as momentum scheduling. However, reducing F also has its drawbacks. It significantly raises resource demands (discussed in section A.8), and there is a risk of reducing the effectiveness or even stalling the training process, as observed for SHD and DVS gesture.

Refractory dropout, can positively affect training outcomes. An explanation for this is that it increases the gradient flow by allowing more spiking activity. However, using full refractory dropout can also reduce performance, likely due to the inability to generalize to inference with refractoriness.

The momentum noise helps by introducing a slight stochastic element into the spike selection process, helping to avoid potential minima that a purely deterministic selection method is prone to get stuck in. This seems to do little for N-MNIST, but for more complex datasets like SHD and DVS gesture it has a significant effect.

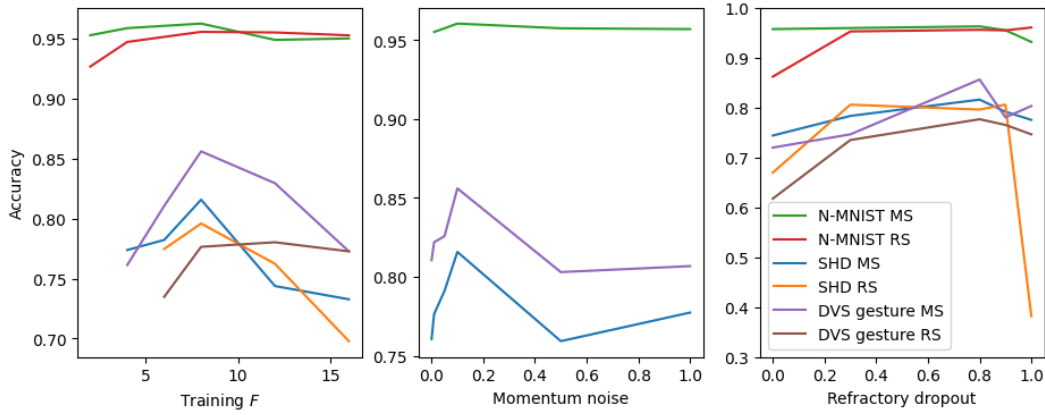


Figure 6: Results from inference with network asynchrony for different hyperparameters used during training. For SHD and DVS gesture with random scheduling, training F s smaller than 6 are not included due to the training failing to converge.

A.8 RESULTS ON RESOURCE USAGE

The increase in processing time is less pronounced for the N-MNIST and DVS gesture datasets compared to the SHD dataset. This discrepancy could be due to computational optimizations that apply specifically to the N-MNIST and DVS gesture datasets (both being vision-based datasets).

Table 7: Resource usage compared between layered and unlayered backpropagation during the second epoch of training on an NVIDIA Quadro RTX 5000.

Method	Time per epoch (s)	VRAM use (MB)
N-MNIST		
Layered	14	210
Unlayered RS, training $F = 16$	19	412
Unlayered RS, training $F = 8$	24	556
Unlayered RS, training $F = 4$	37	986
SHD		
Layered	38	214
Unlayered RS, training $F = 16$	118	2220
Unlayered RS, training $F = 8$	292	4844
Unlayered RS, training $F = 4$	681	10574
DVS gesture		
Layered	120	244
Unlayered RS, training $F = 16$	150	3802
Unlayered RS, training $F = 8$	177	6984
Unlayered RS, training $F = 4$	245	13914

A.9 DETAILS ON CIFAR-10 WITH VGG EXPERIMENTS

The details of network structure of the VGG model that we used to train on the CIFAR-10 dataset is shown in Figure 7a. We have removed (for simplicity) the batch normalization and dropout layers.

To reduce the memory overhead for training with unlayered backpropagation and make it feasible on an NVIDIA RTX 5000 we provided all the input of each sample in a single timestep (thus eliminating the unfolding across timesteps), and we adopted a forward group size $F = 512$ for training and $F = 64$ for testing. Additionally we deployed IF neurons that do not leak.

Output predictions were based on the membrane potential accumulated at the output neurons.

Because of the very large F , MS scheduling policy was also difficult to train requiring rather large amounts of annealing noise (Figure 7b).

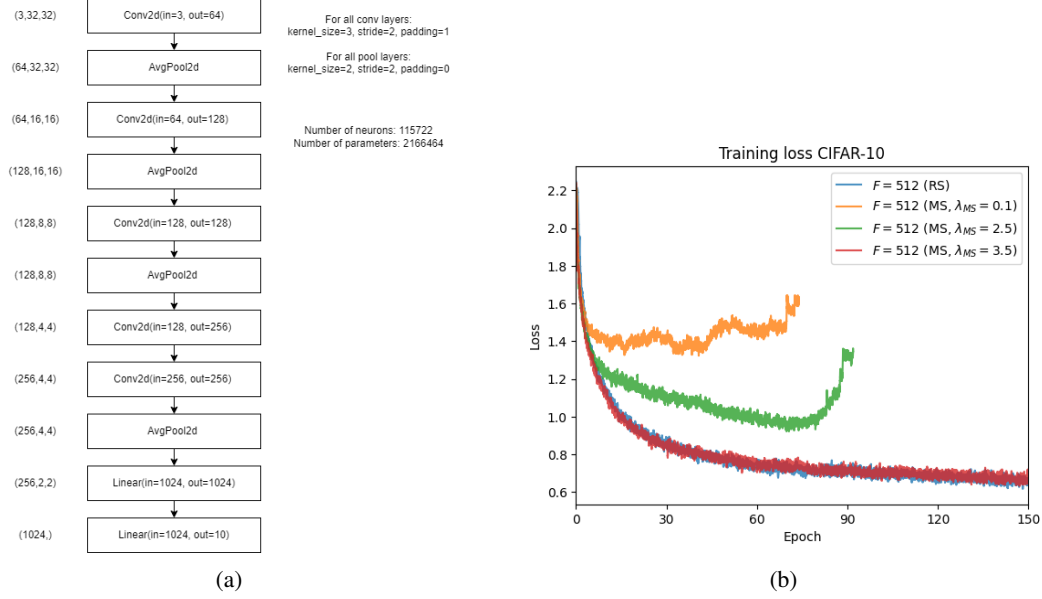


Figure 7: (a) simplified VGG network topology without batch normalization and dropout layers. (b) Training error evolution for different amounts of annealing noise in MS scheduling policy