# Fluctuation-driven initialization for spiking neural network training

Julian Rossbroich<sup>1,2,2</sup>, Julia Gygax<sup>1,2,2</sup>, and Friedemann Zenke<sup>1,2,\*</sup>

<sup>1</sup>Friedrich Miescher Institute for Biomedical Research, Basel, Switzerland <sup>2</sup>Faculty of Science, University of Basel, Switzerland

> <sup>(a)</sup>These authors contributed equally to this work. \*Corresponding author: friedemann.zenke@fmi.ch

#### Abstract

Spiking neural networks (SNNs) underlie low-power, fault-tolerant information processing in the brain and could constitute a power-efficient alternative to conventional deep neural networks when implemented on suitable neuromorphic hardware accelerators. However, instantiating SNNs that solve complex computational tasks in-silico remains a significant challenge. Surrogate gradient (SG) techniques have emerged as a standard solution for training SNNs end-to-end. Still, their success depends on synaptic weight initialization, similar to conventional artificial neural networks (ANNs). Yet, unlike in the case of ANNs, it remains elusive what constitutes a good initial state for an SNN. Here, we develop a general initialization strategy for SNNs inspired by the fluctuation-driven regime commonly observed in the brain. Specifically, we derive practical solutions for data-dependent weight initialization that ensure fluctuation-driven firing in the widely used leaky integrate-andfire (LIF) neurons. We empirically show that SNNs initialized following our strategy exhibit superior learning performance when trained with SGs. These findings generalize across several datasets and SNN architectures, including fully connected, deep convolutional, recurrent, and more biologically plausible SNNs obeying Dale's law. Thus fluctuation-driven initialization provides a practical, versatile, and easy-to-implement strategy for improving SNN training performance on diverse tasks in neuromorphic engineering and computational neuroscience.

### Introduction

Spiking neurons communicate through discrete action potentials, or spikes, thereby enabling energy efficient and reliable information processing in neurobiological and neuromorphic systems [1, 2]. Before using an SNN for any application, their connections need to be task-optimized. In conventional ANNs this step is accomplished through direct end-to-end optimization using back-propagation in combination with suitable parameter initialization [3]. However, the lack of smooth derivatives of neuronal spiking dynamics precludes using gradient-based optimization in SNNs. One increasingly common approach to overcome this issue is SG learning [4–6] which relies on continuous relaxations of the actual gradients for parameter updates. While SGs are a powerful tool for building functional SNN models, they can be adversely affected by poor initial parameter choices. In deep ANNs, suboptimal weight initialization can lead to vanishing or exploding gradients [7–9], thereby creating a major impediment to their use. Optimal weight initialization [10–12] combined with suitable architectural choices such as skip connections [11, 13] largely avoid this issue in ANNs. Similarly, the problem of vanishing gradients has been suggested to affect deep SNNs [14, 15]. However, we still lack a principled strategy for SNN initialization. Here, we close this gap by introducing a practical weight initialization strategy for SNNs. Specifically, we draw inspiration from neurobiology, where neuronal dynamics commonly exhibit fluctuation-driven firing [16, 17]. Since neurons in the fluctuation-driven regime are more sensitive to small changes in the input [18] and thus also to changes in their synaptic weights, we hypothesized that this regime could be advantageous for subsequent SG learning. In the following, we develop a general, yet simple initialization theory for SNNs consisting of LIF neurons, and empirically demonstrate its effectiveness for task-optimizing SNNs using SG techniques.

## Results

Neurons in biological SNNs commonly exhibit irregular and asynchronous firing dynamics [16, 17, 19]. Such dynamics can often be attributed to large sub-threshold fluctuations that can naturally arise through excitatory-inhibitory balance commonly observed in neurobiology [19, 20]. To test whether this fluctuation-driven regime could constitute a suitable initial state for subsequent learning, we proceeded in two steps. First, we derived a set of compact analytical expressions that link the initial synaptic weight distribution with the magnitude of sub-threshold fluctuations. Second, we numerically tested whether initializing SNNs in the fluctuation-driven regime would allow us to rapidly train these networks to high accuracy using SGs.

To arrive at analytical expressions, we note that there are primarily three factors that contribute to the membrane potential fluctuations (Fig. 1a). These are, first, the number and firing statistics of the input neurons, second, the synaptic weight distribution, and third, the postsynaptic and neuronal parameters that govern temporal integration of the inputs. For simplicity, we assume that the presynaptic input arrives from a homogeneous population of independent Poisson neurons and that the initial weight distribution is given by a Gaussian. Further, we limited our derivation to current-based LIF neurons, which are commonly used in SNN models.

To derive an expression that links the synaptic weight distribution to the fluctuation magnitude, we consider a current-based LIF neuron with membrane potential U, whose sub-threshold dynamics are given as the sum of weighted filtered presynaptic spike trains  $S_i$ :

$$U(t) = \sum_{j} w_j \epsilon * S_j(t) , \qquad (1)$$

where  $S_j = \sum_k \delta(t - t_j^k)$  denotes the output spike train of the presynaptic neuron j with firing times  $t_j^k$  and \* is a temporal convolution of the spike train  $S_j(t)$  with  $\epsilon$ , a linear filter kernel with the shape of an evoked postsynaptic potential (PSP). Specifically, we assume a synaptic model with exponentially decaying currents and, therefore, the shape of  $\epsilon$  is fully characterized by the synaptic and membrane time constants  $\tau_{\rm syn}$  and  $\tau_{\rm mem}$  (see Methods and Supplementary Material S1). Since we have many statistically independent inputs, the Central Limit Theorem guarantees that U approaches a normal distribution which is fully specified by its mean  $\mu_U$  and variance  $\sigma_U^2$ . Further assuming that presynaptic spikes are generated by homogeneous Poisson processes with associated firing rates  $\nu_j = \langle S_j \rangle$ , yields the following expressions for the mean and the variance

$$\mu_U \equiv \langle U \rangle = \sum_j w_j \nu_j \int_{-\infty}^{\infty} \epsilon(s) ds = \sum_j w_j \nu_j \bar{\epsilon}$$
<sup>(2)</sup>

$$\sigma_U^2 \equiv \langle U^2 \rangle - \mu_U^2 = \sum_j w_j^2 \nu_j \int_{-\infty}^{\infty} \epsilon(s)^2 ds = \sum_j w_j^2 \nu_j \hat{\epsilon} , \qquad (3)$$

in which  $\bar{\epsilon}$  and  $\hat{\epsilon}$  correspond to definite integrals of the filter kernel and squared filter kernel respectively which can be obtained analytically for many common neuron models (Supplementary



Figure 1. Parameterization of fluctuation-driven spiking serves as an initialization strategy for SNNs. (a) Incoming presynaptic Poisson spike trains (i) are weighted by synaptic strengths  $w_j$  and filtered through a PSP kernel  $\epsilon(t)$  (ii) to yield membrane fluctuations u(t) in a postsynaptic neuron (iii). In the fluctuation-driven regime, the membrane potential crosses the firing threshold  $\theta$  stochastically, resulting in irregular output spike trains. Because the magnitude of membrane potential fluctuations,  $\sigma_U$ , is determined by the parameters of the presynaptic weight distribution,  $\mu_W$  and  $\sigma_W$ , synaptic weights can be initialized from a target value for the fluctuation magnitude. (b) Expected and observed distributions of the membrane potential without considering spike-reset dynamics for different target fluctuation strengths expressed in terms of  $\sigma_U$  and  $\xi$ . (c) As panel (b), but considering the spike reset dynamics in the numerical simulations.

Material S1). For *n* inputs with equal firing rates  $\nu_j = \nu$  and independently drawn normally distributed weights  $W \sim \mathcal{N}(\mu_W, \sigma_W^2)$ , the above expressions further simplify to

$$\mu_U = n\mu_W \nu \bar{\epsilon} \tag{4}$$

$$\sigma_U^2 = n(\sigma_W^2 + \mu_W^2)\nu\hat{\epsilon} .$$
<sup>(5)</sup>

Finally, rewriting Equations (4) and (5) yields the desired expressions linking the synaptic weight distribution with the magnitude of the membrane potential fluctuations:

$$\mu_W = \frac{\mu_U}{n\nu\bar{\epsilon}} \tag{6}$$

$$\sigma_W^2 = \frac{\sigma_U}{n\nu\hat{\epsilon}} - \mu_W^2$$
$$= \frac{1}{n\nu\hat{\epsilon}} \left(\frac{\theta - \mu_U}{\xi}\right)^2 - \mu_W^2 . \tag{7}$$

For a neuron to be in the fluctuation-driven regime requires the bulk of the Gaussian distribution has to lie below the firing threshold (Fig. 1b). At the same time, we require a non-vanishing probability to cross the threshold to ensure some baseline levels of spiking activity. To formalize these requirements, we introduced the target parameter  $\xi$  as

$$\xi \equiv \frac{\theta - \mu_U}{\sigma_U} , \qquad (8)$$

which describes the distance between the mean membrane potential  $\mu_U$  and the spike threshold  $\theta$  in units of the standard deviation  $\sigma_U$  (Fig. 1a and Supplementary Fig. S1). To satisfy the

above requirements,  $\xi$  should be on the order of one. Concretely, we consider the range  $1 \le \xi \le 3$  (Fig. 1b). For zero-mean weight distributions, this directly translates into a desired fluctuation amplitude range  $\frac{1}{3} \le \sigma_U \le 1$ , which, given the above assumptions, is achieved by

$$\sigma_W^2 = \frac{\sigma_U^2}{n\nu\hat{\epsilon}} \tag{9}$$

at initialization.

The expressions are based on a no-spiking assumption. Hence, we expect systematic deviations from the derived membrane potential distribution in the presence of spiking. However, for small sub-threshold fluctuations ( $\sigma_U \ll 1$ ), the systematic contribution of the spike reset becomes negligible (Fig. 1c). Exact membrane potential distributions that take into consideration spike reset dynamics could be obtained using the Fokker-Planck equation [21, 22], however, such an approach does not yield compact analytic expressions and is, thus, less practical for our purposes.

Because Eqs. (4) and (5) are based on an independence assumption that is violated by realworld data, we expected further deviations in numerical simulations with real-world data. To quantify the magnitude of these deviations, we compared the predictions of Eqs. (4) and (5) with observed membrane potential fluctuations in a single LIF neuron exposed to inputs from two realistic datasets. For simplicity, we assumed a zero-mean weight distribution and used Eq. (9) to obtain its standard deviation for different target fluctuation magnitudes  $\sigma_U$ .

First, we considered a synthetic classification dataset based on random manifolds that can flexibly generate SNN benchmarks of arbitrary complexity [5] (Randman; see Methods). We generated a dataset with  $n_{\text{Randman}} = 20$  input neurons and 10 classes in which spike times belonging to the same class are drawn from a smooth random manifold (Fig. 2a) all the while different classes correspond to different manifolds. For each input pattern, each neuron fires precisely one spike during a 100 ms interval. Each 100 ms input interval was followed by 100 ms of inactivity in the input layer to allow for a propagation delay in the hidden layer (Fig. 2a). We then recorded the membrane potential distribution and found, as expected, that it deviated from a Gaussian (Fig. 2b), due to the temporal non-stationarity and structure. Next, we measured the observed membrane potential fluctuations  $\hat{\sigma}_U$  for varying target values of  $\sigma_U$  (Fig. 2c). We found that  $\hat{\sigma}_U$  was systematically smaller than  $\sigma_U$ . However, the magnitude of bias was comparable to the expected variability in the case of Poisson inputs (see Supplementary Material S2).

Next, we considered the Spiking Heidelberg Digits (SHD) speech dataset (Fig. 2d), an SNN benchmark based on real-world auditory data, which consists of approximately 10,000 spoken digits in German and English that have been converted into spikes using a biologically plausible cochlear model [23]. Importantly, SHD has a larger number of input neurons ( $n_{\text{SHD}} = 700$ ) which typically fire more than one spike with an average input firing rate of  $\nu_{\text{SHD}} = 15.8$  Hz. Again, we measured the membrane potential distribution and observed deviations from a Gaussian (Fig. 2e). In contrast to the Randman data, the observed fluctuations  $\hat{\sigma}_U$  were systematically larger than their target  $\sigma_U$  due to heavy tails in the distribution (Fig. 2f). Not surprisingly, real-world data causes systematic deviations from Eqs. (4) and (5), but these differences were on the same order as expected fluctuations due to the finite sample size of the weight and Poisson variability. Hence, we reasoned that our simple theory provides a reasonable approximation for initializing SNNs in the fluctuation-driven regime even when using real-world data.

### Initialization of shallow SNNs

We sought to evaluate whether the fluctuation-driven regime constitutes a good initialization strategy for SNN training. To this end, we trained a fully connected SNN with one hidden layer with 128 units on the Randman dataset (see Tab. 4; Methods). We initialized the weights using the parameters  $\mu_W = 0$  and  $\sigma_W$  given by our theory (Eq. (9) with target  $\sigma_U = 1$ ). This



Figure 2. Real-world datasets induce small systematic biases in fluctuation strength at initialization. (a) Two one-dimensional example manifolds from the Randman dataset, embedded into a three-dimensional space (left) and example spike raster plots corresponding to a sample from each class (right). (b) Theoretically expected distribution and numerically obtained density histogram of the membrane potential of a single neuron without spike reset in response to the Randman dataset. Because of large peaks at u(t) = 0, the x-axes in the first and middle panels have been truncated to 45% and 80% of their maximum, respectively. (c) Numerically observed  $\hat{\sigma}_U$  as a function of the target  $\sigma_U$  for the Randman dataset. The expected relationship corresponds to homogeneous and independent Poisson neurons. Shaded regions indicate standard deviation across neurons. (d) Two spike rasters that correspond to two example inputs from the SHD dataset. Input spikes are obtained by filtering recordings of spoken digits with a biologically inspired cochlear model [23]. (e) As panel (b), for the SHD dataset. X-axes in the first and middle panels have been truncated to 58% and 90% of their maximum, respectively. (f) As panel (c), for the SHD dataset.

choice resulted in asynchronous irregular firing activity consistent with the fluctuation-driven regime (Supplementary Fig. S2a-d). Subsequently, we trained the network in a supervised fashion using SGs with previously established parameters [5], back-propagation through time (BPTT), a maximum-over-time loss defined on ten readout units, and weak spiking activity regularization in the form of a soft upper bound on the population firing rate at the hidden layer (Fig. 3a; Methods). Training resulted in an SNN that accurately solved the task (test accuracy:  $97.3\% \pm 0.2$ ; train accuracy:  $99.6\% \pm 0.0$ ; Fig. 3b).

To test whether our weight initialization strategy confers an advantage over other choices of  $\mu_W$  and  $\sigma_W$ , we performed an extensive parameter search and measured validation accuracy after 200 training epochs. The network achieved the best validation accuracy when  $\mu_W$  was zero or negative and  $\sigma_W$  was close to one (Fig. 3c), well within our suggested regime of  $1 \le \xi \le 3$ . Further, we found a large parameter regime that supported learning at close-to-optimal accuracy for  $-2 \le \mu_W \le 0$  and  $\sigma_W < 10$  which extends beyond the parameter regime suggested by our theory.

To test whether these results would change on a more complex task, we trained a similar SNN on the SHD dataset [23] with weight parameters  $\mu_W = 0$  and  $\sigma_W^{(\text{SHD})} = 0.23$  as suggested by our theory (Eq. (9) with target  $\sigma_U = 1$ ). Due to differences of the number of input neurons and firing rates between the two datasets our theory predicts  $\sigma_W^{(\text{SHD})} \approx 10^{-1} \sigma_W^{(\text{Randman})}$ . After training, the network accurately classified spoken digits (test accuracy:  $65.5\% \pm 0.7$ ; train accuracy:  $100.0\% \pm 0.0$ ; Fig. 3d). As before, we performed an extensive parameter search over different



Figure 3. Initialization in the fluctuation-driven regime results in optimal learning performance. (a) Top: Schematic of the SNN used for training. Bottom: Illustration of the learning dynamics. The supervised loss function  $\mathcal{L}_{sup}$  relies on the maximum membrane potential over time of readout units  $U_i^{(\text{out})}[t]$ , to which a Softmax and cross-entropy loss  $\mathcal{L}_{\text{CE}}$  is applied. All networks were trained by minimizing  $\mathcal{L}_{sup}$  in the direction of negative SGs, computed with BPTT. (b) Snapshot of network activity over time after training on the Randman dataset. Bottom: Spike raster of input layer activity from two different samples corresponding to two different classes is shown. Middle: Spike raster of hidden layer activity. Top: Membrane potential of readout units. The readout units corresponding to the two input classes are highlighted in different shades. (c) Heatmap showing validation accuracy after training on the Randman dataset as a function of the parameters of the synaptic weight distribution at initialization. (d) Same as in panel (b), but for a network trained on the SHD dataset. (e) Same as panel (c), but for the SHD dataset. (f) Validation accuracy as a function of target fluctuation magnitude  $\sigma_U$  for initializations in the balanced state with  $\mu_U = \mu_W = 0$ . The shaded region around the lines indicates the range of values across five random seeds. The sand-colored shaded region corresponds to our suggested target fluctuation magnitude  $\frac{1}{3} \leq \sigma_U \leq 1$ . (g) Average hidden layer firing rate as a function of  $\sigma_U$ . (h) Average magnitude of SGs with respect to the output in the readout layer (top) and hidden layer (bottom) as a function of  $\sigma_U$ . (i) Same as panel (h), but for the average magnitude of SGs with respect to the synaptic weights.

initializations and found that networks initialized in the fluctuation-driven regime  $(1 \le \xi \le 3)$ showed close-to-optimal performance (Fig. 3e, f). Unlike in the Randman case, the parameter regime with good performance was much smaller and tightly constrained around  $\mu_W \approx 0$ . Finally, even though our initialization strategy posits that neurons be in the fluctuation-driven regime, we observed a sizeable fraction of hidden layer neurons with regular firing activity both before (Supplementary Fig. S2e) and after learning (Fig. 3d). We found that our theory predicts these cases (Supplementary Fig. S2d, h) due to the inherent variability in the sampling of synaptic weights (Supplementary Material S2 and Supplementary Fig. S3).

For both datasets, we found that initialization with  $\sigma_W \ll 1$  and  $\mu_W \approx 0$  supported closeto-optimal learning. This result surprised us because the ensuing vanishing membrane potential fluctuations should lead to quiescent hidden layer activity. To check whether this is indeed the case, we initialized networks with different target values for  $\sigma_U$  and recorded their hidden layer activity. As expected, we found that fluctuation magnitudes  $\sigma_U \ll 1$  still supported close-tooptimal learning performance (Fig. 3f), despite an absence of spikes in the hidden layer at the time of initialization (Fig. 3g).

Because vanishing spiking activity should influence gradient magnitudes during backpropagation, we recorded the magnitude of the SG with respect to the output at the readout and hidden layers at the time of the first training epoch. Due to the nature of the loss function, initialization does not affect the magnitude of the gradient in the readout layer but can change the magnitude of the gradient by two orders of magnitude in the hidden layer (Fig. 3h). Consequently, the absolute magnitude of weight changes is also amplified in the hidden layer when fluctuations are large (Fig. 3i). Since the synaptic weight update depends on presynaptic activity, initializations resulting in quiescent hidden layers (Fig. 3g) lead to an absence of weight updates in the readout layer (Fig. 3i). However, as long as SGs do not vanish in the first layer, the network can recover spike propagation and therefore gradient flow during training. That the network is able to learn without problems in this regime may seem surprising at first and is indeed a peculiarity of SGs.

In addition to classification accuracy, the sparsity of neuronal activity is a key SNN performance indicator. To limit firing rates in the hidden layers to a sensible regime, we optimized all networks with activity regularization. Specifically, we added a soft upper bound on the population firing rate (see Methods). This regularization punishes population firing rates in the hidden layers exceeding 10 Hz. To investigate the effect of weight initialization on sparsity, we systematically recorded population firing rates of the above network trained with or without activity regularization. As expected, activity regularization resulted in average population firing rates of  $< 10 \,\text{Hz}$  following training, independent of the target  $\sigma_U$  at initialization. In contrast, networks trained without activity regularization exhibited population firing rates exceeding 60 Hz in the hidden layer and only weak dependence on the target  $\sigma_U$  at initialization (Supplementary Fig. S4a, b). Next, we wanted to ensure that activity regularization does not result in a substantial loss in classification accuracy. To that end, we compared the accuracy of networks trained with and without activity regularization for the given threshold and strength parameters. We found that regularized networks performed only slightly worse than their unregularized counterparts albeit with vastly reduced average firing rates (Supplementary Fig. S4c). Based on these findings, we used activity regularization on the population firing rates in all subsequent experiments.

Thus far, we studied strictly feed-forward SNNs without recurrent hidden layer connections. Recurrent SNNs typically perform better than feed-forward networks on tasks requiring memory such as SHD [5]. To that end, we extended our initialization strategy to networks with recurrent connections (see Methods) and applied it to recurrent SNNs with one hidden layer. As in the case of feed-forward networks, we found recurrent SNNs trained well with sufficiently small target fluctuations  $\sigma_U$  (Supplementary Fig. S5a, b). In summary, shallow SNNs are surprisingly robust to initialization when the absolute magnitude of the weights is small. In practice, initialization with  $\mu_U = 0$  and a target fluctuation magnitude  $\sigma_U \leq 1$  can be used to achieve close-to-optimal learning performance.

### Initialization of deep SNNs

We hypothesized that deep SNNs are more sensitive to initialization, as is the case with deep ANNs [11]. To test this hypothesis, we first extended our initialization strategy to deep and recurrent convolutional spiking neural network (CSNN) architectures (see Methods). We then initialized several CSNN architectures with different numbers of recurrently connected hidden layers according to Eqs. (37-39) with target  $\mu_U = 0$  and different targets  $\sigma_U$ . Subsequently, we trained the resulting networks and measured validation accuracy on held-out data. As expected, sensitivity to the fluctuation magnitude at initialization increased with network depth (Fig. 4a). As in shallow fully connected networks, CSNNs with a single hidden layer were remarkably robust to initialization and close-to-optimal training performance was achieved for  $\sigma_U \leq 1$ . In contrast, deep networks with three hidden layers performed well when the fluctuation magnitude fell into the regime  $0.05 \le \sigma_U \le 3$ . This regime was narrowed further in deeper networks with seven hidden layers, which only achieved high validation accuracy for initializations in the range  $0.5 \leq \sigma_U \leq 2$ . Like in the shallow case, activity regularization ensured sparse activity with a negligible effect on classification accuracy (Supplementary Fig. S4d, e). Finally, although a seven-layer CSNN did not improve classification performance on this task over the three-layer network, we wanted to know whether initialization with  $\sigma_U = 1$  would be conducive for training even deeper networks. To get at this question, we extended our network to ten hidden layers, the deepest possible architecture afforded by our GPU memory while keeping all other training parameters equal to networks with seven layers, and found that initialization with  $\sigma_U = 1$ resulted in reliable training (test accuracy:  $81.1\% \pm 1.6$ ; validation accuracy:  $93.6\% \pm 2.9$ ). Crucially, when we instead trained with Kaiming initialization [11], the standard initialization for non-spiking rectified linear unit (ReLU) networks, learning failed in CSNNs with seven or more hidden layers. In summary, we observed that fluctuation-driven initialization with  $\sigma_U = 1$ supports learning in deep CSNNs.

To check whether depth increases the generalization performance of trained networks, we compared the test error of successfully trained CSNNs with one, three, seven, and ten hidden layers. We found that deeper networks did not show better generalization performance than one-layer networks (Fig. 4b). These findings suggest that the addition of multiple hidden layers does not provide an advantage in recurrently connected CSNNs on the SHD dataset. Since recurrently connected networks can be considered as deep in time, we were wondering whether strictly feed-forward SNNs would benefit from increasing depth. To that end, we repeated training of deep CSNNs with corresponding layer sizes but without recurrent hidden layer connections on the SHD dataset (see Methods). Indeed, we found that deep feed-forward SNNs performed better than shallow feed-forward SNNs (Fig. 4b and Supplementary Fig. S6).

In addition to the classification accuracy, weight initialization affects training speed. To check whether fluctuation-driven initialization is conducive to fast training, we measured the number of required epochs to reach 90% accuracy on the training dataset in CSNNs with one and three hidden layers. We found that networks initialized in the fluctuation-driven regime ( $\sigma_U = 1$ ) trained fastest (Fig. 4c and Supplementary Fig. S7). On average, CSNNs with one hidden layer initialized with target  $\sigma_U = 1$  reached 90% training accuracy after 19.2 epochs, and CSNNs with three hidden layers required 16.8 epochs to reach 90% training accuracy. Thus, fluctuation-driven initialization is conducive to fast training.

Vanishing SGs impair learning in deep SNNs. In deep ANNs initialization is closely related to the problem of vanishing or exploding gradients [7–9]. We wondered whether this mechanism, i.e., vanishing or exploding SGs, prevented training in deep SNNs when  $\sigma_U$  falls outside the optimal regime. To test this idea, we initialized seven-layer CSNNs with different

![](_page_8_Figure_1.jpeg)

Figure 4. Deep CSNNs are sensitive to initialization due to vanishing SGs. (a) Validation accuracy as a function of target fluctuation strength  $\sigma_U$  for recurrent CSNNs of increasing depth. All networks were trained on the SHD dataset. The triangular markers in the right plot correspond to the values of  $\sigma_U$  plotted in panels (d)-(f). The shaded region around the lines indicates the range of values across five random seeds. The sand-colored shaded region corresponds to our suggested target fluctuation magnitude  $\frac{1}{3} \leq \sigma_U \leq 1$ . The dashed line corresponds to Kaiming initialization. (b) Test error of the five best-performing models in terms of validation accuracy, for different numbers of hidden layers and for networks with and without recurrent connections in the hidden layers. \*No initialization parameter sweeps were performed for networks with ten hidden layers. Instead, the data depict results obtained from five networks initialized with target  $\sigma_U = 1$ . (c) Training speed of CSNNs, as illustrated by the number of required epochs to reach 90% training accuracy on the SHD dataset. (d) Population firing rate at initialization (before training) as a function of hidden layers in a CSNN with seven hidden layers, for different values of  $\sigma_U$ . (e) As panel (d), but displaying the magnitude of SGs. (f) As panel (d), but displaying the magnitude of the synaptic weight update. When membrane potential fluctuations are so small that neurons in the previous layer do not spike, the weight update equals zero.

targets  $\sigma_U$  and recorded the neuronal activity in hidden layers. Like in shallow SNNs (Fig. 3g), initializations with small  $\sigma_U$  led to quiescent hidden layers in deep CSNNs, which impaired the activity propagation to deeper layers (Fig. 4d). Specifically, in networks initialized with  $\sigma_U = 0.5$ , only the first four hidden layers exhibited spiking activity. This effect was amplified in networks initialized with  $\sigma_U = 0.2$ , in which all but the first hidden layer were quiescent. In contrast, networks initialized with  $\sigma_U = 2$  exhibited a strong increase in firing rates in deeper layers, and initializations with  $\sigma_U = 1$  led to stable activity propagation with a firing rate of  $\approx 10$  Hz throughout the network.

We next investigated how impaired activity propagation influenced SG magnitudes. To that end, we recorded SG magnitudes at each hidden layer during training. In networks initialized with  $\sigma_U = 0.5$  and  $\sigma_U = 0.2$ , in which spiking activity vanished in deep layers, each quiescent layer decreased SGs by approximately two orders of magnitude (Fig. 4e). As a result, the magnitude of weight updates in early layers decreased by several orders of magnitude consistent with the numerical value of the surrogate derivative for neurons at rest (0.023 for  $\beta = 20$ ; see Methods). Moreover, weight updates vanished in deeper layers, caused by the lack of presynaptic activity (Fig. 4f). In contrast, initializations with  $\sigma_U \ge 1$  led to relatively stable SG and weight update magnitudes across all layers (Fig. 4e, f). Notably, gradients were consistently one to two orders of magnitude smaller in networks initialized with  $\sigma_U = 20$  compared to networks initialized with  $\sigma_U = 1$  or  $\sigma_U = 2$  (Fig. 4e, f).

In summary, the sensitivity to initialization in deep SNNs is caused by impaired activity propagation to deeper layers and associated vanishing SGs. Empirically we found that only initializations with  $\sigma_U \approx 1$  supported both propagation of sparse population activity and stable magnitudes of back-propagating SGs in deep networks.

Since the surrogate derivative used to compute SGs is to some extent freely tunable [5], one might argue that re-scaling it could provide a potential solution to vanishing SGs by ensuring stable gradient magnitudes during back-propagation (see Methods). We tested this approach and found that a re-scaled SG can only prevent vanishing gradients in the absence of spiking at the cost of exploding gradients when the network does exhibit spiking which emerges over training (Supplementary Fig. S8a-c). In strictly feed-forward networks, we found that the gradients were less prone to exploding, hence re-scaling the SG could potentially alleviate the problem of vanishing gradients (Supplementary Fig. S8d-f) and therefore increase robustness to initialization. However, with increasing depth, exploding gradients would likely prevent successful training even in deep feed-forward SNNs.

Seeing that training of deep SNNs was sensitive to the magnitude of SGs [24], we speculated that the robustness to weight initialization we observed in three-layer CSNNs could be attributed to the use of our optimizer with a per-parameter learning rate during training (see Methods). To test this idea, we trained three-layer CSNNs initialized with different  $\sigma_U$  either with a smart optimizer [25, 26] or with stochastic gradient descent (SGD) without an optimizer. We found that networks trained with SGD were indeed more sensitive to the fluctuation magnitude at initialization (Supplementary Fig. S9). This effect was especially prominent in recurrent CSNNs.

Homeostatic plasticity increases robustness to initialization in deep SNNs. Because quiescent hidden layers are closely linked to vanishing SGs and thus to preventing training in deep SNNs, a homeostatically maintained firing rate, as observed in biological neural networks [27–29], could rescue activity propagation and therefore enable training. To test this hypothesis, we implemented homeostatic plasticity as an additional regularization term in the loss function that sets a lower bound on the firing rate of each individual neuron [23], which penalizes quiescent neurons (Fig. 5a; see Methods). We trained three-layer recurrent CSNNs on the SHD dataset, either with or without the additional homeostatic plasticity term. Indeed, homeostatic plasticity rescued training performance for networks initialized with  $\sigma_U \ll 1$  (Fig. 5b).

Next, we investigated whether homeostatic plasticity was necessary throughout the whole training period, or whether rescuing activity propagation before supervised training would be sufficient to enable learning. To this end, we developed a form of dynamic initialization for SNNs involving a homeostatic priming period before training. During the initial priming period, initialized networks were optimized solely on the homeostatic objective to nudge the spiking activity into a regime conducive to learning. After priming, we removed the homeostatic objective and started the supervised training period as usual. Like ongoing homeostatic plasticity, the homeostatic priming period was capable of rescuing learning for initializations with  $\sigma_U \ll 1$  (Fig. 5c). However, in rare cases, the network did not train after successful priming and the restored spiking activity vanished during training on the supervised loss function.

We wondered whether homeostatic plasticity affected the network's generalization performance and thus compared the test error of networks trained with the proposed homeostatic mechanisms. Neither ongoing homeostatic plasticity nor homeostatic priming had a systematic effect on the test error (Fig. 5d). Therefore, we concluded that both biologically inspired homeostatic plasticity and homeostatic priming are effective strategies to increase the robustness towards initialization in deep SNNs without impairing their performance.

![](_page_10_Figure_1.jpeg)

Figure 5. Homeostatic plasticity increases the robustness to initialization in deep SNNs. (a) Illustration of the homeostatic activity mechanism as a firing rate regularizer. Homeostatic plasticity (green) prevents neurons from remaining silent by increasing the synaptic weights when the firing rate is low. In all our simulations, a complementary upper bound activity regularizer (grey), that acts on the population-level, prevents neurons from spiking incessantly. (b) Validation accuracy after training a deep convolutional SNN with three hidden layers on the SHD dataset as a function of  $\sigma_U$ . The colored line corresponds to networks trained with an active homeostatic plasticity mechanism. The black line corresponds to the baseline without homeostatic plasticity. The shaded region around the lines indicates the range of values across five random seeds. (c) As panel (b), for networks that were primed for 10 epochs with a homeostatic plasticity mechanism prior to supervised learning. During supervised learning, the homeostatic mechanism was inactive. (d) Test error of the 5 best-performing models in terms of validation accuracy for models trained with homeostatic plasticity, homeostatic priming, and the baseline model.

Deep SNNs with skip connections are more robust to initialization. In deep ANNs, skip connections are standard practice to facilitate optimization and improve training performance [13, 30, 31]. For instance, residual networks (ResNets) [31] use residual connections, a specific type of identity skip connections whereby the inputs are added directly to the output of a layer or block. We argued that residual connections are ill-defined in SNNs as the spiking non-linearity would only allow adding spikes to the input spike train. Instead, we considered classic skip connections and asked whether they rescue spike propagation in deep CSNNs. We tested this idea in CSNNs with three hidden layers by implementing trainable skip connections between each hidden layer and the readout layer (Supplementary Fig. S10a; Methods). Skip connections indeed increased robustness to initialization, with respect to both large  $\sigma_U > 10$  and small  $\sigma_U \ll 1$  (Supplementary Fig. S10b). However, generalization performance after training did not increase as a result of added skip connections (Supplementary Fig. S10d). Notably, for initializations with small  $\sigma_U \ll 1$ , optimized networks only propagated activity through the skip connection between the first hidden layer and the readout layer, effectively reducing the network to a single hidden layer. As skip connections did not prevent all layers from being quiescent in deep SNNs, we wondered whether homeostatic plasticity and skip connections complement each other and further increase performance for initializations with  $\sigma_U \ll 1$ . Thus, we trained three-layer CSNNs with skip connections and ongoing homeostatic plasticity. Networks with combined skip connections and homeostatic plasticity also exhibited enhanced robustness to initialization but did not show a significantly better generalization performance (Supplementary Fig. S10c, d). We concluded that skip connections are a viable approach to increase the robustness towards initializations with large  $\sigma_U > 10$  in deep CSNNs, but are not able to compensate for vanishing gradients in deep layers when  $\sigma_U \ll 1$ .

Fluctuation-driven initialization performs robustly across datasets. Together, our results suggest that traditional Kaiming initialization used for ANNs is sufficient for training three-layer CSNNs, but breaks down when training seven-layer or deeper CSNNs on the SHD

dataset. In contrast, our proposed initialization strategy with the target fluctuation parameter set to  $\sigma_U = 1$  yields close-to-optimal training performance in all three-, seven-, and even tenlayer networks. To directly compare fluctuation-driven and Kaiming initialization, we measured generalization performance in terms of test accuracy after training. As expected, we found only small differences in test accuracy for three-layer networks (Fig. 6a; Tab. 1).

Specifically, Kaiming initialized three-layer networks achieved an average test accuracy of  $83.1\% \pm 1.2$  (validation accuracy:  $95.9\% \pm 1.6$ ), while the same networks initialized with our proposed strategy reached an average test accuracy of  $82.7\% \pm 1.1$  (validation accuracy:  $94.1\% \pm 1.7$ ). Seven-layer networks initialized with Kaiming initialization performed close to chance level after training (test accuracy:  $4.5\% \pm 0.0$ ; validation accuracy:  $4.7\% \pm 0.5$ ), while networks initialized with  $\sigma_U = 1$  reached  $83.5\% \pm 1.3$  accuracy on the test set (Fig. 6b; Tab. 1; validation accuracy:  $94.9\% \pm 1.0$ ). As homeostatic plasticity was able to compensate for suboptimal initializations by rescuing activity propagation in three-layer CSNNs, we wondered whether these results extend to seven-layer networks. To this end, we trained Kaiming-initialized seven-layer CSNNs with ongoing homeostatic plasticity. Indeed, homeostatic plasticity rescued training, but test accuracy after training (test accuracy:  $77.0\% \pm 3.0$ ; validation accuracy:  $95.0\% \pm 1.8$ ) was worse compared to networks initialized with  $\sigma_U = 1$  that were trained without homeostatic plasticity (Fig. 6b; Tab. 1).

So far, we have limited our investigation to initialization-dependence to deep CSNNs trained on the SHD dataset, which is relatively small and may thus be prone to overfitting. To test whether our findings would generalize to other tasks, we trained deep feed-forward CSNNs on two additional datasets from different input modalities. First, we considered CIFAR-10, a dataset consisting of static images. To translate static image input into spiking, we augmented the networks with an additional layer of simulated sensory neurons into which we injected the individual image pixel values as currents. Both bias currents and current gain were optimized end-to-end with all other network parameters (see Methods). We then constructed deep CSNNs with increasing numbers of hidden layers (see Tab. 6; Methods). As before, networks were either initialized with traditional Kaiming initialization or with a target membrane potential fluctuation magnitude of  $\sigma_U = 1$ . We observed that networks with up to two hidden layers showed good training performance with both initializations (Fig. 6c; Tab. 1). When we increased the number of hidden layers to four, networks initialized with  $\sigma_U = 1$  continued to show good training performance, while networks initialized with Kaiming initialization failed to train (Fig. 6d; Tab. 1). Training on CIFAR-10 with ongoing homeostatic plasticity was able to rescue learning in Kaiming initialized SNNs with four hidden layers.

Since CIFAR-10 is a still image dataset, which lacks temporal dynamics, it is less well suited for assessing SNNs performance. To check whether our results generalize to other commonly used SNN datasets, we considered the DVS-Gesture dataset [32], which consists of short videos

	SHD		CIFAR-10		DVS-Gesture	
	$n_{\rm H} = 3$	$n_{\rm H} = 7$	$n_{\rm H} = 2$	$n_{\rm H} = 4$	$n_{\rm H} = 6$	$n_{\rm H} = 8$
Kaiming	$83.1 \pm 1.2$	$4.5\pm0.0$	$59.5\pm0.8$	$10.0\pm0.0$	$54.6\pm37.1$	$9.1 \pm 0.0$
Kaiming & Hom.	-	$77.0\pm2.9$	-	$70.3 \pm 0.9$	-	$82.3\pm5.3$
Fluctdriven	$82.7 \pm 1.1$	$83.5 \pm 1.3$	$62.4 \pm 0.3$	$65.6 \pm 1.3$	$86.7 \pm 1.2$	$86.4 \pm 1.7$

**Table 1.** Test accuracy in percent after training networks with different numbers of hidden layers and different initializations (Kaiming, Kaiming with homeostatic plasticity and fluctuationdriven initialization with  $\sigma_U = 1$ ) on the SHD, CIFAR-10, and DVS-Gesture datasets. Errors correspond to the standard deviation.

![](_page_12_Figure_1.jpeg)

Figure 6. Fluctuation-driven initialization enables training of deep SNNs across multiple datasets. (a) Test accuracy of three-layer CSNNs trained on the SHD dataset. Networks were initialized either with standard Kaiming initialization (Kaiming) or fluctuation-driven initialization with  $\sigma_{U} = 1.$ All error bars indicate standard deviation across five runs. (b) Test accuracy of seven-layer CSNNs trained on the SHD dataset. Networks with Kaiming initialization were additionally trained with ongoing homeostatic plasticity (Kaiming & Hom. plast.) (c) Same as panel (a), but for two-layer feed-forward CSNNs trained on the CIFAR-10 dataset. (d) Same as panel (b), but for four-layer feed-forward CSNNs trained on the CIFAR-10 dataset. (e) Same as panel (a), but for six-layer feed-forward CSNNs trained on the DVS-Gestures dataset. (f) Same as panel (b), but for eight-layer feed-forward CSNNs trained on the DVS-Gestures dataset.

that depict humans performing different hand gestures. These videos were recorded using an event camera, yielding event-based outputs that can be used to train SNNs on the classification of the performed gestures. As before, we initialized deep CSNNs with an increasing number of hidden layers using either Kaiming initialization or a target  $\sigma_U = 1$  and compared their test accuracy after training (see Tab. 6; Methods). We found that networks with up to six hidden layers could be successfully trained using either Kaiming or our proposed initialization (Fig. 6e; Tab. 1). However, in six-layer networks, initialization with  $\sigma_U = 1$  yielded more reliable training performance and higher accuracy than Kaiming initialization. When we increased the number of hidden layers to eight, networks initialized with Kaiming initialization did not train successfully, while networks initialized with a target  $\sigma_U = 1$  continued to show good learning performance (Fig. 6f; Tab. 1). As already observed on the SHD and CIFAR-10 datasets, training of Kaiming initialized deep networks could be rescued by adding homeostatic plasticity during training.

Taken together, these findings paint a clear pattern of initialization dependencies across datasets: Up to a certain number of hidden layers, which is dataset dependent, Kaiming initialization yields good training performance in SNNs. However, when networks become too deep, vanishing SGs prevent training in networks with Kaiming initialization. In contrast, our proposed initialization strategy enables learning at high performance for deeper networks when the target fluctuation magnitude is set to  $\sigma_U = 1$ . As a complementary data-dependent strategy, homeostatic plasticity can be used to prevent vanishing gradients and rescue learning in deep networks that were initialized in a suboptimal regime.

#### Initializing SNNs that obey Dale's law

Neurons in biological SNNs are separated into excitatory and inhibitory populations, a constraint commonly known as Dale's Law [33]. With added biological constraints, functional SNNs constitute an important in-silico model system for computational neuroscience. To advance the development of biologically constrained SNNs, we extended our initialization theory to SNNs obeying Dale's law (see Methods), i.e., in which each hidden layer consists of recurrently connected but separate excitatory and inhibitory populations (Fig. 7a). At initialization, we require a balance between excitatory and inhibitory currents ( $\mu_U = 0$ ), as is commonly observed in biology [34, 35]. To accomplish such balance, we assume that excitatory and inhibitory synaptic weights are drawn from independent exponential distributions, whose mean values are set according to our theory to ensure the desired membrane potential dynamics (Supplementary Fig. S11). This strategy allowed us to initialize Dalian networks with the same target  $\sigma_U$  as non-Dalian networks.

To test whether Dalian networks in the fluctuation-driven regime could be trained to high accuracy like their non-Dalian counterparts, we first considered fully connected recurrent Dalian SNNs with one hidden layer trained on the SHD dataset (see Methods). Dalian networks initialized with  $\sigma_U = 1$  accurately solved the SHD task after training for 200 epochs (99.8%±0.0 train & 82.2%±1.2 test accuracy; Fig. 7b). Next, to test the robustness to initialization in Dalian networks, we initialized Dalian SNNs with different targets  $\sigma_U$  and trained them on the SHD dataset. For direct comparison between Dalian SNNs and non-Dalian SNNs, we constructed SNNs with a total of  $n_{\rm h} = 160$  hidden layer neurons, which were further split into  $n_{\rm exc} = 128$  and  $n_{\rm inh} = 32$  neurons for the Dalian case (see Tab. 4; Methods). After training, the Dalian networks exhibited similar robustness to initialization as non-Dalian networks (Fig. 7c). While we did not observe a large difference between Dalian and non-Dalian networks in validation accuracy,

![](_page_13_Figure_3.jpeg)

Initialization of Dalian SNNs in the fluctuation-driven regime. Figure 7. (a) Schematic of a shallow SNN obeying Dale's law. Excitatory (red) and inhibitory (blue) populations are recurrently connected, but separate. (b) Snapshot of network activity over time after training a shallow SNN obeying Dale's law on the SHD dataset. Bottom: Spike raster of input layer activity from two samples corresponding to two different classes. Middle: Spike raster of excitatory (red) and inhibitory (blue) activity in the hidden layer. Top: Membrane potential of readout units. The readout units corresponding to the two input classes are highlighted in different shades. (c) Performance comparison of Dalian and Non-Dalian shallow SNNs. Left: Validation accuracy after training on the SHD dataset as a function of initialization target  $\sigma_{U}$ . The shaded region around the lines indicates the range of values across five random seeds. The sand-colored shaded region corresponds to our suggested target fluctuation magnitude  $1 \le \xi \le 3$ . Right: Test error of the five best-performing models in terms of validation accuracy, for Dalian and Non-Dalian SNNs. Error bars mark  $\pm$  one standard deviation. (d) As panel (c), for Dalian and Non-Dalian three-layer CSNNs.

Dalian networks exhibited higher accuracy on the SHD test dataset. This result suggests that the separation into excitatory and inhibitory populations could provide a functionally beneficial constraint for SNNs with one recurrently connected hidden layer trained on the SHD dataset.

We wondered whether the better generalization performance of shallow Dalian SNNs would extend to deeper CSNN network architectures. To address this question, we constructed Dalian CSNNs with three hidden layers (see Methods). Again, networks were initialized with different targets  $\sigma_U$  and trained on the SHD dataset. We found that Dalian CSNNs with three hidden layers were more sensitive to initialization than their non-Dalian counterparts (Fig. 7d). However, when successfully trained, Dalian and Non-Dalian CSNNs resulted in similar test accuracies.

In summary, our initialization strategy extends to Dalian SNNs with different network architectures and enables robust training on the SHD dataset. Unexpectedly, constraining networks with Dale's law increased generalization accuracy by 7.1% in shallow networks. However, this effect did not generalize to deep CSNNs. Thus initializing Dalian networks in the fluctuationdriven regime is beneficial for their training and it will be interesting future work to study whether and how these findings generalize to larger datasets.

# Discussion

We have introduced a general and easy-to-implement initialization strategy for SNNs and shown that it yields close-to-optimal training speed and classification performance across different SNN architectures and datasets. To that end, we developed a simple and general theory based on the notion of fluctuation-driven firing and tested it empirically in numerical simulations. We found that shallow SNN architectures are surprisingly robust to initialization with small synaptic weight magnitudes, whereas deep CSNNs require carefully chosen initial weight distributions that our theory accurately predicts. Further, our analysis showed that suboptimal initial weight choices result in vanishing or exploding SGs, similar to ANNs. Importantly, for all network architectures, including deep convolutional, recurrent, and Dalian SNNs, and the different datasets we considered, we found that fluctuation-driven initialization with given target membrane fluctuations of  $\sigma_{II} = 1$ , resulted in stable activity propagation and close-to-optimal learning performance. Based on our results, we recommend initializing SNNs in the fluctuationdriven regime using a target  $\sigma_U = 1$  for all practical purposes. If activity propagation remains limited after training, a problem we observed in deeper network architectures, we recommend the addition of firing rate homeostasis to the training loss either for the entire training process or transiently during an initial priming period.

Functional SNNs are most commonly obtained by converting a previously trained ANN [36–40] or through direct training using timing-based methods [41–45] or SGs [5, 6, 46]. While both approaches can result in well-performing networks, direct training typically leads to sparser activity levels while also leveraging spike timing which can be beneficial for energy efficiency [47]. The initialization strategy developed in this article mainly applies to direct training approaches and specifically for SNNs trained with SGs.

Most previous SNN studies relied on weight initialization strategies that were established for ANNs in which they aim at keeping the variance of gradients constant through time or layers. For example, Xavier (Glorot) initialization [10] achieves stable variance in the backward pass by appropriately scaling the initial weight distribution. While the Xavier initialization was originally developed for linear networks, the Kaiming (He) initialization [11] extends this approach by explicitly taking into account the ReLU non-linearity, thereby enabling the training of deeper ReLU networks. While both Xavier and Kaiming initialization posit a scaling of weights by the number of input neurons as  $\sim 1/n_{\rm in}$ , their profound effects on learning performance largely result from the different choice of the absolute weight scale, which differs by a factor of two, a direct consequence of neuronal non-linearity.

Alternatively, the weight scale in the case of SNNs is often determined empirically, however, there are some proposed initialization strategies, although they often lack a sound theoretical foundation. For example, Lee et al. [14] proposed to normalize the magnitude of backpropagated errors across layers by initializing synaptic weights from a uniform distribution  $W_{(l)} \sim \mathcal{U}[-\sqrt{3/n_l}, \sqrt{3/n_l}]$ , where  $n_l$  is the number of incoming synapses. However, this approach was only validated in networks with two hidden layers trained on an event-based version of the MNIST dataset [48] and requires manual tuning of a per-layer weight scale to define the spiking threshold. Bellec et al. [49] in turn initialized weights for spiking LSTM models from a Normal distribution as  $W_{(l)} \sim \mathcal{N}(0, 1/n_{l-1})$ , whereas Zenke et al. [5] used a uniform distribution  $W_{(l)} \sim \mathcal{U}[-\sqrt{1/n_{l-1}}, \sqrt{1/n_{l-1}}]$ . A more intricate approach was developed by Herranz-Celotti et al. [50], who suggested several conditions on the initial weights that aim, e. g., to balance the variance of the gradients across time and layers. Based on those conditions, the authors derived a way to determine the weight scale for a Uniform distribution. While initialization with an ad-hoc chosen weight scale can support successful training in shallow networks, none of these studies applied their initialization strategies to network architectures with more than two hidden layers. However, as shown in this article, shallow network architectures are intrinsically robust to initialization as long as the weights are small enough while the need for SNN-specific initialization mainly arises when training deep SNNs. It thus remains an open question whether these results generalize to deep SNNs.

Recently, Ding et al. [51] proposed an initialization strategy that generalized to deep SNN architectures. The authors related the magnitude of backpropagating gradients in feed-forward SNNs to the synaptic weight distribution and proposed a weight scale for normally distributed synaptic weights that takes into account some parameters of neuronal dynamics, but does not consider dataset-dependent input parameters. While similar to the approach outlined here, this initialization strategy is limited to centered weight distributions and feed-forward networks. In addition, this particular study limited the forward pass to 20 time steps and delta synapses, compared to 100-500 time steps and current-based synapses in our simulations. Using delta synapses and a smaller number of time steps can increase the performance of SNNs but does so at the cost of biologically realistic membrane potential dynamics. How well these results generalize to recurrently connected SNNs or more biologically plausible membrane dynamics is unclear.

Our fluctuation-driven initialization strategy follows a similar approach to Glorot et al. [10] and He et al. [11] by setting a target variance for neuronal activity. However, due to the non-continuous nature of the spiking non-linearity, we formulated the goal in terms of the membrane potential variance  $\sigma_U$  instead of the post-non-linearity activation. Our theory results in weight scaling that not only accounts for the number of hidden layer neurons but also data-and architecture-dependent parameters.

In contrast to the above approaches, Mishkin et al. [12] proposed an iterative initialization strategy to achieve unit variance of neuronal activations at each layer during a pre-training period. The implementation of a pre-training period is similar to the homeostatic priming period we applied here. However, instead of setting an explicit target for the population variance, our homeostatic regularizer tuned per-neuron spiking activity to enable activity propagation.

Our work has several limitations. First, our theory is limited to LIF neurons with currentbased synapses. Although the current-based LIF is by far the most commonly used neuron model in SNNs, its synaptic dynamics can allow for biologically implausible and undesirable membrane potential values. Indeed, we found that some neurons exhibit exceptionally small  $(u(t) \ll 0)$  or large  $(u(t) \gg \theta)$  membrane potential values after training, which were not intended when designing the SGs. Future work could explore the possibility of using conductance-based synapses or additional regularization losses to constrain membrane potentials to a biologically plausible range while still allowing for large simulation time steps and thus rapid training.

Second, we performed numerical simulations with a relatively large time step of  $\Delta t = 2$  ms.

Choosing the simulation time step marks a trade-off between computational efficiency on one side and sensitivity of the membrane potential to quickly changing inputs on the other. Indeed, better performing deep SNNs have been trained using a time step on the order of the membrane potential time constant [51, 52]. Our choice of simulation time step reflects a compromise between minimizing computation time and allowing for sufficiently realistic membrane dynamics.

Third, our initialization theory for recurrent SNNs and SNNs following Dale's law, rests on the assumption of balanced input currents, i.e.,  $\mu_U = 0$ , similar to what is observed in neurobiology [19, 20]. Whether and how this balanced state contributes to initial learning phases in the brain remains an open question for experimental and theoretical neuroscience. However, in our numerical simulations, sweeps across the parameters of initial weight distribution in shallow SNNs (Fig. 3) suggest that a slight dominance of inhibition over excitation may represent a similarly favorable or even more advantageous initial state for learning. Therefore, it equally remains to be clarified whether unbalanced currents, for example by a slight dominance of inhibition at initialization, could further support learning in functional SNNs models.

Fourth, our fluctuation-driven initialization theory makes several assumptions that could be violated in some use cases using real-world data. Our theory assumed that all input neurons are independent of each other and fire according to a homogeneous Poisson process with a common firing rate  $\nu$ . Although we have shown that the systematic bias from violating this assumption in the Randman and SHD datasets is not too large (Fig. 2), other datasets with a different spatiotemporal structure could lead to destructive deviations from the theory. As a result, the current initialization strategy could be improved by taking into account more complex firing statistics of the input data. Additionally, our derivations neglected the spike reset of LIFs neurons. While mathematically more complex, it would be possible to consider the reset dynamics in our derivations using a Fokker-Plank approach [21]. However, given that the deviations from the theory due to spatiotemporal structure in the data likely outweigh the contribution of the spike reset, it is questionable whether this extension would confer an advantage.

Finally, we assumed equal firing rates  $\nu = \nu_{dataset}$  for all neurons in a layer in deep SNNs, and for both excitatory and inhibitory populations in Dalian SNNs. Despite being violated for most initialization targets (cf. Fig 4), this simplification allowed for effective initialization with a common target  $\sigma_{U} = 1$  across multiple datasets with vastly different average firing rates (see Methods). Interestingly, for initialization with target  $\sigma_U = 1$  on the SHD dataset, we indeed observed relatively constant firing rates across layers. A consistent method to estimate the firing rate distribution in deep layers at the time of initialization could improve the performance of other initialization targets and could potentially enable training of deeper SNNs. As an alternative approach, dynamic initialization during a pre-training priming period could be extended to adjust weights by regularizing the output firing rate to a target value in an iterative fashion. Similar to approaches that have been proposed for ANNs [12], such an iterative and dynamic initialization strategy could enable activity propagation and learning in even deeper SNNs. However, increasing the number of layers in recurrently connected SNNs did not lead to significant performance improvements in our study. Given the success of deep ANNs, this suggests either that the datasets used to evaluate SNNs are too simple, or that deep SNN architectures and learning algorithms are still in their infancy and could be significantly improved.

In conclusion, the fluctuation-driven initialization proposed in this article facilitates training of diverse SNN architectures in neuromorphic engineering and computational neuroscience by striking a balance between seamless applicability and learning performance. Our work also adds further support to the idea that the fluctuation-driven firing regime, which is widely observed in the brain, may serve as an optimal initial state for future learning, and specifically for scenarios in which learning can be seen as an end-to-end optimization problem [53–55]. While our work only provides the first step toward more effective SNN initialization, it opens up several future exciting directions such as initialization in the presence of sparse connectivity or neuronal cell-type diversity, and suggests that we should take a deeper look at the role of homeostatic plasticity in dynamically preparing networks for optimal learning performance.

# Methods

# Learning tasks

We trained SNNs on several synthetic and real-world classification problems with increasing computational complexity and from different input modalities (auditory, static images, video) to test our initialization strategy. We chose one synthetic dataset and three real-world datasets covering different input modalities to generalize our results across different datasets. In the following, we briefly describe each dataset. The exact specifications of each dataset after preprocessing are summarized in Table 2.

Synthetic random manifolds (Randman). We used a versatile synthetic classification dataset based on precise input spike timings drawn from smooth random manifolds as previously described ([5]; https://github.com/fzenke/randman). The approach allows for flexible dataset generation with different degrees of complexity by varying the number of classes, the intrinsic manifold dimension D, the smoothness parameter  $\alpha$ , and the embedding space dimension M.

Here, we chose parameters to ensure the problem could not be solved by an ANN without a hidden layer. Specifically, we set the embedding space dimension  $M = n_{\rm randman} = 20$ ,  $D = \alpha = 1$  and generated spike trains of 100 ms duration with 10 different classes for all our simulation experiments. To account for delays in activity propagation through the network, we appended 100 ms of no spiking activity to the generated inputs, resulting in a total duration of  $T_{\rm randman} = 200$  ms and hence an average input firing rate of  $\nu_{\rm randman} = 5$  Hz. Further, we used the same random seed to generate the dataset for all experiments in which we compare different initializations to avoid variability due to differences in the dataset. Specifically, we generated a 10-way classification dataset with 1000 samples for each class, 800 of which served as training data and two sets of 100 samples each served as validation and testing data, respectively.

Spiking Heidelberg Digits (SHD). The SHD dataset [23] is a real-world auditory dataset containing recordings of spoken digits (0 - 9) in both German and English from different speakers. It is freely available for download at https://ieee-dataport.org/open-access/ heidelberg-spiking-datasets. To obtain input spikes, the raw audio data was pre-processed by a biologically inspired cochlear model [23] and mapped into an  $n_{\rm SHD} = 700$  dimensional input space. As individual input samples are of different duration, we considered only the first  $T_{\rm SHD} = 700 \,\mathrm{ms}$  of each sample, which corresponds to a fraction > 98 % of all input spikes. Spliced inputs were binned into  $\frac{T_{\rm SHD}}{\Delta t}$  time steps and fed directly into the SNNs. We used a random subset corresponding to 10% of the training data as a validation set. To evaluate generalization performance, we finally used the standard SHD test dataset which contains data from separate speakers that were not included in the training dataset.

**CIFAR-10.** The CIFAR-10 dataset consists of 3x32x32 pixel images belonging to 10 different classes (6000 images for each class) and is commonly used as a visual classification dataset for neural networks [56]. The first dimension of the input data corresponds to the three RGB color channels. As an image dataset, it does not have an intrinsic time dimension in the input. To translate static images into temporal spiking input, we designed an additional sensory neuron encoding layer, placed in between the input layer and the first hidden layer, that converts static

images into spike trains. First, each input pixel data was repeated with a fan-out factor of five along the channel dimension, thereby creating an effective input dimension of 15x32x32 for each image. Second, each pixel value was multiplied by a gain factor to which a bias term was added before using the result as a current input to an encoding layer consisting of 15x32x32 LIF units. The encoding weights for height and width dimensions were tied across all encoding units, leading to an encoding weight matrix of shape 15x1x1. Thus, each encoding neuron receives the weighted pixel value of a single color channel as a constant synaptic current and transduces this value into an output spike train. Synaptic weights of the encoding layer were not subject to our initialization strategy, but randomly drawn from a normal distribution with mean 0 and standard deviation  $\frac{\Delta t}{\tau_{\text{syn}}\sqrt{n_{l-1}}}$ . Biases of encoding units were randomly drawn from a normal distribution with mean 0 and standard deviation  $\frac{1}{\sqrt{n_{enc}}}$ . Both encoding weights and biases were optimized end-to-end during training. For training, CIFAR-10 images were transformed to a normalized range [-1,1] and presented as input to the encoding layer for a duration of  $T_{\text{CIFAR-10}} = 100 \,\text{ms}$ . To obtain an estimate of the firing rate required for the initialization of hidden layers, we measured the average population firing rate of the encoding layer in response to the CIFAR-10 training dataset at the time of initialization, resulting in  $\nu_{\text{CIFAR-10}} = 14.3$  Hz.

**DVS128 Gesture Dataset.** The DVS-gesture dataset [32] is a standard benchmark for event-based processing. It consists of 1342 videos of 11 different hand and/or arm gestures that were recorded with a biologically inspired Dynamic Vision Sensor (DVS), yielding sparse and asynchronous input spike trains. The data from 23 recorded subjects serve as training data, while the data from 6 separate subjects serve as test data. Before training, we applied data augmentation and down-sampling, more specifically (1) random omission of events, (2) down-sampling of the original recordings, and (3) random temporal crop. First, recorded (binary) events were dropped with a probability of p = 0.5. Second, the original 2x128x128 pixels recordings were down-sampled to 2x32x32 pixels. Third, a random 1-second fragment was extracted from each sample. These 1-second long segments were then binned into  $\frac{1000 \text{ ms}}{\Delta t}$  time steps and used as input to the SNN for  $T_{\text{DVS-Gestures}} = 1000 \text{ ms}$ .

	Randman	SHD	CIFAR-10	DVS-Gestures
Duration $T_{\text{dataset}}$ [ms]	200	700	100	1,000
Input dimensions	1	1	2	3
Input neurons $n_{\text{dataset}}$	20	700	$32 \ge 32$	2 x 32 x 32
Firing rate $\nu_{\text{dataset}}$ [Hz]	5	15.8	14.3	9.2
Classes	10	20	10	11
Total training samples	8,000	$7,\!340$	$54,\!000$	1,077

 Table 2. Dataset specifications after pre-processing.

### Network models

All SNN models were trained with SGs using PyTorch [57]. To this end, we used custom software written in Python 3.6.9, which is available on https://github.com/fmi-basel/stork. It includes the fluctuation-driven initialization methods discussed in this paper and example notebooks to replicate all main findings. For numerical simulations, all models were implemented in discrete time with a time step  $\Delta t = 2 \text{ ms}$ . This time step was a compromise between numerical integration accuracy and computational and memory efficiency during training.

**Neuron model.** All units were implemented as simple LIF neurons with exponential currentbased synapses [22]. In discrete time, the membrane potential of neuron i in layer l is characterized by the update equation

$$U_i^{(l)}[n+1] = \left(\lambda_{\rm mem} U_i^{(l)}[n] + (1-\lambda_{\rm mem}) I_i^{(l)}[n]\right) \left(1 - S_i^{(l)}[n]\right) , \qquad (10)$$

where  $U_i^{(l)}[n]$  is this neuron's membrane potential at time step n and  $S_i^{(l)}[n]$  is the associated binary (spiking) output of this neuron defined as  $S_i^{(l)}[n] = \Theta \left( U_i^{(l)}[n] - \theta \right)$  with spike threshold  $\theta$ , where  $\Theta$  is the Heaviside step function. For simplicity, we set  $\theta = 1$ , so that the resting membrane potential is zero and the firing threshold is equal to one. The membrane decay variable  $\lambda_{\text{mem}}$  is determined by the membrane time constant  $\tau_{\text{mem}}$  through  $\lambda_{\text{mem}} \equiv \exp\left(-\frac{\Delta t}{\tau_{\text{mem}}}\right)$ . Lastly,  $I_i^{(l)}[n]$  denotes the incoming synaptic current to neuron i at time step n and is defined as

$$I_i^{(l)}[n+1] = \lambda_{\rm syn} I_i^{(l)}[n] + \sum_j w_{ij}^{(l)} S_j^{(l-1)}[n] + \sum_j v_{ij}^{(l)} S_j^{(l)}[n]$$
(11)

with the feed-forward weight matrix W and optional recurrent weight matrix V. The synaptic decay variable  $\lambda_{\text{syn}}$  is related to the synaptic time constant through  $\lambda_{\text{syn}} \equiv \exp\left(-\frac{\Delta t}{\tau_{\text{syn}}}\right)$ . The neuronal parameters used throughout our simulations can be found in Table 3.

At the beginning of each mini-batch, all neurons were reset to their resting membrane potential of  $U_i^{(l)}[0] = 0$  and a non-spiking state  $S_i^{(l)}[0] = 0$ .

	Non-Dalian SNNs	Dalian SNNs (exc. / inh.)
$\tau_{\rm mem}  [{\rm ms}]$	20	20 / 20
$\tau_{\rm syn} \ [{\rm ms}]$	10	10 / 20

Table 3. Neuronal parameters  $\tau_{\rm mem}$  and  $\tau_{\rm syn}$  used in the numerical simulations of SNNs.

**Readout units.** The units in the readout layer are identical to the above neuron model but were not allowed to spike. Additionally, the membrane time constant of readout units  $\tau_{out}$  could be different from the hidden layer units. Unless otherwise mentioned, we set  $\tau_{out} = T_{data}$  for all simulations to allow readout units to integrate inputs over the entire stimulus duration.

**Dale's Law.** In SNNs obeying Dale's law (cf. Fig. 7), each hidden layer consists of independent excitatory (*E*) and inhibitory (*I*) populations of LIF neurons with membrane time constants  $\tau_{\text{mem}}^{\text{E}}$  and  $\tau_{\text{mem}}^{\text{I}}$ , respectively. In discrete time, the membrane potential of each excitatory or inhibitory neuron *i* in layer *l* is identical to Eq. (10), where  $\lambda_{\text{mem}}$  is replaced with the population-specific decay variables  $\lambda_{\text{mem}}^{E}$  and  $\lambda_{\text{mem}}^{I}$ , respectively (cf. Tab. 3). Like in the non-Dalian case, the decay variables are related to the membrane time constants as  $\lambda_{\text{mem}}^{E} \equiv \exp\left(-\frac{\Delta t}{\tau_{\text{mem}}^{E}}\right)$  and  $\lambda_{\text{mem}}^{I} \equiv \exp\left(-\frac{\Delta t}{\tau_{\text{mem}}^{I}}\right)$ . Contrary to the non-Dalian case, the input currents in Dalian SNNs consist of separate excitatory and inhibitory components originating from distinct presynaptic populations. For both excitatory and inhibitory populations, the input current can therefore be decomposed into feed-forward excitatory (*F*), recurrent excitatory (*R*), and recurrent inhibitory (*I*) components, such that

$$I_i^{(l),E}[n] = I_i^{(l),FE}[n] + I_i^{(l),RE}[n] - I_i^{(l),IE}[n]$$
(12)

$$I_i^{(l),I}[n] = I_i^{(l),FI}[n] + I_i^{(l),RI}[n] - I_i^{(l),II}[n] .$$
(13)

Thus, both the excitatory and inhibitory populations receive two sources of excitatory and one source of inhibitory input. In discrete time, the incoming synaptic currents to the excitatory neuron i of layer l are given as

$$I_{i}^{(l),FE}[n+1] = \lambda_{\text{syn}}^{E} I_{i}^{(l),FE}[n] + \sum_{j} w_{ij}^{(l),FE} S_{j}^{(l-1),E}[n]$$
(14)

$$I_{i}^{(l),RE}[n+1] = \lambda_{\text{syn}}^{E} I_{i}^{(l),RE}[n] + \sum_{j} w_{ij}^{(l),RE} S_{j}^{(l),E}[n]$$
(15)

$$I_i^{(l),IE}[n+1] = \lambda_{\text{syn}}^I I_i^{(l),IE}[n] + \sum_j^J w_{ij}^{(l),IE} S_j^{(l),I}[n] , \qquad (16)$$

where  $\lambda_{\text{syn}}^E$  and  $\lambda_{\text{syn}}^I$  are the decay variables of excitatory and inhibitory currents, which are related to their respective synaptic time constants  $\tau_{\text{syn}}^E$  and  $\tau_{\text{syn}}^I$  (cf. Tab. 3) as described before. Similarly, the synaptic currents into inhibitory neuron *i* of layer *l* are defined as

$$I_{i}^{(l),FI}[n+1] = \lambda_{\text{syn}}^{E} I_{i}^{(l),FI}[n] + \sum_{j} w_{ij}^{(l),FI} S_{j}^{(l-1),E}[n]$$
(17)

$$I_{i}^{(l),RI}[n+1] = \lambda_{\text{syn}}^{E} I_{i}^{(l),RI}[n] + \sum_{j}^{c} w_{ij}^{(l),RI} S_{j}^{(l),E}[n]$$
(18)

$$I_i^{(l),II}[n+1] = \lambda_{\text{syn}}^I I_i^{(l),II}[n] + \sum_j^{J} w_{ij}^{(l),II} S_j^{(l),I}[n] .$$
(19)

Together, the dynamics of each Dalian hidden layer are therefore determined by two feed-forward weight matrices  $W^{FE}$  and  $W^{FI}$  and four recurrent weight matrices  $W^{RE}$ ,  $W^{RI}$ ,  $W^{IE}$ , and  $W^{II}$ .

**Connectivity.** Feed-forward and recurrent networks were all-to-all connected without bias terms unless mentioned otherwise.

We used two types of convolutional networks with 1-dimensional and 2-dimensional convolutional kernels (cf. Tabs. 5 and 6). All CSNNs have recurrently connected layers unless mentioned otherwise. Recurrent connections in CSNNs were implemented as convolutions with filter kernels of size five and a stride of one. For weight initialization of CSNNs, we set  $n = fan_{in}$ , the number of inputs to each filter.

In SNNs and CSNNs obeying Dale's law, weights between consecutive layers and recurrent weights within hidden layers were constrained to be positive both at initialization and continuously during training, except for the readout weights, which were not sign constrained. Both excitatory and inhibitory populations in hidden layers received feed-forward inputs from the excitatory population of the previous layer. All networks obeying Dale's law were fully recurrent, featuring recurrent connections within and between excitatory and inhibitory populations in each layer ( $E \rightarrow E$ ,  $E \rightarrow I$ ,  $I \rightarrow I$ , and  $I \rightarrow E$ ).

**Skip connections.** We implemented skip connections as additional all-to-all connections with trainable weights between each except the last hidden layer and the readout layer, such that the readout units receive a separate input from every hidden layer (cf. Supplementary Fig. S10).

**Supervised loss function.** All networks were trained by minimizing a standard cross-entropy loss

$$\mathcal{L}_{\text{sup}} = -\frac{1}{K} \sum_{k=1}^{K} \sum_{c=1}^{C} y_c^k \log\left(p_c^k\right) , \qquad (20)$$

where the one-hot encoded target for input k is denoted by  $y_c^k$ , K is the number of input samples and C is the number of classes. The associated output probabilities  $p_c^k$  are given by the Softmax function

$$p_{c}^{k} = \frac{\exp(a_{c}^{k})}{\sum_{i=1}^{C} \exp(a_{i}^{k})} .$$
(21)

The scores  $a_c^k$  for each input k are dependent on the membrane potential of the associated readout units  $U_c^{(\text{out})}$  and can take different forms. For all simulations in this paper, we defined the score as the maximum value over all time steps  $a_c^k = \max_n \left( U_c^{(\text{out})}[n] \right)$ .

Activity regularization. Unless otherwise mentioned, all networks were subject to activity regularization to constrain spiking activity. To that end, we added loss terms corresponding to a soft upper bound on the population-level spiking activity for each layer l as

$$g_{\rm upper}^{(l),k} = \left( \left[ \frac{1}{M^{(l)}} \sum_{i}^{M^{(l)}} \zeta_{i}^{(l),k} - v_{\rm upper} \right]_{+} \right)^{2} , \qquad (22)$$

where  $\zeta_i^{(l),k} = \left(\sum_n^N S_i^{(l),k}[n]\right)$  is the spike count of neuron *i* in layer *l* given input sample *k* and  $M^{(l)}$  is the number of neurons in hidden layer *l*. In 1-dimensional CSNNs receiving auditory inputs, we set  $M^{(l)} = n_{\text{features}}^{(l)} \times n_{\text{neurons}}^{(l)}$  and, similarly, in 2-dimensional CSNNs receiving visual inputs, we set  $M^{(l)} = n_{\text{features}}^{(l)} \times n_x^{(l)} \times n_y^{(l)}$  with  $n_x^{(l)}$  and  $n_y^{(l)}$  denoting the number of *x* and *y* coordinates in the layer, respectively. This activity regularization effectively prevents the population-level activity from exceeding the threshold spike count  $v_{\text{upper}}$ , which we set to  $v_{\text{upper}} = \frac{T_{\text{data}}}{100}$  to achieve an upper bound average population firing rate of 10 Hz per layer. The regularization loss in case of population-level upper bound for spiking activity  $\mathcal{L}_{\text{UB}}$  would thus be

$$\mathcal{L}_{\rm UB} = -\lambda_{\rm upper} \sum_{l}^{L} g_{\rm upper}^{(l),k} , \qquad (23)$$

where  $\lambda_{upper}$  denotes the strength of the regularization.

**Homeostatic plasticity.** In networks with homeostatic plasticity (cf. Figs. 5 and 6) we added an additional term to the total loss acting as a per-neuron lower bound on the spiking activity. This per-neuron lower bound loss on spiking activity  $\mathcal{L}_{HP}$  was defined as

$$g_{\text{lower}}^{(l),k} = \frac{1}{M^{(l)}} \sum_{i}^{M^{(l)}} \left( \left[ -\left(\zeta_{i}^{(l),k} - v_{\text{lower}}\right) \right]_{+} \right)^{2}$$
(24)

$$\mathcal{L}_{\rm HP} = -\lambda_{\rm lower} \sum_{l}^{L} g_{\rm lower}^{(l),k} , \qquad (25)$$

where the first equation describes the per-neuron loss term for each layer l.  $\zeta_i^{(l),k}$  corresponds to the spike count of neuron i in layer l,  $M^{(l)}$  is the number of neurons in layer l,  $v_{\text{lower}}$  denotes a lower bound on the spike count and  $\lambda_{\text{lower}}$  is the regularizer strength.

With  $v_{\text{lower}} = 1$ , this additional regularization term penalizes neurons that do not spike and thus ensures spiking activity in each neuron. By setting  $v_{\text{lower}}$  to other positive values one may achieve a desired lower bound on the per-neuron firing rate.

**Surrogate gradient descent.** To minimize the loss  $\mathcal{L}$ , we adjusted network parameters in the direction of the negative SG. We computed SGs for the parameter updates using BPTT and the automatic differentiation capabilities of PyTorch [57]. Because the spiking non-linearity of the spiking neuron model is not differentiable, we approximate its derivative

$$S'\left(U_i^{(l)}[n]\right) = \Theta'\left(U_i^{(l)}[n] - \theta\right)$$
(26)

with the surrogate

$$\tilde{S}'\left(U_i^{(l)}[n]\right) = h\left(U_i^{(l)}[n] - \theta\right) .$$
(27)

Throughout this study, we use the SuperSpike surrogate non-linearity [46]

$$h(x) = \frac{1}{(\beta|x|+1)^2}$$
(28)

with steepness parameter  $\beta = 20$ . For the simulations of deep CSNNs with a rescaled SG non-linearity reported in Supplementary Figure S8, we used the re-scaled surrogate derivative

$$\tilde{S}'\left(U_i^{(l)}[n]\right) = \frac{h\left(U_i^{(l)}[n] - \theta\right)}{h\left(\theta\right)} .$$
<sup>(29)</sup>

In this case, the surrogate derivative at rest is equal to one, i.e.,  $\tilde{S}'(0) = 1$ , where "at rest" refers to the absence of input to the corresponding neuron, causing its membrane potential to remain at zero. Thus, using this rescaled non-linearity, and in the absence of any membrane potential fluctuations, gradient magnitudes do not decay during backpropagation over the inactive layers.

**Optimizer.** We used the SMORMS3 optimizer [25] unless mentioned otherwise. Given a parameter  $\theta$ , SMORMS3 performs the following update step after every mini-batch:

$$g_1^{(\theta)} := \left(1 - r^{(\theta)}\right) g_1^{(\theta)} + r^{(\theta)} \left(\frac{\partial \mathcal{L}}{\partial \theta}\right)$$
(30)

$$g_2^{(\theta)} := \left(1 - r^{(\theta)}\right) g_2^{(\theta)} + r^{(\theta)} \left(\frac{\partial \mathcal{L}}{\partial \theta}\right)^2 \tag{31}$$

$$m^{(\theta)} := 1 + m^{(\theta)} \left( 1 - \frac{\left(g_1^{(\theta)}\right)^2}{g_2^{(\theta)} + \epsilon} \right) ,$$
 (32)

where  $r^{(\theta)} = \frac{1}{m^{(\theta)}+1}$  and  $\epsilon = 1 \times 10^{-16}$  is a small positive value to avoid division by zero. Before the first training epoch, the optimizer state variables are initialized as  $g_1^{(\theta)} = g_2^{(\theta)} = 0$  and  $m^{(\theta)} = 1$ . The parameter update after each mini-batch is then performed as

$$\Delta \theta = -\left(\frac{\partial \mathcal{L}}{\partial \theta}\right) \min\left(\eta, \frac{\left(g_1^{(\theta)}\right)^2}{g_2^{(\theta)} + \epsilon}\right) \frac{1}{\sqrt{g_2^{(\theta)}} + \epsilon} , \qquad (33)$$

where  $\eta$  is the base learning rate.

**Overview of network architectures.** The diverse network architectures considered in this article were chosen as a compromise between memory requirements, performance, and training time. In the following, we provide a brief overview, while giving the precise parameter values in Tables 4-6.

Fully connected shallow SNNs. In Figure 3 and Supplementary Figures S2, S5 we implemented fully connected feed-forward SNNs with a single hidden layer. These networks were trained either on the Randman or the SHD dataset. For both tasks, we used 128 neurons in the hidden layer and adjusted the sizes of the input and readout layers to match the dataset requirements. For the fully connected recurrent SNNs used in Figure 7 and Supplementary Figure S5, we added fully connected recurrent weights to the hidden layer. The recurrent SNNs in Figure 7 had a wider hidden layer with 160 neurons, to match the number of hidden neurons in the Dalian network. The exact network specifications are summarized in Table 4.

	Randman	SHD
No. input neurons	20	700
No. output neurons	10	20
No. hidden neurons	128	128 or 160
Dalian SNNs: No. hidden neurons	-	128 exc. / 32 inh.
Mini-batch size	400	400
No. training epochs	200	200

**Table 4.** Network and training parameters used for simulations of fully connected SNNs with a single hidden layer on the Randman and SHD datasets.

Fully connected shallow SNNs following Dale's law. Shallow SNNs following Dale's law (cf. Fig. 7) had one hidden layer with 160 neurons, 128 of which were excitatory and 32 inhibitory, following a four-to-one ratio between excitatory and inhibitory neurons commonly observed in the mammalian cortex. Dalian networks were fully recurrently connected within the hidden layer through  $I \to E$ ,  $I \to I$ ,  $E \to I$ , and  $E \to E$  connections. Feed-forward connections from the input to the hidden layer populations were constrained to be excitatory. Readout units received inputs solely from the excitatory population of the hidden layer, but readout weights were not subject to a sign constraint.

Deep feed-forward convolutional SNNs. Deep feed-forward CSNNs (cf. Figs. 4, 6; Supplementary Figs. S6, S8), were implemented with up to ten consecutive hidden layers and trained on the SHD, CIFAR-10, or the DVS-Gesture datasets. CSNNs trained on the CIFAR-10 or DVS-Gesture datasets additionally implemented a max pooling operation with a kernel size of  $2 \times 2$  after every second hidden layer. Network sizes and parameters of the convolutional operations are summarized in Tables 5 and 6. These architectures were chosen to create networks of different depths with similar widths while ensuring that deeper layers still contained a reasonable number of neurons.

Deep recurrent convolutional SNNs. Unless mentioned otherwise, all CSNNs trained on the SHD dataset additionally implemented recurrent connections in each hidden layer (cf. Figs. 4-7; Supplementary Figs. S4, S7-S10). Recurrent connections in CSNNs were implemented as convolutional operations with a kernel size of five and a stride of one. The exact parameters of recurrent CSNNs trained on the SHD dataset can be found in Table 5.

Deep recurrent CSNNs following Dale's law. CSNNs following Dale's law (cf. Fig. 7) were implemented with separate excitatory and inhibitory populations in each hidden layer. Except for readout connections, which were not sign constrained, all feed-forward connections were constrained to be excitatory. As in the case of shallow SNNs following Dale's law, each hidden layer consisted of separate excitatory and inhibitory populations. Dalian CSNNs followed the same architecture as non-Dalian CSNNs (Tab. 5) for excitatory neurons, and each hidden layer was augmented with an additional population of inhibitory neurons of size  $N_{exc}/4$ . Excitatory  $E \to I$  and  $E \to E$  recurrent connections were implemented as convolutional operations with a kernel size of 5 and a stride of 1. Inhibitory  $I \to I$  and  $I \to E$  recurrent connections were implemented as convolutional operations with a kernel size of 3 and a stride of 1.

Dataset			SHD	
No. input neurons			700	
No. output neurons			20	
No. training epochs			200	
No. hidden layers	1	3	7	10
Mini-batch size	400	400	100	100
No. hidden neurons	16	16-32-64	$16-32-64-\ldots-64$	$16-32-64-\ldots-64$
Kernel size (ff)	21	21-7-7	7-55	5
Stride (ff)	10	10-3-3	3-22	2
Padding (ff)	2	2	2	2
No. parameters (ff)	24,858	$24,\!656$	99,952	$157,\!520$
Kernel size (rec)	5	5	5	5
Stride (rec)	1	1	1	1
No. parameters (rec)	52,734	$51,\!536$	208,752	327,760

**Table 5.** Network and training parameters used for simulations of deep convolutional SNNson the SHD dataset.

Dataset	CIFAR-10		DVS-Gestures		
No. input neurons	3	$2 \times 32$	$2 \times 32 \times 32$		
No. output neurons		10		11	
No. training epochs		50		20	
No. hidden layers	<b>2</b>	4	6	8	
Mini-batch size	128	128	16	8	
No. hidden neurons	32-32	32-32-64-64	32-32-64-64-128-128	32-32-64-64-128128	
Kernel size	$3 \times 3$	$3 \times 3$	3  imes 3	3  imes 3	
Stride	1	1	1	1	
Padding	2	2	2	2	
No. parameters	$95,\!456$	109,792	$308,\!800$	$586,\!816$	

**Table 6.** Network and training parameters used for simulations of deep convolutional SNNs on the CIFAR-10 and DVS-Gestures datasets.

### Weight initialization

For fluctuation-driven initialization of synaptic weight parameters, the PSP-kernel parameters  $\bar{\epsilon}$  and  $\hat{\epsilon}$  introduced in Equations (2) and (3) can be computed analytically or numerically (see Supplementary Material S1). Because we used a relatively large time step of  $\Delta t = 2 \text{ ms}$  for which there are non-negligible differences between the two, we used the numerical integration values for all simulations as they are closer to the actual simulation (Tab. 7).

For strictly feed-forward networks, the fluctuation-driven initialization strategy was already covered in the main text. In the following, we derive the extensions to deep convolutional SNNs, recurrent SNNs, and SNNs obeying Dale's law.

Fluctuation-driven initialization of recurrent networks. For the initialization of recurrent layers, we introduce the additional parameter  $0 < \alpha < 1$ , that determines the proportion of membrane potential fluctuations caused by *feed-forward* connections in contrast to *recurrent* 

	Non-Dalian SNNs	Dalian SNNs (exc. / inh.)
$\bar{\epsilon}$	0.0110	0.0110 / 0.0061
$\hat{\epsilon}$	0.0020	$0.0020\ /\ 0.0012$

**Table 7.** Values of the PSP-kernel integrals  $\bar{\epsilon}$  and  $\hat{\epsilon}$  used for weight initialization in the numerical simulations, rounded to four decimal places. Due to the large simulation time step of  $\Delta t = 2 \text{ ms}$ ,  $\bar{\epsilon}$  and  $\hat{\epsilon}$  were obtained numerically. The analytical expressions for  $\bar{\epsilon}$  and  $\hat{\epsilon}$  can be found in Supplementary Table S1.

connections:

$$\alpha = \frac{\text{Part of } \sigma_U^2 \text{ caused by feed-forward connections}}{\text{Total } \sigma_U^2} .$$
(34)

To this end, we consider a postsynaptic LIF neuron receiving feed-forward input from  $n_F$  neurons with firing rate  $\nu_F$  and recurrent input from  $n_R$  hidden layer neurons with average firing rate  $\nu_R$ . Feed-forward weights are initialized as  $W \sim \mathcal{N}(\mu_W, \sigma_W^2)$  and recurrent weights are initialized as  $V \sim \mathcal{N}(\mu_V, \sigma_V^2)$  The mean  $\mu_U$  and variance  $\sigma_U^2$  of the membrane potential are then given by

$$u_U = n_F \mu_W \nu_F \bar{\epsilon} + n_R \mu_V \nu_R \bar{\epsilon} \tag{35}$$

$$\sigma_U^2 = n_F (\sigma_W^2 + \mu_W^2) \nu_F \hat{\epsilon} + n_R (\sigma_V^2 + \mu_V^2) \nu_R \hat{\epsilon} .$$
(36)

In practice the firing rate  $\nu_R$  of the hidden layer is difficult to predict due to finite-size effects. Hence, we make the simplifying assumption  $\nu = \nu_F = \nu_R = \nu_{\text{dataset}}$ . In other words, we assume that the average firing rate of the hidden layers is equal to the input firing rate.

Since we want  $\alpha$  to control the membrane potential *fluctuations* only, which are determined by  $\sigma_W^2$  and  $\sigma_V^2$ , we can initialize recurrent and feed-forward weights with a common mean, i.e.,  $\mu_{WV} = \mu_W = \mu_V$  defined as

$$\mu_{WV} = \frac{\mu_U}{(n_F + n_R) \,\nu_{\text{dataset}} \bar{\epsilon}} \tag{37}$$

and subsequently solve for  $\sigma_W^2$  and  $\sigma_V^2$  independently:

 $\mu$ 

$$\sigma_W^2 = \frac{\alpha}{n_F \nu \hat{\epsilon}} \left(\frac{\theta - \mu_U}{\xi}\right)^2 - \mu_{WV}^2 \tag{38}$$

$$\sigma_V^2 = \frac{1-\alpha}{n_R \nu \hat{\epsilon}} \left(\frac{\theta - \mu_U}{\xi}\right)^2 - \mu_{WV}^2 .$$
(39)

In this article, we used  $\alpha = 0.9$  for all simulations with recurrently connected hidden layers unless stated otherwise, so that the majority of membrane potential fluctuations originate from feed-forward input.

Fluctuation-driven initialization of Dalian networks. Networks following Dale's law consist of separate excitatory and inhibitory populations whose output weights are sign constrained. To initialize the sign constrained connections, we relied on exponential or log-normal weight distributions instead of normally distributed weights, where the choice of a log-normal distribution is inspired by findings from neurobiology [58].

Parameterizing the excitatory and inhibitory weight distributions with  $\lambda$  for the exponential and  $\mu$  for the log-normal distribution, respectively, allows us to obtain explicit expressions for the initial weight distributions leading to the target membrane potential fluctuations with mean  $\mu_U$  and variance  $\sigma_U^2$ . Unless stated otherwise, the weights in Dalian SNNs were initialized using the exponential distribution throughout the numerical simulations. While we provide here the expression for weight initialization using the exponential distribution, a derivation for lognormally distributed initial weights can be found in the Supplementary Material S3.

We start by observing that, regardless of the weight distribution from which synaptic weights are sampled, the mean and the variance of the membrane potential of a neuron i in a Dalian network are defined as

$$\mu_U^{(i)} = \sum_j^{n_E} w_{ij}^E \nu_E \bar{\epsilon}_E - \sum_k^{n_I} w_{ik}^I \nu_I \bar{\epsilon}_I$$
(40)

$$\left(\sigma_{U}^{(i)}\right)^{2} = \sum_{j}^{n_{E}} (w_{ij}^{E})^{2} \nu_{E} \hat{\epsilon}_{E} + \sum_{k}^{n_{I}} (w_{ik}^{I})^{2} \nu_{I} \hat{\epsilon}_{I} , \qquad (41)$$

where we assume equal firing rates  $\nu_E$  and  $\nu_I$  for all excitatory and inhibitory neurons in our experiments, respectively.

For weights drawn from exponential distributions, i.e.,  $w^E \sim \text{Exp}(\lambda_E)$  and  $w^I \sim \text{Exp}(\lambda_I)$ with mean  $\frac{1}{\lambda}$  and variance  $\frac{1}{\lambda^2}$ , we can rewrite the mean and the variance of the membrane potential for each neuron as

$$\mu_U = \frac{n_E \nu_E \bar{\epsilon}_E}{\lambda_E} - \frac{n_I \nu_I \bar{\epsilon}_I}{\lambda_I}$$

$$\sigma_U^2 = n_E \left( \frac{1}{\lambda_E^2} + \left( \frac{1}{\lambda_E} \right)^2 \right) \nu_E \hat{\epsilon}_E + n_I \left( \frac{1}{\lambda_I^2} + \left( \frac{1}{\lambda_I} \right)^2 \right) \nu_I \hat{\epsilon}_I$$

$$= \frac{2n_E \nu_E \hat{\epsilon}_E}{\lambda_E^2} + \frac{2n_I \nu_I \hat{\epsilon}_I}{\lambda_I^2} .$$
(42)
(43)

We further assume that the target  $\mu_U = 0$ , as would be expected in balanced networks. From the definition of  $\mu_U$ , we obtain an explicit relationship between  $\lambda_I$  and  $\lambda_E$ 

$$\lambda_I = \lambda_E \frac{n_I \nu_I \bar{\epsilon}_I}{n_E \nu_E \bar{\epsilon}_E} , \qquad (44)$$

which we use to define a combined E/I ratio based on network parameters

$$\Delta_{EI} = \frac{n_I \nu_I \bar{\epsilon}_I}{n_E \nu_E \bar{\epsilon}_E} \ . \tag{45}$$

Substitution of this relationship into equation (43) gives us

$$\sigma_U^2 = \frac{2n_E\nu_E\hat{\epsilon}_E}{\lambda_E^2} + \frac{2n_I\nu_I\hat{\epsilon}_I}{(\lambda_E\Delta_{EI})^2} .$$
(46)

Finally, we can solve the above for  $\lambda_E$ 

$$\lambda_E = \frac{\sqrt{2(\Delta_{EI}^2 n_E \nu_E \hat{\epsilon}_E + n_I \nu_I \hat{\epsilon}_I)}}{\sigma_U \Delta_{EI}} . \tag{47}$$

Together, equations (44) and (47) allow us to parameterize excitatory and inhibitory weights as a function of  $\sigma_U$ , taking into account data- and network-dependent parameters, which is summarized in Table 8. Note that this initialization relies on a target membrane potential mean  $\mu_U = 0$ .

Network architecture	Weight distribution	Weight parameters	Good regime for initialization
Feed-forward	Centered: $W \sim \mathcal{N}\left(0, \sigma_W^2\right)$	$\sigma_W^2 = \frac{\sigma_U^2}{n\nu\hat{\epsilon}}$	$\frac{1}{3} \le \sigma_U \le 1$
recu-tor ward	Non-centered: $W \sim \mathcal{N}\left(\mu_W, \sigma_W^2\right)$	$\mu_W = \frac{\mu_U}{n\nu\bar{\epsilon}}$ $\sigma_W^2 = \frac{1}{n\nu\hat{\epsilon}} \left(\frac{\theta - \mu_U}{\xi}\right)^2 - \mu_W^2$	$\mu_U < \theta$ $1 \le \xi \le 3$
	Centered: $W \sim \mathcal{N}\left(0, \sigma_W^2\right)$ $V \sim \mathcal{N}\left(0, \sigma_V^2\right)$	$\sigma_W^2 = \alpha \frac{\sigma_U^2}{n_F \nu \hat{\epsilon}}$ $\sigma_V^2 = (1 - \alpha) \frac{\sigma_U^2}{n_R \nu \hat{\epsilon}}$	$\frac{1}{3} \le \sigma_U \le 1$ $0 < \alpha < 1$
Recurrent	Non-centered: $W \sim \mathcal{N}\left(\mu_{WV}, \sigma_W^2\right)$ $V \sim \mathcal{N}\left(\mu_{WV}, \sigma_V^2\right)$	$\mu_{WV} = \frac{\mu_U}{(n_F + n_R)\nu\bar{\epsilon}}$ $\sigma_W^2 = \frac{\alpha}{n_F\nu\hat{\epsilon}} \left(\frac{\theta - \mu_U}{\xi}\right)^2 - \mu_{WV}^2$ $\sigma_V^2 = \frac{1 - \alpha}{n_R\nu\hat{\epsilon}} \left(\frac{\theta - \mu_U}{\xi}\right)^2 - \mu_{WV}^2$	$\mu_U < \theta$ $1 \le \xi \le 3$ $0 < \alpha < 1$

Table 8. Summary of strategies for fluctuation-driven initialization of SNNs.

Fluctuation-driven initialization of Dalian networks with excitatory recurrence. Dalian layers are always recurrently connected, as they require a connection between the separate excitatory and inhibitory populations in each layer. For this reason, the Dalian network from the above paragraph has recurrent inhibitory connections  $(I \rightarrow E \text{ and } I \rightarrow I)$ . Here, we consider the case of additional excitatory recurrence, i.e.  $E \rightarrow I$  and  $E \rightarrow E$  connections.

Again, we require inhibitory currents to balance excitatory currents on average to achieve a mean membrane potential  $\mu_U = 0$ . Additionally, similar to non-Dalian SNNs with recurrent connections, the parameter  $\alpha$  describes the proportion of excitatory membrane potential fluctuations that are caused by feed-forward excitatory connections, whereas the proportion of recurrent excitation is given by  $(1 - \alpha)$ . For the derivation, we consider a single neuron in a Dalian layer, receiving one recurrent inhibitory (I), one feed-forward excitatory (F), and one recurrent excitatory (R) input connection. In this setting, mean  $\mu_U$  and variance  $\sigma_U^2$  of the membrane potential of that neuron are given by

$$u_U = \frac{N_F \nu_F \bar{\epsilon}_E}{\lambda_F} + \frac{N_R \nu_R \bar{\epsilon}_E}{\lambda_R} - \frac{N_I \nu_I \bar{\epsilon}_I}{\lambda_I}$$
(48)

$$\sigma_U^2 = \frac{2N_F \nu_F \hat{\epsilon}_E}{\lambda_F^2} + \frac{2N_R \nu_R \hat{\epsilon}_E}{\lambda_R^2} + \frac{2N_I \nu_I \hat{\epsilon}_I}{\lambda_I^2} .$$
(49)

For the sake of simpler notation, we assume  $\nu = \nu_F = \nu_R = \nu_I$  in this derivation. We also made this assumption of equal firing rates in the application of this initialization strategy in our numerical simulations. Since it is not possible to estimate the firing rates of excitatory and inhibitory hidden neuron populations in advance, we chose  $\nu = \nu_{dataset}$ .

The ratio of membrane potential fluctuations caused by the excitatory feed-forward connections compared to the total excitatory input, which we defined as  $\alpha$ , can explicitly be written as

$$\alpha = \frac{\text{Part of } \sigma_U^2 \text{ caused by excitatory feed-forward connections}}{\text{Part of } \sigma_U^2 \text{ caused by all excitatory connections}} = \frac{\frac{2N_F \nu \epsilon_E}{\lambda_F^2}}{\frac{2N_F \nu \hat{\epsilon}_E}{\lambda_F^2}} + \frac{2N_R \nu \hat{\epsilon}_E}{\lambda_R^2}, \quad (50)$$

which we can solve for  $\lambda_R$  to obtain

$$\lambda_R = \lambda_F \sqrt{\frac{\alpha N_R}{N_F - \alpha N_F}} = \lambda_F \Delta_R , \qquad (51)$$

where we introduced the scalar  $\Delta_R$  to make subsequent notation easier. We can then insert Eq. (51) into Eq. (48)

$$\mu_U = \frac{N_F \nu \bar{\epsilon}_E}{\lambda_F} + \frac{N_R \nu \bar{\epsilon}_E}{\lambda_F \Delta_R} - \frac{N_I \nu \bar{\epsilon}_I}{\lambda_I}$$
(52)

to receive an expression for  $\lambda_I$ 

$$\lambda_I = \lambda_F \frac{\Delta_R \bar{\epsilon}_I N_I}{\Delta_R \bar{\epsilon}_E N_F + \bar{\epsilon}_E N_R} = \lambda_F \Delta_{EI}^R .$$
(53)

We introduce here again a network-parameter dependent scalar  $\Delta_{EI}^{R}$ . Using the scalars  $\Delta_{R}$  and  $\Delta_{EI}^{R}$ , we can now substitute both  $\lambda_{I}$  and  $\lambda_{R}$  in Eq. (49) to obtain

$$\sigma_U^2 = \frac{2N_F \nu \hat{\epsilon}_E}{\lambda_F^2} + \frac{2N_R \nu \hat{\epsilon}_E}{\left(\lambda_F \Delta_R\right)^2} + \frac{2N_I \nu \hat{\epsilon}_I}{\left(\lambda_F \Delta_{EI}^R\right)^2} , \qquad (54)$$

which can be solved for  $\lambda_F$ 

$$\lambda_F = \frac{\sqrt{2\nu \left( (\Delta_{EI}^R)^2 \hat{\epsilon}_E N_R + \Delta_R^2 \left( \Delta_{EI}^R N_F \hat{\epsilon}_E + N_I \hat{\epsilon}_I \right) \right)}}{\sigma_U \Delta_R \Delta_{EI}^R} \ . \tag{55}$$

Equations (51), (53) and (55) let us parameterize the initial feed-forward excitatory (F), recurrent excitatory (R), and recurrent inhibitory (I) weight distributions as a function of the target membrane potential fluctuations  $\sigma_U$  with a mean membrane potential  $\mu_U = 0$ . The suggested weight distributions including their parameters and a range of values for a good initialization are summarized in Table 9.

Kaiming (He) initialization. We implemented Kaiming (He) initialization as described by He et al. [11]. This commonly used strategy was originally derived for ANNs with ReLU non-linearities and purports drawing the initial weights from a centered normal distribution

$$W \sim \mathcal{N}\left(0, \frac{2}{n}\right)$$
, (56)

where n is the number of neurons and the weights have mean zero and variance  $\sigma_W^2 = \frac{2}{n}$ .

e for	
on	

Published as: Rossbroich, Gygax, and Zenke (2022). Neuromorph Comput Eng 10.1088/2634-4386/ac97bb

Network architecture	Weight distribution	Weight parameters		Good regime for initialization
Feed-forward	$W_E \sim \operatorname{Exp}(\lambda_E)$ $W_I \sim \operatorname{Exp}(\lambda_I)$	$\lambda_E = \frac{\sqrt{2(\Delta_{EI}^2 n_E \nu_E \hat{\epsilon}_E + n_I \nu_I \hat{\epsilon}_I)}}{\sigma_U \Delta_{EI}}$ $\lambda_I = \Delta_{EI} \lambda_E$	$\Delta_{EI} = \frac{n_I \nu_I \bar{\epsilon}_I}{n_E \nu_E \bar{\epsilon}_E}$	$\mu_U = 0$ $\frac{1}{3} \le \sigma_U \le 1$
Recurrent	$W_F \sim \operatorname{Exp}(\lambda_F)$ $W_R \sim \operatorname{Exp}(\lambda_R)$ $W_I \sim \operatorname{Exp}(\lambda_I)$	$\lambda_F = \frac{\sqrt{2\nu \left( (\Delta_{EI}^R)^2 \hat{\epsilon}_E N_R + \Delta_R^2 \left( \Delta_{EI}^R N_F \hat{\epsilon}_E + N_I \hat{\epsilon}_I \right) \right)}}{\sigma_U \Delta_R \Delta_{EI}^R}$ $\lambda_R = \lambda_F \Delta_R$ $\lambda_I = \lambda_F \Delta_{EI}^R$	$\Delta_R = \sqrt{\frac{\alpha N_R}{N_F - \alpha N_F}}$ $\Delta_{EI}^R = \frac{\Delta_R \bar{\epsilon}_I N_I}{\Delta_R \bar{\epsilon}_E N_F + \bar{\epsilon}_E N_R}$	$\mu_U = 0$ $\frac{1}{3} \le \sigma_U \le 1$ $0 < \alpha < 1$

Table 9. Summary of strategies for fluctuation-driven initialization of Dalian SNNs.

# Acknowledgments

This work was supported by the Novartis Research Foundation and the Swiss National Science Foundation [grant number PCEFP3\_202981].

# Author contributions

F.Z. conceived the study. J.R., J.G., and F.Z. wrote simulation code. J.R. and J.G. performed simulations and analyses. J.R., J.G., and F.Z. wrote the manuscript.

# **Competing interests**

The authors declare no competing interests.

# References

- [1] Sterling, P. and Laughlin, S. Principles of Neural Design. The MIT Press, 2017.
- [2] Indiveri, G., Linares-Barranco, B., Hamilton, T., Schaik, A. van, Etienne-Cummings, R., Delbruck, T., Liu, S.-C., Dudek, P., Häfliger, P., Renaud, S., Schemmel, J., Cauwenberghs, G., Arthur, J., Hynna, K., Folowosele, F., Saïghi, S., Serrano-Gotarredona, T., Wijekoon, J., Wang, Y., and Boahen, K. "Neuromorphic Silicon Neural Circuits". In: *Frontiers in Neuroscience* 5 (2011).
- [3] Poole, B., Lahiri, S., Raghu, M., Sohl-Dickstein, J., and Ganguli, S. "Exponential expressivity in deep neural networks through transient chaos". In: Advances in Neural Information Processing Systems. Ed. by Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R. Vol. 29. Curran Associates, Inc., 2016.
- [4] Hunsberger, E. and Eliasmith, C. Spiking Deep Networks with LIF Neurons. Version 1. 2015. arXiv: 1510.08829 [cs.LG].
- [5] Zenke, F. and Vogels, T. P. "The Remarkable Robustness of Surrogate Gradient Learning for Instilling Complex Function in Spiking Neural Networks". In: *Neural computation* 33.4 (2021), pp. 899–925.
- [6] Neftci, E. O., Mostafa, H., and Zenke, F. "Surrogate Gradient Learning in Spiking Neural Networks: Bringing the Power of Gradient-Based Optimization to Spiking Neural Networks". In: *IEEE Signal Processing Magazine* 36.6 (2019), pp. 51–63.
- [7] Hochreiter, S. "Untersuchungen zu dynamischen neuronalen Netzen". MA thesis. Technische Universität München, 1991.
- [8] Hochreiter, S. and Schmidhuber, J. "Long short-term memory". In: Neural computation 9.8 (1997), pp. 1735–1780.
- [9] Pascanu, R., Mikolov, T., and Bengio, Y. "On the difficulty of training recurrent neural networks". In: *ICML*. 2013.
- [10] Glorot, X. and Bengio, Y. "Understanding the difficulty of training deep feedforward neural networks". In: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. Ed. by Teh, Y. W. and Titterington, M. Vol. 9. Proceedings of Machine Learning Research. PMLR, 2010, pp. 249–256.
- [11] He, K., Zhang, X., Ren, S., and Sun, J. "Delving deep into rectifiers: Surpassing humanlevel performance on ImageNet classification". In: 2015 IEEE International Conference on Computer Vision (ICCV). Santiago, Chile: IEEE, 2015, pp. 1026–1034.

- [12] Mishkin, D. and Matas, J. All you need is a good init. Version 7. 2015. arXiv: 1511.06422 [cs.LG].
- [13] Srivastava, R. K., Greff, K., and Schmidhuber, J. Training Very Deep Networks. Version 2. 2015. arXiv: 1507.06228 [cs.LG].
- [14] Lee, J. H., Delbruck, T., and Pfeiffer, M. "Training deep spiking neural networks using backpropagation". In: *Frontiers in Neuroscience* 10 (2016).
- [15] Ledinauskas, E., Ruseckas, J., Juršėnas, A., and Buračas, G. Training Deep Spiking Neural Networks. Version 1. 2020. arXiv: 2006.04436 [cs.CV].
- [16] Tiesinga, P. H., José, J. V., and Sejnowski, T. J. "Comparison of current-driven and conductance-driven neocortical model neurons with Hodgkin-Huxley voltage-gated channels". In: *Physical review. E, Statistical physics, plasmas, fluids, and related interdisciplinary topics* 62.6 Pt B (2000), pp. 8413–8419.
- [17] "Neuronal integration of synaptic input in the fluctuation-driven regime". In: *The Journal of neuroscience: the official journal of the Society for Neuroscience* 24.10 (2004), pp. 2345–2356.
- [18] Petersen, P. C. and Berg, R. W. "Lognormal firing rate distribution reveals prominent fluctuation-driven regime in spinal motor networks". In: *eLife* 5 (2016), e18805.
- [19] Vogels, T. P., Rajan, K., and Abbott, L. F. "Neural network dynamics". In: Annual Review of Neuroscience 28 (2005), pp. 357–376.
- [20] Brunel, N. "Dynamics of sparsely connected networks of excitatory and inhibitory spiking neurons". In: *Journal of Computational Neuroscience* 8.3 (2000), pp. 183–208.
- [21] Amit, D. J. and Brunel, N. "Model of global spontaneous activity and local structured activity during delay periods in the cerebral cortex". In: *Cerebral cortex* 7.3 (1997), pp. 237– 252.
- [22] Gerstner, W., Kistler, W. M., Naud, R., and Paninski, L. Neuronal Dynamics: From Single Neurons to Networks and Models of Cognition. Cambridge University Press, 2014.
- [23] Cramer, B., Stradmann, Y., Schemmel, J., and Zenke, F. "The Heidelberg Spiking Data Sets for the Systematic Evaluation of Spiking Neural Networks". In: *IEEE Transactions* on Neural Networks and Learning Systems (2020), pp. 1–14.
- [24] Yin, B., Corradi, F., and Bohte, S. M. Accurate and efficient time-domain classification with adaptive spiking recurrent neural networks. 2021. DOI: 10.48550/ARXIV.2103.12593.
- [25] Funk, S. RMSprop loses to SMORMS3 Beware the Epsilon! https://sifter.org/ simon/journal/20150420.html. Accessed: 2022-4-20. 2015.
- [26] Kingma, D. P. and Ba, J. "Adam: A Method for Stochastic Optimization". en. In: arXiv:1412.6980 (Jan. 2017). arXiv:1412.6980 [cs].
- [27] Turrigiano, G. G. and Nelson, S. B. "Homeostatic plasticity in the developing nervous system". en. In: *Nature reviews. Neuroscience* 5.2 (Feb. 2004), pp. 97–107.
- [28] Gjorgjieva, J., Evers, J. F., and Eglen, S. J. "Homeostatic Activity-Dependent Tuning of Recurrent Networks for Robust Propagation of Activity". In: *The Journal of Neuroscience* 36.13 (2016), pp. 3722–3734.
- [29] Zenke, F. and Gerstner, W. "Hebbian plasticity requires compensatory processes on multiple timescales". In: *Philosophical Transactions of the Royal Society B* 372.1715 (2017), p. 20160259.
- [30] Srivastava, R. K., Greff, K., and Schmidhuber, J. Highway Networks. 2015. arXiv: 1505. 00387 [cs.LG].

- [31] He, K., Zhang, X., Ren, S., and Sun, J. "Deep residual learning for image recognition". In: Proceedings of the IEEE conference on computer vision and pattern recognition. 2016, pp. 770–778.
- [32] Amir, A., Taba, B., Berg, D., Melano, T., McKinstry, J., Di Nolfo, C., Nayak, T., Andreopoulos, A., Garreau, G., Mendoza, M., Kusnitz, J., Debole, M., Esser, S., Delbruck, T., Flickner, M., and Modha, D. "A Low Power, Fully Event-Based Gesture Recognition System". In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, 2017, pp. 7388–7397.
- [33] Eccles, J. C., Fatt, P., and Koketsu, K. "Cholinergic and inhibitory synapses in a pathway from motor-axon collaterals to motoneurones". In: *The Journal of physiology* 126.3 (1954), pp. 524–562.
- [34] Rupprecht, P. and Friedrich, R. W. "Precise Synaptic Balance in the Zebrafish Homolog of Olfactory Cortex". In: *Neuron* 100.3 (2018), 669–683.e5.
- [35] Spiegel, I., Mardinly, A. R., Gabel, H. W., Bazinet, J. E., Couch, C. H., Tzeng, C. P., Harmin, D. A., and Greenberg, M. E. "Npas4 regulates excitatory-inhibitory balance within neural circuits through cell-type-specific gene programs". In: *Cell* 157.5 (2014), pp. 1216–1229.
- [36] Esser, S. K., Appuswamy, R., Merolla, P., Arthur, J. V., and Modha, D. S. "Backpropagation for Energy-Efficient Neuromorphic Computing". In: Advances in Neural Information Processing Systems. Ed. by Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and Garnett, R. Vol. 28. Curran Associates, Inc., 2015.
- [37] Hunsberger, E. and Eliasmith, C. Training spiking deep networks for neuromorphic hardware. Version 1. 2016. arXiv: 1611.05141 [cs.LG].
- [38] Cao, Y., Chen, Y., and Khosla, D. "Spiking Deep Convolutional Neural Networks for Energy-Efficient Object Recognition". In: *International journal of computer vision* 113.1 (2015), pp. 54–66.
- [39] O'Connor, P., Neil, D., Liu, S.-C., Delbruck, T., and Pfeiffer, M. "Real-time classification and sensor fusion with a spiking deep belief network". In: *Frontiers in neuroscience* 7 (2013), p. 178.
- [40] Bu, T., Ding, J., Yu, Z., and Huang, T. Optimized Potential Initialization for Low-latency Spiking Neural Networks. 2022. arXiv: 2202.01440 [cs.NE].
- [41] Bohte, S. M., Kok, J. N., and La Poutré, H. "Error-backpropagation in temporally encoded networks of spiking neurons". In: *Neurocomputing* 48.1-4 (2002), pp. 17–37.
- [42] Booij, O. and Nguyen, H. "A gradient descent rule for spiking neurons emitting multiple spikes". In: *Information Processing Letters* 95.6 (2005), pp. 552–558.
- [43] Mostafa, H. "Supervised Learning Based on Temporal Coding in Spiking Neural Networks." In: *IEEE transactions on neural networks and learning systems* 29.7 (2018), pp. 3227–3235.
- [44] Kheradpisheh, S. R. and Masquelier, T. "Temporal Backpropagation for Spiking Neural Networks with One Spike per Neuron". In: *International Journal of Neural Systems* 30.06 (2020), p. 2050027.
- [45] Comsa, I. M., Potempa, K., Versari, L., Fischbacher, T., Gesmundo, A., and Alakuijala, J. "Temporal Coding in Spiking Neural Networks with Alpha Synaptic Function". In: 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2020, pp. 8529–8533.
- [46] Zenke, F. and Ganguli, S. "SuperSpike: Supervised Learning in Multilayer Spiking Neural Networks". In: *Neural computation* 30.6 (2018), pp. 1514–1541.

- [47] Davidson, S. and Furber, S. B. "Comparison of Artificial and Spiking Neural Networks on Digital Hardware". In: *Frontiers in Neuroscience* 15 (2021).
- [48] Orchard, G., Jayawant, A., Cohen, G. K., and Thakor, N. "Converting Static Image Datasets to Spiking Neuromorphic Datasets Using Saccades". In: *Frontiers in neuroscience* 9 (2015), p. 437.
- [49] Bellec, G., Salaj, D., Subramoney, A., Legenstein, R., and Maass, W. "Long short-term memory and learning-to-learn in networks of spiking neurons". In: *Proceedings of the 32nd International Conference on Neural Information Processing Systems*. Curran Associates Inc., 2018, pp. 795–805.
- [50] Herranz-Celotti, L. and Rouat, J. Surrogate Gradients Design. 2022. arXiv: 2202.00282 [cs.AI].
- [51] Ding, J., Zhang, J., Yu, Z., and Huang, T. Accelerating Training of Deep Spiking Neural Networks with Parameter Initialization. 2022. URL: https://openreview.net/forum? id=T8BnDXDTcFZ.
- [52] Na, B., Mok, J., Park, S., Lee, D., Choe, H., and Yoon, S. AutoSNN: Towards Energy-Efficient Spiking Neural Networks. 2022. arXiv: 2201.12738 [cs.NE].
- [53] Marblestone, A. H., Wayne, G., and Kording, K. P. "Toward an Integration of Deep Learning and Neuroscience". In: *Frontiers in Computational Neuroscience* 10 (2016).
- [54] Richards, B. A., Lillicrap, T. P., Beaudoin, P., Bengio, Y., Bogacz, R., Christensen, A., Clopath, C., Costa, R. P., Berker, A. d., Ganguli, S., Gillon, C. J., Hafner, D., Kepecs, A., Kriegeskorte, N., Latham, P., Lindsay, G. W., Miller, K. D., Naud, R., Pack, C. C., Poirazi, P., Roelfsema, P., Sacramento, J., Saxe, A., Scellier, B., Schapiro, A. C., Senn, W., Wayne, G., Yamins, D., Zenke, F., Zylberberg, J., Therien, D., and Kording, K. P. "A deep learning framework for neuroscience". In: *Nature Neuroscience* 22.11 (2019), pp. 1761–1770.
- [55] Lillicrap, T. P., Santoro, A., Marris, L., Akerman, C. J., and Hinton, G. "Backpropagation and the brain". In: *Nature Reviews Neuroscience* 21 (2020), pp. 335–346.
- [56] Krizhevsky, A. Learning Multiple Layers of Features from Tiny Images. Tech. rep. 2009.
- [57] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. "PyTorch: An Imperative Style, High-Performance Deep Learning Library". In: Advances in Neural Information Processing Systems 32. Ed. by Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. Curran Associates, Inc., 2019, pp. 8024–8035.
- [58] Buzsáki, G. and Mizuseki, K. "The log-dynamic brain: how skewed distributions affect network operations". In: *Nature Reviews Neuroscience* 15 (2014), pp. 264–278.

**Supplementary Figures** 

![](_page_34_Figure_1.jpeg)

Figure S1. Initialization in the fluctuation-driven regime with non-zero  $\mu_{U}$ . (a) Different fluctuation targets  $\xi$  plotted in the space spanned by the parameters of a non-centered Gaussian weight distribution  $W \sim \mathcal{N}(\mu_W, \sigma_W)$ . The red region indicates the regime of meandriven initialization, where  $\mu_U > \theta$ . The border between fluctuation- and mean-driven regime is dependent on data and network parameters. (b) Expected and observed distributions of the membrane potential for different values of the target membrane potential mean  $\mu_U$ . All three displayed initializations lead to fluctuations of the same magnitude  $\xi = 2$ . (c) The standard deviation of the weights  $\sigma_W$  as a function of the target  $\mu_U$  for different fluctuation targets  $\xi$ .

#### 35

![](_page_35_Figure_1.jpeg)

Figure S2. Activity of shallow SNNs at time of initialization in the fluctuationdriven regime. (a) Snapshot of activity over time before training on the Randman dataset for an SNN with one hidden layer. Bottom: spike raster of input layer activity from two different samples corresponding to two different classes. Middle: Spike raster of hidden layer activity. Top: Membrane potential of readout units. The readout units corresponding to the two input classes are highlighted in different shades. The network was initialized with a target  $\sigma_U = 1$ . (b) Distribution of the coefficient of variation (CV) of inter-spike intervals (ISIs) of hidden layer neurons. (c) Standard deviation of the population firing rate filtered with a time constant of 5 ms. The histogram depicts the distribution of  $\sigma_{\text{Rate}}$  across different input samples. (d) Theoretically expected (Supplementary Material S2) and numerically observed proportion of mean-driven neurons at the time of initialization as a function of the target fluctuation magnitude  $\sigma_U$ , for networks that are initialized to be trained on the Randman dataset (±1 standard deviation). The sand-colored shaded region indicates the target regime  $1/3 \leq \sigma_U \leq 1$ . (e)-(h) Same as panels (a)-(d), for the SHD dataset.

![](_page_36_Figure_1.jpeg)

Figure S3. Population-level variability induced by random sampling of synaptic weights. (a) Poisson input spike trains. (b) Membrane potentials  $u_i(t)$  of three example neurons that were initialized with the same target  $\mu_U = 0$  and  $\sigma_U = 1/2$ . (c) Corresponding distributions of the membrane potentials  $U_i$  for each of the three example neurons in panel (B). The black dashed line indicates the target membrane potential distribution  $U \sim \mathcal{N}(\mu_U, \sigma_U^2)$ . Note that the observed means  $\hat{\mu}_{U_i}$  of the three membrane potential distributions deviate from the target. (d) The analytically expected and numerically observed distribution of  $\hat{\mu}_U$  follows a Gaussian. Numerical simulations consider 5000 postsynaptic neurons initialized with the same target  $\mu_U = 0$  and  $\sigma_U = 1/2$ . Even when the target is set clearly in the fluctuationdriven regime, a proportion of neurons can be expected to be mean-driven. (e) The observed membrane potential variances  $\hat{\sigma}_U^2$  follow a Gamma distribution. (f) The observed standard deviations of the membrane potential  $\hat{\sigma}_U$  are Nakagami distributed. For the derivations of analytical solutions in panels (d)-(e), see Supplementary Material S2.

![](_page_37_Figure_1.jpeg)

Figure S4. Sparsity of spiking activity before and after training. (a) Population firing rate of SNNs with one hidden layer before training (black) and after training with (blue) and without (green) regularization of the population firing rate. The shaded region around the lines indicates the standard deviation across five random seeds. The sand-colored shaded region corresponds to our suggested target fluctuation magnitude  $\frac{1}{3} \leq \sigma_U \leq 1$ . The horizontal dashed line at 10 Hz indicates the upper bound imposed by the activity regularization. (b) Distribution of firing rates of single neurons in an example network initialized with  $\sigma_U = 1$  before training (black) and after training with (blue) and without (green) activity regularization through an upper bound on the population firing rate. The dashed line indicates the position of the firing rate regularizer at 10 Hz. (c) Test error of the five best-performing networks trained with and without the firing rate regularizer. (d) As panel (a), but depicting average population firing rates in each hidden layer of CSNNs with three hidden layers. (e) As panel (c), for CSNNs with three hidden layers.

![](_page_38_Figure_1.jpeg)

Figure S5. Fluctuation-driven initialization of recurrent SNNs. (a) Validation accuracy after training as a function of the target fluctuation magnitude  $\sigma_U$  at initialization for SNNs with one hidden layer featuring only forward or additional recurrent connections. The shaded region around the lines indicates the range of values across five random seeds. The sand-colored shaded region corresponds to our suggested target fluctuation magnitude  $\frac{1}{3} \leq \sigma_U \leq 1$ . (b) Population firing rate of hidden layer neurons at the time of initialization as a function of  $\sigma_U$ . Recurrent connections cause the firing rate to increase through a positive feedback loop when the initial fluctuation magnitude is large. (c) As panel (a), for recurrent SNNs initialized with different relative magnitudes of recurrent connections to membrane potential fluctuations. (d) As panel (b), for for recurrent SNNs initialized with different values of  $\alpha$ .

![](_page_38_Figure_3.jpeg)

Figure S6. Deep feed-forward CSNNs are sensitive to initialization. Validation accuracy as a function of target membrane potential fluctuation strength  $\sigma_U$  for strictly feed-forward CSNNs of increasing depth. All networks were trained on the SHD dataset. The shaded region around the lines indicates the range of values across five random seeds. The sand-colored shaded region corresponds to our suggested target fluctuation magnitude  $\frac{1}{3} \leq \sigma_U \leq 1$ . The dashed line indicates Kaiming initialization.

![](_page_39_Figure_1.jpeg)

Figure S7. Fluctuation-driven initialization accelerates learning. (a) Training accuracy as a function of training epoch for CSNNs with a single hidden layer trained on the SHD dataset. Networks were initialized with different target fluctuation magnitudes  $\sigma_U$ . The dashed line indicates 90% training accuracy. (b) As panel (a), for CSNNs with three hidden layers. (c) Number of required epochs to reach 90% training accuracy on the SHD dataset as a function of target fluctuation magnitude  $\sigma_U$  for CSNNs with a single hidden layer. Triangular markers correspond to the values of  $\sigma_U$  plotted in panel (a) and the sand-colored region corresponds to our suggested target fluctuation magnitude  $\frac{1}{3} \leq \sigma_U \leq 1$ . (d) As panel (c), for CSNNs with three hidden layers. In all panels, shaded regions indicate  $\pm 1$  standard deviation across 5 random initializations.

![](_page_40_Figure_1.jpeg)

Figure S8. Re-scaled surrogate gradients can prevent vanishing gradients at the cost of exploding gradients. (a) Population firing rate at time of initialization as a function of hidden layers in a CSNN with seven hidden layers, for different target fluctuation magnitudes  $\sigma_U$ . The shaded region around the lines indicates one standard deviation across five random seeds. (b) The magnitude of surrogate gradients  $\partial L/\partial \tilde{s}$  as a function of hidden layer. The network is being trained with a re-scaled version of the SuperSpike non-linearity that ensures propagation of gradients even in the absence of spikes (see Methods). For initializations with  $\sigma_U \geq 1$ , the gradients explode. (c) As panel (a), but displaying the magnitude of weight changes  $|\partial L/\partial W|$ . When neurons in the previous layer are quiescent, the weight update equals zero. (d)-(f) As panels (a)-(c), for a CSNN without recurrent connections in hidden layers.

![](_page_41_Figure_1.jpeg)

Figure S9. Robustness to weight initialization is sensitive to the choice of optimizer. (a) Validation accuracy after training as a function of the target fluctuation magnitude  $\sigma_U$  at initialization for feed-forward CSNNs with three hidden layers trained on the SHD dataset. Networks were either trained with SGD with learning rate  $\eta$  or with the Adam or SMORMS3 optimizers (see Methods). The shaded region around the lines indicates one standard deviation across five random seeds. (b) As Panel (a), for recurrently connected CSNNs with three hidden layers. (c) Validation loss as a function of training epochs for the best-performing feed-forward networks of each optimization scheme plotted in panel (a). (d) As panel (c), for the recurrently connected CSNNs plotted in panel (b).

![](_page_41_Figure_3.jpeg)

Figure S10. Skip connections can increase robustness to initialization in deep CSNNs. (a) Illustration of skip connections in deep CSNNs. Readout units receive inputs from every hidden layer separately. (b) Validation accuracy after training CSNNs with three hidden layers either with (colored line) or without skip connections on the SHD dataset. The shaded region around the lines indicates the range of values across five random seeds. (c) As panel (b), but the colored line corresponds to networks with both skip connections and homeostatic plasticity. (d) Test accuracy of the 5 best-performing models in terms of validation accuracy for models trained with skip connections, combined skip connections and homeostatic plasticity and the baseline model.

![](_page_42_Figure_1.jpeg)

Figure S11. Parameterization of fluctuation-driven spiking as an initialization strategy for Dalian SNNs. Incoming presynaptic Poisson spike trains (i) from separate excitatory (red) and inhibitory (blue) populations are weighted by the respective synaptic strengths  $W^E$ and  $W^I$  and filtered through the respective PSP kernels  $\epsilon_E(t)$  and  $\epsilon_I(t)$  (ii) to yield membrane potential fluctuations u(t) in a postsynaptic neuron (iii). As in the non-Dalian case, the magnitude of membrane potential fluctuations,  $\sigma_U$ , is determined by the parameters of the presynaptic weight distributions,  $\lambda^E$  and  $\lambda^I$ . Synaptic weights can thus be initialized from a target  $\sigma_U$  (see Methods).

## Supplementary Material

#### S1 PSP kernel parameters

**Current-based synapses.** The parameters  $\bar{\epsilon}$  and  $\hat{\epsilon}$  characterize the integral of the PSP kernel  $\epsilon(t)$ , and the integral of the squared PSP kernel  $\epsilon(t)^2$ , respectively:

$$\bar{\epsilon} = \int_{-\infty}^{\infty} \epsilon(s) ds \tag{57}$$

$$\hat{\epsilon} = \int_{-\infty}^{\infty} \epsilon(s)^2 ds .$$
(58)

To arrive at analytical expressions for  $\bar{\epsilon}$  and  $\hat{\epsilon}$ , we first derive the kernel  $\epsilon(t)$  from the differential equations of a LIF neuron [1]

$$\tau_{\rm mem} \frac{du(t)}{dt} = -u(t) + I(t) \tag{59}$$

$$\frac{dI(t)}{dt} = -\frac{I(t)}{\tau_{\rm syn}} + \sum_{j} w_j S_j(t) , \qquad (60)$$

where  $\tau_{\text{mem}}$  and  $\tau_{\text{syn}}$  are the membrane and synaptic time constants,  $S_j(t)$  are the input spike trains from the presynaptic neuron j weighted by the synaptic weight  $w_j$ , u(t) is the membrane potential and I(t) is the current.

To obtain an expression for the kernel  $\epsilon(t)$ , we consider that neuron *i* receives a single spike from one presynaptic neuron *j* at time t = 0 with a synaptic weight of  $w_j = 1$ . In this case, the synaptic current is simply

$$I(t) = \exp\left(-\frac{t}{\tau_{\rm syn}}\right),\tag{61}$$

which can be inserted into equation (59) to obtain an explicit solution for the membrane potential

$$u(t) = \frac{1}{1 - \frac{\tau_{\text{mem}}}{\tau_{\text{syn}}}} \left( \exp\left(-\frac{t}{\tau_{\text{syn}}}\right) - \exp\left(-\frac{t}{\tau_{\text{mem}}}\right) \right) \Theta(t) .$$
 (62)

Eq. (62) corresponds to the PSP-kernel  $\epsilon(t)$  evoked by a single presynaptic spike, provided that the membrane potential stays in the sub-threshold regime. Note that in the limit of  $\tau_{\rm mem} \to \tau_{\rm syn}$ , the membrane potential follows a scaled alpha function

$$\lim_{\tau_{\rm mem}\to\tau_{\rm syn}} u(t) = \frac{t}{\tau_{\rm syn}} \exp\left(-\frac{t}{\tau_{\rm mem}}\right) \Theta(t) .$$
(63)

We can now solve the integrals in Eqs. (57) and (58) that define the kernel parameters  $\bar{\epsilon}$  and  $\hat{\epsilon}$ , starting with  $\tau_{\text{mem}} \neq \tau_{\text{syn}}$ :

$$\bar{\epsilon} = \int_{-\infty}^{\infty} \epsilon(t) dt$$

$$= \int_{0}^{\infty} \frac{1}{1 - \frac{\tau_{\text{mem}}}{\tau_{\text{syn}}}} \left( \exp(-\frac{t}{\tau_{\text{syn}}}) - \exp(-\frac{t}{\tau_{\text{mem}}}) \right) dt$$

$$= \frac{1}{1 - \frac{\tau_{\text{mem}}}{\tau_{\text{syn}}}} \left( -\tau_{\text{syn}} \exp\left(-\frac{t}{\tau_{\text{syn}}}\right) + \tau_{\text{mem}} \exp\left(-\frac{t}{\tau_{\text{mem}}}\right) \right) \Big|_{0}^{\infty}$$

$$= \tau_{\text{syn}}$$
(64)

$$\hat{\epsilon} = \int_{-\infty}^{\infty} \epsilon^{2}(t) dt$$

$$= \int_{0}^{\infty} \left[ \frac{1}{1 - \frac{\tau_{\text{mem}}}{\tau_{\text{syn}}}} \left( \exp\left(-\frac{t}{\tau_{\text{syn}}}\right) - \exp\left(-\frac{t}{\tau_{\text{mem}}}\right) \right) \right]^{2} dt$$

$$= \left( \frac{1}{1 - \frac{\tau_{\text{mem}}}{\tau_{\text{syn}}}} \right)^{2} \left( \frac{2\tau_{\text{mem}}\tau_{\text{syn}} \exp\left(-\frac{t}{\tau_{\text{syn}}} - \frac{t}{\tau_{\text{mem}}}\right)}{\tau_{\text{syn}} + \tau_{\text{mem}}} \frac{\tau_{\text{syn}} \exp\left(-2\frac{t}{\tau_{\text{syn}}}\right)}{2} - \frac{\tau_{\text{mem}} \exp\left(-2\frac{t}{\tau_{\text{mem}}}\right)}{2} \right) \Big|_{0}^{\infty}$$

$$= \frac{\tau_{\text{syn}}^{2}}{2(\tau_{\text{syn}} + \tau_{\text{mem}})}.$$
(65)

When solving for  $\bar{\epsilon}$  and  $\hat{\epsilon}$  in the case of  $\tau_{\text{mem}} = \tau_{\text{syn}}$  (taking the  $\lim_{\tau_{\text{mem}}\to\tau_{\text{syn}}}$  and applying de L'Hospital's rule), one will arrive at the same solutions as in the above case (see Tab. S1).

**Delta-Synapses.** For reasons of simplicity, current-based synaptic transmission is often replaced with 'delta synapses' in SNNs. In this case, the membrane potential dynamics in response to a single input spike can be described as a mono-exponential decay

$$u(t) = \exp\left(-\frac{t}{\tau_{\rm mem}}\right) \tag{66}$$

and the kernel parameters simplify to

$$\bar{\epsilon} = \int_{-\infty}^{\infty} \epsilon(t) dt = \int_{0}^{\infty} \exp\left(-\frac{t}{\tau_{\rm mem}}\right) dt = \tau_{\rm mem}$$
(67)

$$\hat{\epsilon} = \int_{-\infty}^{\infty} \epsilon^2(t) dt = \int_0^{\infty} \left[ \exp\left(-\frac{t}{\tau_{\rm mem}}\right) \right]^2 dt = \frac{\tau_{\rm mem}}{2} .$$
(68)

A summary of analytical expressions for  $\bar{\epsilon}$  and  $\hat{\epsilon}$  can be found in Tab. S1.

Numerical estimation of  $\bar{\epsilon}$  and  $\hat{\epsilon}$ . Note that the analytical solutions summarized in Tab. S1 might not reflect the values of  $\bar{\epsilon}$  and  $\hat{\epsilon}$  obtained in numerical simulations of the SNNs employing the neuronal dynamics. Specifically, numerical simulations of SNNs often operate with a large time step and integrate neuronal dynamics using simple forward Euler approaches. A more accurate approach with regards to the fluctuations in numerical simulations would therefore be to solve the integrals using the same numerical approximation as employed during the simulations. To quantify the degree to which the analytical solutions for  $\bar{\epsilon}$  and  $\hat{\epsilon}$  deviate from the numerical approximations of the integrals, we numerically approximated the integrals with different values for the simulation time step. Indeed, if the simulation time step was large ( $\geq 1 \text{ ms}$ ), the numerical solutions in this paper, we calculated  $\bar{\epsilon}$  and  $\hat{\epsilon}$  through numerical approximation (forward Euler) using the same time step as during SNN training. The numerical values for  $\bar{\epsilon}$  and  $\hat{\epsilon}$  used throughout our numerical simulations can be found in Tab. 7 of the main text.

![](_page_45_Figure_1.jpeg)

Figure S12. Numerical calculation of  $\bar{\epsilon}$  and  $\hat{\epsilon}$ . (a) Difference (in % of the analytical solution) between numerically calculated and analytically calculated  $\bar{\epsilon}$  and  $\hat{\epsilon}$  as a function of the simulation time step. (b) The resulting difference in the standard deviation of synaptic weights  $\sigma_W$  using fluctuation-driven initialization.

	Delta synapses	Current-based synapses
$\overline{\epsilon}$	$ au_{ m mem}$	$ au_{ m syn}$
$\hat{\epsilon}$	$rac{ au_{ ext{mem}}}{2}$	$rac{ au_{ m syn}^2}{2( au_{ m syn}+ au_{ m mem})}$

**Table S1.** Analytical expressions for  $\bar{\epsilon}$  and  $\hat{\epsilon}$  for LIF neurons with delta synapses or currentbased synapses.

#### S2 Population-level variability in membrane potential fluctuations

The initialization strategy discussed in this article is based on target values for the membrane potential mean  $\mu_U$  and its standard deviation  $\sigma_U$ . Due to the inherent stochasticity arising from random sampling of synaptic weights, we can expect deviations from these targets in the numerically observed membrane potential with mean  $\hat{\mu}_U$  and  $\hat{\sigma}_U$ . As illustrated in Figs. 2 and S2 in the main text, this variability can cause some neurons to fire in the mean-driven regime, even if the initialization targets are set in the fluctuation-driven regime.

Here, we analyze the expected variability of membrane potential fluctuations across a population of m postsynaptic neurons with independently drawn weight vectors  $\{\vec{w}_1, \vec{w}_2, ... \vec{w}_m\}$ . Specifically, we are interested in the sampling distributions of  $\hat{\mu}_U$  and  $\hat{\sigma}_U^2$ . For simplicity, we ignore the spiking dynamics of the LIF neuron model and assume that synaptic weights are independently drawn from the zero-mean Normal distribution  $W \sim \mathcal{N}(0, \sigma_W^2)$ .

Sampling distribution of  $\mu_W$  and  $\sigma_W$ . In order to derive the sampling distributions of the membrane potentials, we first need to derive the sampling distribution of the synaptic weights. That is, for independently drawn weight vectors  $\{\vec{w}_1, \vec{w}_2, ..., \vec{w}_m\}$  of size n, we are interested in the distributions of the sample mean  $\hat{\mu}_W$  and the sample variance  $\hat{\sigma}_W$ . For weights drawn from the zero-mean distribution  $W \sim \mathcal{N}(0, \sigma_W^2)$ , the former is simply

$$\hat{\mu}_W \sim \mathcal{N}\left(0, \frac{\sigma_W^2}{n}\right) \ . \tag{69}$$

To obtain the sampling distribution of the variance, we first observe that, according to Cochran's Theorem [2],

$$n\frac{\hat{\sigma}_W^2}{\sigma_W^2} \sim \chi_{n-1}^2 , \qquad (70)$$

which we can alternatively express as a special case of a Gamma distribution  $\Gamma(k, \theta)$  with shape parameter  $k = \frac{n-1}{2}$  and scale parameter  $\theta = 2$ . Hence the above equation can be written as

$$n\frac{\hat{\sigma}_W^2}{\sigma_W^2} \sim \Gamma\left(\frac{n-1}{2}, 2\right) \ . \tag{71}$$

Using the scaling property of the gamma function, we can solve Eq. (71) for  $\hat{\sigma}_W^2$  under the requirement  $\frac{\sigma_W^2}{n} > 0$ . This holds as long as we have variance in the weight initialization. Hence we can express the distribution of the sample variance  $\hat{\sigma}_W^2$  as

$$\hat{\sigma}_W^2 \sim \Gamma\left(\frac{n-1}{2}, \frac{2\sigma_W^2}{n}\right)$$
 (72)

**Sampling distribution of**  $\hat{\mu}_U$ . We start by observing that the membrane potential U is normally distributed according to the Central Limit Theorem and its mean  $\mu_U$  and variance  $\sigma_U^2$  were given in Eqs. (4) and (5) in the main text (see Fig. 1). To derive the sampling distribution of  $\hat{\mu}_U$ , we observe that the sample  $\hat{\mu}_U$  is related to  $\hat{\mu}_W$  through

$$\hat{\mu}_U = n\nu\bar{\epsilon}\hat{\mu}_W . \tag{73}$$

From the above equation and our derivation of the sampling distribution of  $\hat{\mu}_W$ , it becomes apparent that

$$\hat{\mu}_U \sim \mathcal{N}\left(0, n^2 \nu^2 \bar{\epsilon}^2 \left(\frac{\sigma_W^2}{n}\right)\right) , \qquad (74)$$

which can be further simplified to

$$\hat{\mu}_U \sim \mathcal{N}\left(0, \sigma_U^2 \frac{\nu \bar{\epsilon}^2}{\hat{\epsilon}}\right) \tag{75}$$

by expressing it in terms of the initialization target  $\sigma_U$ , for which we inserted Eq. (9) from the main text (Fig. S3a-d).

We can conclude that random sampling of the weights induces systematic variance in the expected membrane potential means  $\hat{\mu}_U$  of neurons receiving inputs from *n* homogeneous Poisson processes with firing rate  $\nu$ . Specifically, the expected variance on the population level is independent of the number of inputs *n*, but scales with the target fluctuation magnitude  $\sigma_U$  and the input firing rate  $\nu$ .

**Sampling distribution of**  $\hat{\sigma}_U$ . We follow a similar approach to derive the sampling distribution of fluctuation magnitudes in the population. Formally, we are looking for the distribution of  $\hat{\sigma}_U^2$ , which can be derived from observing that for a neuron *i* 

$$\left(\hat{\sigma}_{U}^{(i)}\right)^{2} = n\left(\left(\hat{\sigma}_{W}^{(i)}\right)^{2} + \left(\hat{\mu}_{W}^{(i)}\right)^{2}\right)\nu\hat{\epsilon} .$$

$$(76)$$

To get the distribution of  $\hat{\sigma}_U$ , we first need to determine the distribution of the right hand side. Starting with the distribution of  $\hat{\mu}_W^2$ , we observe that the standardized form of  $\hat{\mu}_W^2$  follows a chi-square distribution with one degree of freedom:

$$\frac{\hat{\mu}_W^2}{\frac{\sigma_W^2}{n}} \sim \chi_1^2 . \tag{77}$$

Similarly to the first paragraph, we can rewrite the chi-square distribution as a Gamma distribution and use its scaling properties to obtain the distribution of  $\hat{\mu}_W^2$ :

$$\hat{\mu}_W^2 \sim \Gamma\left(\frac{1}{2}, \frac{2\sigma_W^2}{n}\right) \ . \tag{78}$$

Note that both  $\hat{\mu}_W^2$  and  $\hat{\sigma}_W^2$  are Gamma distributed with a shared scale parameter  $\theta = \frac{2\sigma_W^2}{n}$  and that these random variable are independent. We can therefore use the summation property of the Gamma distribution to obtain

$$\left(\hat{\sigma}_W^2 + \hat{\mu}_W^2\right) \sim \Gamma\left(\frac{n}{2}, \frac{2\sigma_W^2}{n}\right) \tag{79}$$

and finally plug this result back into Eq. (76) to obtain the distribution

$$\hat{\sigma}_U^2 \sim \Gamma\left(\frac{n}{2}, 2\nu\hat{\epsilon}\sigma_W^2\right) \ . \tag{80}$$

As we did in the previous paragraph, we can insert the solution for  $\sigma_W^2$  in the case of centered weights, given by Eq. (9), to simplify the distribution to

$$\hat{\sigma}_U^2 \sim \Gamma\left(\frac{n}{2}, \frac{2\sigma_U^2}{n}\right) \tag{81}$$

as displayed in Fig. S3f. We can alternatively express the expected variability as the distribution of standard deviations, which follows the Nakagami distribution [3]

$$\hat{\sigma}_U \sim \text{Nakagami}\left(\frac{n}{2}, \sigma_U^2\right)$$
(82)

with shape parameter  $m = \frac{n}{2}$  and spread parameter  $\Omega = \sigma_U^2$  (Fig. S3f). Thus, random sampling of synaptic weights induces a systematic variance in  $\sigma_U$  that scales with  $\sigma_U^2$  and inversely with the number of inputs.

# S3 Fluctuation-driven initialization of Dalian networks using log-normally distributed weights

The initialization of the weights in separate excitatory and inhibitory neuronal populations, as is the case with Dalian networks, requires weights sampled from one-sided distributions. Inspired by neurobiological evidence [4], we consider weights sampled from the log-normal distribution parameterized by  $\mu$  and  $\sigma > 0$  that gives rise to a random variable  $X = e^{\mu + \sigma Z}$  where  $Z \sim N(0, 1)$ is a standard normal random variable. The expected value and the variance of X are then defined as

$$\mathbb{E}[X] = \exp\left(\mu + \frac{\sigma^2}{2}\right) \tag{83}$$

$$\mathbb{V}[X] = \left(\exp\left(\sigma^2\right) - 1\right)\exp\left(2\mu + \sigma^2\right) . \tag{84}$$

To parameterize synaptic weights, we use the log-normal distribution to obtain excitatory weights  $\ln(W^E) \sim \mathcal{N}(\mu_E, \sigma_E)$  and inhibitory weights  $\ln(W^I) \sim \mathcal{N}(\mu_I, \sigma_I)$ . As we only have two equations, which we use to describe the membrane potential of a neuron in the fluctuationdriven regime (mean and variance of the membrane potential), we also need to restrict the parameterization to two parameters in total. Hence, we set  $\sigma_E = \sigma_I = 1$ . We can therefore simplify the above to

$$\mathbb{E}[W] = e^{\mu + 1/2} \tag{85}$$

$$\mathbb{V}[W] = e^{2\mu+2} - e^{2\mu+1} \tag{86}$$

and also note that the second moment of the distribution is

$$\mathbb{E}[W^2] = e^{2\mu + 2} . (87)$$

Given these definitions, we can write down the values of  $\mu_U$  and  $\sigma_U^2$  as

$$\mu_{U} = N_{E}\nu_{E}\bar{\epsilon}_{E}e^{(\mu_{E}+1/2)} - N_{I}\nu_{I}\bar{\epsilon}_{I}e^{(\mu_{I}+1/2)}$$
(88)

$$\sigma_U^2 = N_E \nu_E \hat{\epsilon}_E e^{(2\mu_E + 2)} + N_I \nu_I \hat{\epsilon}_I e^{(2\mu_I + 2)} .$$
(89)

We once more set a target mean membrane potential  $\mu_U = 0$  to achieve a balanced state at initialization. Using this, we solve Eq. (88) to receive an expression for  $\mu_I$ 

$$\mu_I = \mu_E + \log\left(\frac{N_E \nu_E \bar{\epsilon}_E}{N_I \nu_I \bar{\epsilon}_I}\right) \quad , \tag{90}$$

which can be further simplified to

$$\mu_I = \mu_E + \log\left(\frac{1}{\Delta_{EI}}\right) \tag{91}$$

by using the definition of  $\Delta_{EI}$  from Eq. (45). Substituting this result into Eq. (89) allows us to solve for  $\mu_E$  and we receive

$$\mu_E = \frac{1}{2} \log \left( \frac{\sigma_U^2}{N_E \nu_E \hat{\epsilon}_E + N_I \nu_I \hat{\epsilon}_I \left(\frac{1}{\Delta_{EI}}\right)^2} \right) - 1 .$$
(92)

Finally, Equations (91) and (92) together with  $\sigma_E = \sigma_I = 1$  and  $\mu_U = 0$  provide us the parameters to initialize the inhibitory and excitatory weights sampled from a log normal distribution. **Dalian networks with excitatory recurrence.** We start from the mean and variance of the membrane potential for a neuron receiving a feed-forward excitatory, recurrent excitatory and recurrent inhibitory input, where we assume  $\nu_R = \nu_F = \nu_I = \nu$ ,

$$\mu_U = (N_F \nu \bar{\epsilon}_E) e^{\mu_F + 1/2} + (N_R \nu \bar{\epsilon}_E) e^{\mu_R + 1/2} - (N_I \nu \bar{\epsilon}_I) e^{\mu_I + 1/2}$$
(93)

$$\sigma_U^2 = (N_F \nu \hat{\epsilon}_E) e^{2\mu_F + 2} + (N_R \nu \hat{\epsilon}_E) e^{2\mu_R + 2} + (N_I \nu \hat{\epsilon}_I) e^{2\mu_I + 2}$$
(94)

and the definition of  $\alpha$ 

$$\alpha = \frac{\text{Part of } \sigma_U^2 \text{ caused by excitatory feed-forward connections}}{\text{Part of } \sigma_U^2 \text{ caused by all excitatory connections}}$$
(95)

$$= \frac{(N_F \nu \hat{\epsilon}_E) e^{2\mu_F + 2}}{(N_F \nu \hat{\epsilon}_E) e^{2\mu_F + 2} + (N_R \nu \hat{\epsilon}_E) e^{2\mu_R + 2}}.$$
(96)

Here we are again requiring a mean membrane potential  $\mu_U = 0$  and we set the variances of the log-normal distributions to one,  $\sigma_R = \sigma_F = \sigma_I = 1$ .

Performing the same sequence of solving and substituting as in the above paragraph, we find explicit equations for the three weight distribution parameters:

$$\mu_{R} = \mu_{F} + \frac{1}{2}\log(N_{F} - \alpha N_{F}) - \log(\alpha N_{R}) = \mu_{F} + \Delta_{R}$$
(97)

$$\mu_I = \mu_F + \frac{1}{2} \log \left( \frac{\bar{\epsilon}_E \left( e^{\Delta_R} N_R + N_F \right)}{N_I \bar{\epsilon}_I} \right) = \mu_F + \Delta_{EI}^R \tag{98}$$

$$\mu_F = \frac{1}{2} \log \left( \frac{\sigma_U^2}{\nu \left( e^{2\Delta_R} N_R \hat{\epsilon}_E + e^{2\Delta_{EI}^R} \hat{\epsilon}_I N_I + N_F \hat{\epsilon}_E \right)} \right) - 1 .$$
(99)

## References

- [1] Gerstner, W. and Kistler, W. M. *Spiking Neuron Models*. Cambridge University Press, 2002.
- "The distribution of quadratic forms in a normal system, with applications to the analysis of covariance". In: Mathematical Proceedings of the Cambridge Philosophical Society 30.2 (1934), pp. 178–191.
- [3] Huang, L.-F. "The Nakagami and its related distributions". In: WSEAS Transactions on Mathematics 15 (2016), pp. 477–485.
- [4] Buzsáki, G. and Mizuseki, K. "The log-dynamic brain: how skewed distributions affect network operations". In: *Nature Reviews Neuroscience* 15 (2014), pp. 264–278.