

---

# Appendices

## A MORE TRAINING DETAILS

We collect the training hyper-parameters in the following table.

	sl-MNIST	SHD	PTB
batch size	256	256	32
weight decay	0.0	0.1	0.1
gradient norm	1.0	1.0	1.0
train/val/test	45k/5k/10k	8k/1k/2k	930k/74k/82k
	samples	samples	words
learning rate	$3.16 \cdot 10^{-4}$	$3.16 \cdot 10^{-4}$	$3.16 \cdot 10^{-5}$
layers width	128, 128	256, 256	1700, 300
label smoothing	0.1	0.1	0.1
time step repeat	2	2	2
SELT factor	0.8	0.436	4.595

The learning rates were chosen after a grid search fixing dampening and sharpness to 1. The learning rates considered are in the set  $\{10^{-2}, 3.16 \cdot 10^{-3}, 10^{-3}, 3.16 \cdot 10^{-4}, 10^{-4}, 3.16 \cdot 10^{-5}, 10^{-5}\}$ . The results of the grid search are reported in figure 2. The learning rate chosen for the rest of the paper was the one that made all the shapes perform reasonably well, rectangular included. This mostly resulted in a suboptimal learning rate only for the derivative of the fast sigmoid, which still out-performed the rest in the sl-MNIST and SHD, and performed comparatively on the PTB.

We train with crossentropy loss, the AdaBelief optimizer (Zhuang et al., 2020), Stochastic Weight Averaging (Izmailov et al., 2018) and Decoupled Weight Decay (Loshchilov and Hutter, 2019). For the BiGamma distribution, we choose  $\alpha = 5$  and  $\beta = \sqrt{\alpha(\alpha + 1)} = 5.47$  to have a variance of 1. For the PTB task, the input passes through an embedding layer before passing to the first layer, and the output of the last layer is multiplied by the embedding to produce the output, removing the need for the readout (Woźniak et al., 2020; Radford et al., 2018).

Notice that we do not implement forced refractory periods that would prevent the neuron from firing too fast, as sometimes done in the neuromorphic literature, since we want to reduce the non differentiable steps in the system. Thus,  $p = 1$  is possible if the inputs are strong and frequent enough.

## B NEURON MODEL COMPLEXITY

The energy consumed per layer can be used as a metric of neuron complexity, as done in (Yin et al., 2021; Hunger, 2005).

Neural model	Energy (Complexity)
<b>LIF</b>	$(mnp_{l-1} + np_l)E_{AC} + nE_{MAC}$
<b>ALIF</b>	$(mnp_{l-1} + np_l + 2np_l)E_{AC} + 3nE_{MAC}$
<b>LSTM</b>	$4(mn + nn)E_{MAC} + 17nE_{MAC}$
<b>sLSTM</b>	$4(mnp_{l-1} + np_l)E_{AC} + 3np_lE_{AC}$

Table 1: **Neuron complexity.** We use the energy consumed per layer as a metric of neuron complexity (Yin et al., 2021; Hunger, 2005). We use  $n = n_l$  and  $m = n_{l-1}$  as the width of the layer and its input,  $p_l$  for the firing rate of the layer  $l$ .  $E_{MAC}$  is the energy cost of a multiply-accumulate operation and  $E_{AC}$  of an accumulate operation. As shown, *ALIF* always results in a larger number of operations and energy consumption than *LIF*. For large networks,  $n, m \gg 1$ , the square terms dominates, and the *sLSTM* results in 4 times more energy consumption.

## C INTERPLAY BETWEEN SG AND INITIALIZATION

We show how each initialization has a different preferred SG, and viceversa, how each SG has a different preferred initialization in figure 6. Best mean across SG is achieved by the Orthogonal Normal initialization. Best mean across initializations is achieved by the derivative of the fast-sigmoid and the exponential SG.

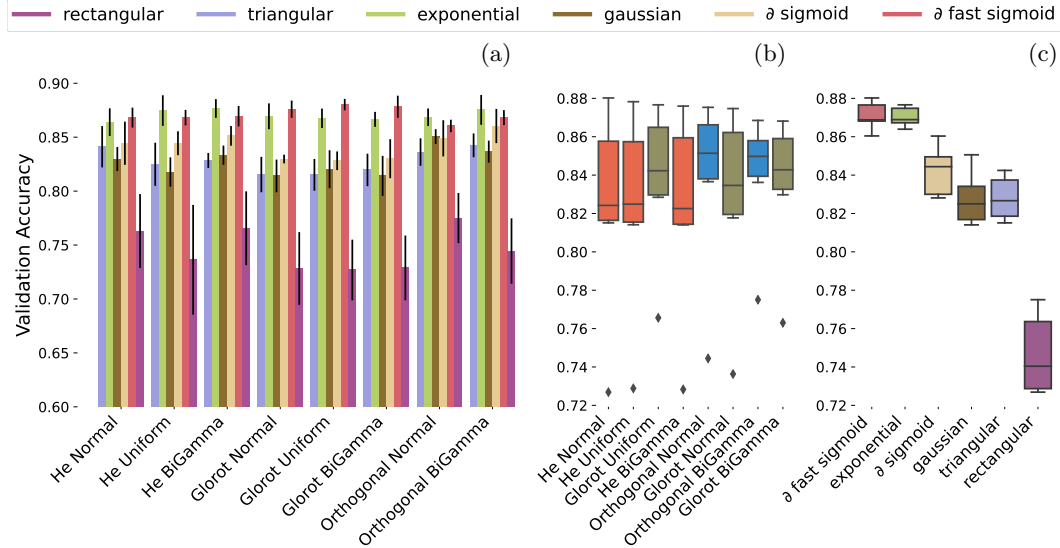


Figure 6: **Orthogonal initialization leads to higher accuracy on the LIF network.** Results for a LIF network trained on the SHD task. (a) Best accuracy when different SG trainings are initialized with different schemes. (b) Results aggregated across SG shapes. Overall, *He* initialization, in orange, achieves the best extreme values, but the *Orthogonal* achieves the best mean, in blue. The BiGamma distribution decreases the accuracy variance. (c) Results aggregated across initializations. Both, the derivative of the fast sigmoid and the exponential SG achieve the best mean accuracy.

## D LIST OF SURROGATE GRADIENTS SHAPES

We list here the shapes that we used in this article as surrogate gradients.

---








	SG name	$f(v)$
	<b>triangular</b>	$\max(1 -  v , 0)$
	<b>exponential</b>	$e^{-2 v }$
	<b>gaussian</b>	$e^{-\pi v^2}$
	<b><math>\partial</math> sigmoid</b>	$4 \operatorname{sigmoid}(4v) (1 - \operatorname{sigmoid}(4v))$
	<b><math>\partial</math> fast-sigmoid</b>	$\frac{1}{(1+ 2v )^2}$
	<b>rectangular</b>	$\mathbb{1}_{ v  < \frac{1}{2}}$
	<b><math>q</math>-PseudoSpike</b> ( $q > 1$ )	$\frac{1}{(1 + \frac{2}{q-1} v )^q}$

Table 2: **Mathematical definitions of the surrogate gradients studied in this article.** Our Heaviside activation  $\sigma(v) = \tilde{H}(v)$ , where  $v$  is the centered voltage, has the SG  $\sigma'(v) = \gamma f(\beta \cdot v)$ , where  $\beta$  is the SG sharpness,  $\gamma$  the SG dampening, and  $f$  is the shape of choice. The constants, are chosen for the SG to have a maximal value of 1 and an area under the curve of 1.

## E DETAILED DERIVATION OF THE CONDITIONS

We derive the constraints on the hyper-parameters that will lead the LIF to meet the conditions proposed at initialization. The LIF we will be using is defined by

$$\mathbf{y}_t = \alpha_{decay} \mathbf{y}_{t-1} (1 - \mathbf{x}_{t-1}) + \mathbf{i}_t \quad (1)$$

where  $\mathbf{i}_t = W_{rec} \mathbf{x}_{t-1} + W_{in} \mathbf{z}_t + \mathbf{b}$ , as described in the main text, and the multiplicative factor  $(1 - \mathbf{x}_{t-1})$  represents the reset mechanism.

### E.1 RECURRENT MATRIX MEAN SETS THE FIRING RATE (I)

We show how condition (I) leads to a constraint on the mean of the recurrent connectivity with a LIF neuron model, that will lead the network to meet that condition at initialization.

**Lemma 1.** *Applying condition (I), which states that we want  $\operatorname{Median}[v] = 0$ , to an LIF network, and further assuming  $\bar{w}_{in} = 0, b = 0$ , the approximation  $\operatorname{Mean}[v] \approx \operatorname{Median}[v]$ , and constant  $\bar{i}_t$  over time, it results in the constraint*

$$\bar{w}_{rec} = \frac{1}{n_{rec} - 1} (2 - \bar{\alpha}_{decay}) \bar{\vartheta} \quad (2)$$

*Proof.* First we show that  $\operatorname{Median}[v] = 0 \implies \operatorname{Mean}[x] = 1/2$ , where  $x = \tilde{H}(v)$ . In equation 4 we write the marginal distribution of  $p(x) = \int p(x|v)p(v)dv$ , and the double integral is represented with one integration symbol. Then, we notice that  $x$  has a deterministic dependence on  $v$ ,  $x = H(v)$ , which probabilistically is described by the delta function  $p(x|v) = \delta(x - H(v))$ . Then, we integrate over  $x$ , and in the last equation we notice that integrating with respect to the Heaviside is equivalent to restricting the integration limits from zero to infinity.

$$Mean[x] = \int xp(x)dx \quad (3)$$

$$= \int xp(x|v)p(v)dx dv \quad (4)$$

$$= \int xp(v)\delta(x - H(v))dx dv \quad (5)$$

$$= \int p(v)dvH(v) \quad (6)$$

$$= \int_0^\infty p(v)dv \quad (7)$$

If  $Median[v] = 0$ , half of it's probability mass is on each side of 0, so the last integral is equal to  $1/2$ , QED.

Since working with medians is mathematically harder than working with means, we assume that  $Mean[v] \approx Median[v]$ , with the caveat that it will make the result approximate. To justify that they are similar, it can be shown that for a unimodal distribution  $v \sim p(v)$  with the first two moments defined, we have  $|Mean[v] - Median[v]| \leq \sqrt{0.6Var[v]}$

We use the notation  $\bar{x} = Mean[x]$  interchangeably. We calculate how the mean of the voltage elements is propagated through time, assuming the mean input current to remain constant over time  $\bar{i}_t = \bar{i}$  at initialization, to simplify the mathematical development, and assuming per condition (I), that  $\bar{x} = 1 - x = 1/2$  we have

$$\bar{y}_t = \bar{\alpha}_{decay} \overline{(1 - \mathbf{x}_{t-1})} \bar{y}_{t-1} + \bar{i} \quad (8)$$

$$= \frac{1}{2} \bar{\alpha}_{decay} \bar{y}_{t-1} + \bar{i} \quad (9)$$

$$= \frac{1}{2} \bar{\alpha}_{decay} \left( \frac{1}{2} \bar{\alpha}_{decay} \bar{y}_{t-2} + \bar{i} \right) + \bar{i} \quad (10)$$

$$= \frac{1}{2^{t-1}} \bar{\alpha}_{decay}^{t-1} \bar{y}_1 + \left( \sum_{t'=0}^{t-2} \frac{1}{2^{t'}} \bar{\alpha}_{decay}^{t'} \right) \bar{i} \quad (11)$$

$$= \frac{1}{2^{t-1}} \bar{\alpha}_{decay}^{t-1} \bar{y}_1 + \frac{1 - \frac{1}{2^{t-1}} \bar{\alpha}_{decay}^{t-1}}{1 - \frac{1}{2} \bar{\alpha}_{decay}} \bar{i} \quad (12)$$

where we used the fact that the same LIF definition applies to different time steps, the geometric series formula, and the fact that for independent random variables  $E[XY] = E[X]E[Y]$ . For  $t \rightarrow \infty$  and using  $0 < \bar{\alpha}_{decay} < 1$

$$\bar{y}_t = \frac{1}{1 - \frac{1}{2} \bar{\alpha}_{decay}} \bar{i} \quad (13)$$

$$\bar{y}_t - \bar{\vartheta} = \frac{1}{1 - \frac{1}{2} \bar{\alpha}_{decay}} \bar{i} - \bar{\vartheta} \quad (14)$$

Assuming we want this condition to hold independently of the dataset, we set  $Mean[W_{in}] = 0$ , and assuming that we do not want to promote this behavior with fixed internal currents, but with the recurrent activity instead, then  $\mathbf{b} = 0$ .

We remark that we denote  $Mean[W\mathbf{x}]$  as the mean vector whose element  $i$  is

$$Mean[W\mathbf{x}]_i = Mean[\sum_{\substack{j=1 \\ j \neq i}} w_{ij}x_j] \quad (15)$$

$$= \sum_{\substack{j=1 \\ j \neq i}} Mean[w_{ij}x_j] \quad (16)$$

$$= \sum_{\substack{j=1 \\ j \neq i}} Mean[wx] \quad (17)$$

$$= (n_{rec} - 1)Mean[wx] \quad (18)$$

where the condition  $j \neq i$  in the summand reminds that neurons are not connected to themselves in our recurrent architecture. In the first equality, the index  $i$  denoting the element in the vector, is equivalent as choosing the row  $i$  of  $W$ , so it is not necessary to specify it outside the square brackets. The equality before the last one is a consequence of considering any neuron as mutually independent to any other at initialization, as done by

Then,

$$Mean[y_t - \vartheta] = \frac{1}{1 - \frac{1}{2}\bar{\alpha}_{decay}} \left( (n_{rec} - 1)\bar{w}_{rec}\bar{x}_{t-1} \right) - \bar{\vartheta} \quad (19)$$

$$0 = \frac{1}{1 - \frac{1}{2}\bar{\alpha}_{decay}} (n_{rec} - 1)\bar{w}_{rec}\bar{x}_{t-1} - \bar{\vartheta} \quad (20)$$

$$0 = \frac{1}{1 - \frac{1}{2}\bar{\alpha}_{decay}} (n_{rec} - 1)\bar{w}_{rec}\frac{1}{2} - \bar{\vartheta} \quad (21)$$

$$(n_{rec} - 1)\bar{w}_{rec} = (2 - \bar{\alpha}_{decay})\bar{\vartheta} \quad (22)$$

$$\bar{w}_{rec} = \frac{1}{n_{rec} - 1} (2 - \bar{\alpha}_{decay})\bar{\vartheta} \quad (23)$$

where in the second line we applied condition (I) in the form of  $Mean[v_t] \approx Median[v_t] = 0$ , so  $Mean[y_t - \vartheta] = 0$ , and in the third line we applied again condition (I),  $\bar{x}_t = 1/2$ . In the main text we turn  $\bar{\vartheta}, \bar{\alpha}_{decay} \rightarrow \vartheta, \alpha_{decay}$ , since here we consider the more general case where those are as well random variables, and we simplify it in the main text for cleanliness, assuming they are constant.

□

We therefore found a constraint on the mean of the recurrent matrix initialization, that leads the LIF network to satisfy condition I at initialization. The constraint is equation 23 with  $\bar{w}_{in} = 0$ , and  $\mathbf{b} = 0$ .

## E.2 RECURRENT MATRIX VARIANCE CAN MAKE RECURRENT AND INPUT VOLTAGES COMPARABLE (II)

We apply condition (II) to the LIF network, that gives us a constraint that the recurrent matrix has to meet at initialization for the condition to be true.

**Lemma 2.** *Applying condition (II), which states that we want  $Var[W_{rec}x_{t-1}] = Var[W_{in}z_t]$ , to an LIF network, and further assuming  $\bar{x} = 1/2$ , and  $\bar{w}_{in} = 0$ , it results in the constraint*

$$Var[w_{rec}] = 2(Var[z_t] + \bar{z}_t^2) \frac{n_{in}}{n_{rec} - 1} Var[w_{in}] - \frac{1}{2}\bar{w}_{rec}^2 \quad (24)$$

*Proof.* The second condition, is that the recurrent and the input contribution to the variance need to match

$$Var[W_{rec}x_{t-1}] = Var[W_{in}z_t] \quad (25)$$

where the variance is computed at each element, after the matrix multiplication is performed, following the method described in

$$Var[W\mathbf{x}]_i = Var[\sum_{j=1} w_{ij}x_j] \quad (26)$$

$$= \sum_{j=1} Var[w_{ij}x_j] \quad (27)$$

$$= \sum_{j=1} Var[w_{ij}x_j] \quad (28)$$

$$= n_W Var[w_{ij}x_j] \quad (29)$$

The second and third equality are a consequence of considering any neuron as mutually independent to any other at initialization, as done by

Therefore the vector-wise condition II is equivalent to the element-wise

$$(n_{rec} - 1)Var[w_{rec}x_{t-1}] = n_{in}Var[w_{in}z_t] \quad (30)$$

Since the time dimension is averaged out, the time axis can be randomly shuffled, and the LIF activity is indistinguishable from a Bernoulli process through the mean and variance of the activity. Therefore if  $\bar{x}_t = p$ , we have  $Var[x_t] = p(1-p)$  when averaged over time, with  $p$  the probability of firing. Therefore it is as well true that  $\overline{x_t^2} = Var[x_t] + \bar{x}_t^2 = p$ .

We apply the fact that for independent  $w, x$

$$Var[w_{ij}x_j] = \overline{w_{ij}^2} \overline{x_j^2} - \overline{w_{ij}}^2 \bar{x}_j^2 \quad (31)$$

and assuming  $\bar{w}_{in} = 0$  and  $p = 1/2$  we have

$$Var[w_{rec}x_{t-1}] = (Var[w_{rec}] + \bar{w}_{rec}^2)p - \bar{w}_{rec}^2 p^2 \quad (32)$$

$$= \frac{1}{4}(2Var[w_{rec}] + \bar{w}_{rec}^2) \quad (33)$$

$$Var[w_{in}z_t] = (Var[z_t] + \bar{z}_t^2)Var[w_{in}] \quad (34)$$

Substituting in equation 30 implies

$$\frac{1}{4}(2Var[w_{rec}] + \bar{w}_{rec}^2) = (Var[z_t] + \bar{z}_t^2) \frac{n_{in}}{n_{rec} - 1} Var[w_{in}] \quad (35)$$

$$Var[w_{rec}] = 2(Var[z_t] + \bar{z}_t^2) \frac{n_{in}}{n_{rec} - 1} Var[w_{in}] - \frac{1}{2} \bar{w}_{rec}^2 \quad (36)$$

□

Therefore condition (II) led us to the constraint that  $W_{rec}$  has to meet at initialization, equation 36, for the condition to be true. The final equation further assumes that  $\bar{w}_{in} = 0$  and  $p = 1/2$ .

### E.3 SG DAMPENING CONTROLS GRADIENT MAXIMUM (III)

We apply condition (III) to the LIF network, which gives us a constraint that the dampening has to meet at initialization for the condition to be true.

**Lemma 3.** *Applying condition (III), which states that we want  $Max[\frac{\partial}{\partial\theta}y_t] = Max[\frac{\partial}{\partial\theta}y_{t-1}]$ , to an LIF network, and assuming that (1)  $\sigma'$  and  $\frac{\partial}{\partial\theta}y_{t-1}$  are statistically independent and (2) we do not pass the gradient through the reset, it results in the constraint*

$$\gamma = \frac{1}{(n_{rec} - 1)\hat{w}_{rec}} \left( 1 - \hat{\alpha}_{decay} - \xi \cdot n_{in}\hat{w}_{in}\gamma_{in} \right) \quad (37)$$

where  $\xi$  is zero for the first layer and it's one for the other layers in the stack.

*Proof.* We want the maximal value of the gradient to remain stable, without exploding, when transmitted through time and through different layers

$$Max[\frac{\partial}{\partial\theta}y_t] = Max[\frac{\partial}{\partial\theta}y_{t-1}] \quad (38)$$

where when we write  $\partial/\partial\theta$ , we use  $\theta$  as a placeholder for any quantity that we want to propagate through gradient descent. Taking the derivative of the LIF definition and stopping the gradient from going through the reset we have

$$\frac{\partial}{\partial\theta}y_t = \alpha_{decay} \frac{\partial}{\partial\theta}y_{t-1}(1 - x_{t-1}) + W_{rec} \frac{\partial}{\partial\theta}x_{t-1} + \xi W_{in} \frac{\partial}{\partial\theta}z_t \quad (39)$$

Here we introduce the symbol  $\xi \in \{0, 1\}$ , where  $\xi = 1$  is used when  $z_t$  comes from a trainable layer below, and  $\xi = 0$  when  $z_t$  represents the data. We consider as well that

$$\frac{\partial}{\partial\theta}z_t = \frac{\partial}{\partial\theta}\tilde{H}_{in}(y_t^{in} - \vartheta_{in}) = \sigma'_{in} \frac{\partial}{\partial\theta}y_t^{in} \quad (40)$$

$$\frac{\partial}{\partial\theta}x_{t-1} = \frac{\partial}{\partial\theta}\tilde{H}(y_{t-1} - \vartheta) = \sigma' \frac{\partial}{\partial\theta}y_{t-1} \quad (41)$$

where  $\tilde{H}_{in}, y_t^{in}, \vartheta_{in}$  are the Heaviside, the voltage and the threshold of the layer below,  $\sigma' = \frac{\partial\tilde{H}}{\partial v}$  is the surrogate gradient, and  $\sigma'_{in}$  is the surrogate gradient from the layer below. Substituting in equation 39, then

$$\frac{\partial}{\partial\theta}y_t = \alpha_{decay} \frac{\partial}{\partial\theta}y_{t-1}(1 - x_{t-1}) + W_{rec}\sigma' \frac{\partial}{\partial\theta}y_{t-1} \quad (42)$$

$$+ \xi W_{in}\sigma'_{in} \frac{\partial}{\partial\theta}y_t^{in} \quad (43)$$

We use  $Max$  and  $Min$  in a statistical ensemble sense, as the maximum/minimum value that a variable could take if sampled over and over again

$$Max[X] = \sup_{x \sim p(x)} x \quad (44)$$

$$Min[X] = \inf_{x \sim p(x)} x \quad (45)$$

With this definition, if  $X, Y$  are independent random variables  $Max[X + Y] = Max[X] + Max[Y]$  and if they are positive  $Max[XY] = Max[X]Max[Y]$ . We observe, as we did before for the variance and the mean of  $Wx$ , that

$$Max[Wx] = n_W Max[w_x] \quad (46)$$

$$Min[Wx] = n_W Min[w_x] \quad (47)$$

We take the maximal value of  $\frac{\partial}{\partial \theta} y_t$ , we make the assumption that  $\sigma'$  and  $\frac{\partial}{\partial \theta} y_{t-1}$  are statistically independent, we use the fact that the highest value that the surrogate gradient can take is given by the dampening factor  $Max[\sigma'] = \gamma$ , we denote as  $\gamma_{in}$  the dampening factor of the layer below in the stack, and we take  $Max[1 - x_{t-1}] = 1$ :

$$\begin{aligned} Max[\frac{\partial}{\partial \theta} y_t] &= Max[\alpha_{decay} \frac{\partial}{\partial \theta} y_{t-1}] \\ &\quad + (n_{rec} - 1) Max[w_{rec} \sigma' \frac{\partial}{\partial \theta} y_{t-1}] \\ &\quad + \xi n_{in} Max[w_{in} \sigma'_{in} \frac{\partial}{\partial \theta} y_t^{in}] \end{aligned} \quad (48)$$

$$\begin{aligned} &= Max[\alpha_{decay}] Max[\frac{\partial}{\partial \theta} y_{t-1}] \\ &\quad + (n_{rec} - 1) Max[w_{rec}] Max[\sigma'] Max[\frac{\partial}{\partial \theta} y_{t-1}] \\ &\quad + \xi n_{in} Max[w_{in}] Max[\sigma'_{in}] Max[\frac{\partial}{\partial \theta} y_t^{in}] \end{aligned} \quad (49)$$

$$\begin{aligned} &= Max[\alpha_{decay}] Max[\frac{\partial}{\partial \theta} y_{t-1}] \\ &\quad + (n_{rec} - 1) Max[w_{rec}] \gamma Max[\frac{\partial}{\partial \theta} y_{t-1}] \\ &\quad + \xi n_{in} Max[w_{in}] \gamma_{in} Max[\frac{\partial}{\partial \theta} y_t^{in}] \end{aligned} \quad (50)$$

where we used the fact that  $\sigma'$  is positive in the second equality. We apply condition (III), which states that all maximal gradients are equivalent, and for cleanliness we use the notation  $Max[x] = \hat{x}$

$$1 = \hat{\alpha}_{decay} + (n_{rec} - 1) \hat{w}_{rec} \gamma + \xi n_{in} \hat{w}_{in} \gamma_{in} \quad (51)$$

$$\gamma = \frac{1}{(n_{rec} - 1) \hat{w}_{rec}} \left( 1 - \hat{\alpha}_{decay} - \xi \cdot n_{in} \hat{w}_{in} \gamma_{in} \right) \quad (52)$$

where we only had to rearrange terms. □

We set  $\xi = 0$  in the main text for readability and because we observed better performance with it. This final equation 52 gives the value that the dampening has to take to keep the maximal gradient value stable, namely, condition (III) true at initialization.

#### E.4 SG SHARPNESS CONTROLS GRADIENT VARIANCE (IV)

We apply condition (IV) to the LIF network to constrain the choice of surrogate gradient variance.

**Lemma 4.** *Applying condition (IV), which states that we want  $Var[\frac{\partial}{\partial \theta} y_t] = Var[\frac{\partial}{\partial \theta} y_{t-1}]$ , to an LIF network, and assuming that (1) we do not pass the gradient through the reset, and (2) zero mean gradients at initialization, it results in the constraint*



---


$$\overline{\sigma'^2} = \frac{1 - \frac{1}{2}\overline{\alpha^2}_{decay} - \xi \cdot n_{in}\overline{w_{in}^2} \overline{\sigma'^2}_{in}}{(n_{rec} - 1)\overline{w^2}_{rec}} \quad (53)$$

where  $\xi$  is zero for the first layer and is one for the other layers in the stack.

*Proof.* Condition (IV) states that we want the variance of the gradient to remain stable across time and layers. Taking the derivative of the LIF we arrive at equation 43:

$$\begin{aligned} \frac{\partial}{\partial \theta} y_t = & \alpha_{decay} \frac{\partial}{\partial \theta} y_{t-1} (1 - x_{t-1}) \\ & + W_{rec} \sigma' \frac{\partial}{\partial \theta} y_{t-1} + \xi W_{in} \sigma'_{in} \frac{\partial}{\partial \theta} y_t^{in} \end{aligned} \quad (54)$$

Taking the variance and assuming that the monomials in the polynomial are statistically independent, we can consider the variance of the sum to be the sum of the variances:

$$\begin{aligned} Var[\frac{\partial}{\partial \theta} y_t] = & Var[\alpha_{decay} \frac{\partial}{\partial \theta} y_{t-1} (1 - x_{t-1})] \\ & + Var[W_{rec} \sigma' \frac{\partial}{\partial \theta} y_{t-1}] \\ & + Var[\xi W_{in} \sigma'_{in} \frac{\partial}{\partial \theta} y_t^{in}] \end{aligned} \quad (55)$$

$$\begin{aligned} Var[\frac{\partial}{\partial \theta} y_t] = & Var[\alpha_{decay} \frac{\partial}{\partial \theta} y_{t-1} (1 - x_{t-1})] \\ & + (n_{rec} - 1) Var[w_{rec} \sigma' \frac{\partial}{\partial \theta} y_{t-1}] \\ & + n_{in} Var[\xi w_{in} \sigma'_{in} \frac{\partial}{\partial \theta} y_t^{in}] \end{aligned} \quad (56)$$

where  $\xi = 0$  if  $w_{in}$  connects to the data and  $\xi = 1$  if it connects to the layer below in the stack. We denote by  $\sigma'_{in}$  the surrogate gradient of the layer below.

Assuming gradients  $g$  with mean zero, and weights and gradients  $w, g$  to be independent random variables at initialization:

$$Var[wg] = (Var[g] + E[g]^2)(Var[w] + E[w]^2) - E[g]^2 E[w]^2 \quad (57)$$

$$= Var[g](Var[w] + E[w]^2) \quad (58)$$

$$= Var[g]E[w^2] \quad (59)$$

which gives

$$\begin{aligned} Var[\frac{\partial}{\partial \theta} y_t] = & E[\alpha_{decay}^2 (1 - x_{t-1})^2] Var[\frac{\partial}{\partial \theta} y_{t-1}] \\ & + (n_{rec} - 1) E[(w_{rec} \sigma')^2] Var[\frac{\partial}{\partial \theta} y_{t-1}] \\ & + \xi \cdot n_{in} E[(w_{in} \sigma'_{in})^2] Var[\frac{\partial}{\partial \theta} y_t^{in}] \end{aligned} \quad (60)$$

We apply condition IV, we want gradients to have the same variance, irrespective of the time step, or the neuron in the stack, which results in

$$1 = \frac{1}{2}E[\alpha_{decay}^2] + (n_{rec} - 1)E[(w_{rec}\sigma')^2] + \xi \cdot n_{in}E[(w_{in}\sigma'_{in})^2] \quad (61)$$

$$1 = \frac{1}{2}E[\alpha_{decay}^2] + (n_{rec} - 1)E[w_{rec}^2]E[\sigma'^2] + \xi \cdot n_{in}E[w_{in}^2]E[\sigma'_{in}^2] \quad (62)$$

where we used the fact that for independent variables  $X, Y$  we have  $E[X^p Y^q] = E[X^p]E[Y^q]$  in the third and fourth line. Using the notation  $E[x] = \bar{x}$ , the implied condition on the SG is

$$\overline{\sigma'^2} = \frac{1 - \frac{1}{2}\overline{\alpha_{decay}^2} - \xi \cdot n_{in}\overline{w_{in}^2} \overline{\sigma'^2_{in}}}{(n_{rec} - 1)\overline{w_{rec}^2}} \quad (63)$$

□

We therefore found the constraint that the second non-centered moment of the SG has to satisfy, equation 63, if we want condition IV to hold. We set  $\xi = 0$  in the main text for readability and because we observed better performance with it. We show how to relate it to the sharpness of the exponential SG in Appendix E.5.

#### E.5 APPLYING CONDITION IV TO THE EXPONENTIAL SG

We show how we apply equation 63, to choose the sharpness of an exponential SG. For that we need to define the dependence of the variance of the SG with its sharpness. We use as equivalent notation for the surrogate gradient

$$\sigma'(v) = \frac{\partial \tilde{H}(v)}{\partial v} = \gamma f(\beta \cdot v)$$

We denote no dependency with the voltage in  $\sigma'$ , when we consider it as a random variable, and we introduce the dependency  $\sigma'(v)$  when we assume the voltage dependence is known. The moments of the surrogate gradient are given by

$$E[\sigma'^m] = \int \sigma'^m p(\sigma') d\sigma' \quad (64)$$

$$= \iint \sigma'^m p(\sigma'|v) p(v) dv d\sigma' \quad (65)$$

$$= \iint \sigma'^m \delta(\sigma' - \sigma'(v)) d\sigma' p(v) dv \quad (66)$$

$$= \int \sigma'(v)^m p(v) dv \quad (67)$$

where we used the marginalization rule in the second equality and in the third equality we used the fact that  $\sigma$  is a deterministic function of  $v$ , so it inherits its randomness from  $v$ . We are going to assume as the non-informative prior a uniform distribution between the minimal and maximal values of  $y_t - \vartheta$ .

---


$$E[\sigma'(v)^m] = \int \sigma'(v)^m p(v) dv \quad (68)$$

$$= \frac{1}{y_{max} - y_{min}} \int_{y_{min}-\vartheta}^{y_{max}-\vartheta} \sigma'(v)^m dv \quad (69)$$

$$= \frac{\gamma^m}{\beta(y_{max} - y_{min})} \int_{\beta(y_{min}-\vartheta)}^{\beta(y_{max}-\vartheta)} f(v')^m dv' \quad (70)$$

where we used the non informative uniform prior assumption in the second equality and we used  $\sigma' = \gamma f(\beta v)$  followed by the change of variable  $v' = \beta v$  in the third equality. Considering the exponential SG we have that, calling  $v_i$  one of the integration limits above, if  $v_i$  is positive

$$\int_0^{v_i} g(|v|)^m dv = \int_0^{v_i} g(v)^m dv \quad (71)$$

and if  $v_i$  is negative

$$\int_0^{v_i} g(|v|)^m dv = \int_0^{v_i} g(-v)^m dv \quad (72)$$

$$= - \int_0^{-v_i} g(v)^m dv \quad (73)$$

$$= - \int_0^{|v_i|} g(v)^m dv \quad (74)$$

where we made the change of variable  $v \rightarrow -v$  in the second equality. Therefore

$$\int_{v_-}^{v_+} g(|v|)^m dv = \text{sign}(v_+) \int_0^{|v_+|} g(v)^m dv \quad (75)$$

$$- \text{sign}(v_-) \int_0^{|v_-|} g(v)^m dv \quad (76)$$

Given that for  $v_i > 0$  we have

$$\int_0^{v_i} \text{exponential}(v)^m dv = \int_0^{v_i} e^{-2m|v|} dv \quad (77)$$

$$= \int_0^{v_i} e^{-2mv} dv \quad (78)$$

$$= -\frac{1}{2m} e^{-2mv_i} + \frac{1}{2m} \quad (79)$$

$$= -\frac{1}{2m} e^{-2m|v_i|} + \frac{1}{2m} \quad (80)$$

then, for  $v_+ > 0$  and  $v_- < 0$

$$\begin{aligned} \int_{v_-}^{v_+} \text{exponential}(v)^m dv = \\ -\frac{1}{2m} e^{-2m|v_+|} - \frac{1}{2m} e^{-2m|v_-|} + \frac{2}{2m} \end{aligned} \quad (81)$$

where  $v_+ = \beta(y_{max} - \vartheta)$  and  $v_- = \beta(y_{min} - \vartheta)$  and we show how to compute  $y_{max} = \text{Max}[y_t]$  and  $y_{min} = \text{Min}[y_t]$  in section E.7. Since the dependence with  $\beta$  is quite complex, we find the  $\beta$  that satisfies the last equation and equation 63 through random search. Notice how equation 68, shows a dependence of the SG variance proportional to the square of the dampening and inversely proportional to the sharpness, which recalls our numerical results, where a high sharpness and a low dampening were preferred. This is how condition (IV) is used to fix the sharpness of the exponential SG.

#### E.6 APPLYING CONDITION IV TO THE $q$ -PSEUDOSPIKE SG

Instead, when using (IV) to determine the tail-fatness of the SG, we set  $\beta = 1$  and use

$$\int_{v_-}^{v_+} \text{q-PseudoSpike}(|v|)^2 dv = \int_{v_-}^{v_+} \text{q-PseudoSpike}(|v|)^2 dv \quad (82)$$

$$= \int_0^{|v_+|} \text{q-PseudoSpike}(v)^2 dv + \int_0^{|v_-|} \text{q-PseudoSpike}(v)^2 dv \quad (83)$$

$$= -\frac{q+2|v_+|-1}{2(2q-1)} \frac{1}{\left(1 + \frac{2}{q-1}|v_+|\right)^{2q}} - \frac{q+2|v_-|-1}{2(2q-1)} \frac{1}{\left(1 + \frac{2}{q-1}|v_-|\right)^{2q}} + \frac{q-1}{(2q-1)} \quad (84)$$

When inserted in equation 63, we use gradient descent to optimize  $q$  and find the value that satisfies (IV).

#### E.7 MAXIMAL AND MINIMAL VOLTAGE VALUES ACHIEVABLE BY THE NETWORK AT INITIALIZATION

We calculate the maximum and minimum value that the voltage  $y$  can take, to be able to complete the argument about the variance of the backward pass of section E.5. First, we use  $\text{Max}$  and  $\text{Min}$  in a statistical ensemble sense, as the maximum/minimum value that a variable could take if sampled over and over again

$$\text{Max}[X] = \sup_{x \sim p(x)} x \quad (85)$$

$$\text{Min}[X] = \inf_{x \sim p(x)} x \quad (86)$$

When applied to the definition of LIF

$$\begin{aligned} \text{Max}[y_t] &= \text{Max}[\alpha_{decay} y_{t-1} (1 - x_{t-1})] + \text{Max}[W_{rec} x_{t-1}] \\ &\quad + \text{Max}[b] + \text{Max}[W_{in} z_t] \end{aligned} \quad (87)$$

$$\begin{aligned} &= \text{Max}[\alpha_{decay}] \text{Max}[y_{t-1}] + (n_{rec} - 1) \text{Max}[w_{rec}] \\ &\quad + \text{Max}[b] + n_{in} \text{Max}[w_{in}] \end{aligned} \quad (88)$$

$$\begin{aligned} \text{Max}[y_t] &= \frac{1}{1 - \text{Max}[\alpha_{decay}]} \left( (n_{rec} - 1) \text{Max}[w_{rec}] \right. \\ &\quad \left. + \text{Max}[b] + n_{in} \text{Max}[w_{in}] \right) \end{aligned} \quad (89)$$

where we used the fact that if  $x_t, z_t$  were sampled over and over, the maximum value that they could take is all neurons having fired at the same time, we used the fact that

$\alpha_{decay}, \vartheta > 0$ , and we assumed that the maximum is going to stay constant through time  $Max[y_{t-1}] = Max[y_t]$ . Notice that the maximal voltage is achieved when all neurons in the layer fired at  $t - 1$ , equation 88, except for the neuron under study, that stayed silent at  $t - 1$ , to have 89. Similarly for the bound to the minimal voltage:

$$Min[y_t] = Min[\alpha_{decay}y_{t-1}(1 - x_{t-1})] + (n_{rec} - 1)Min[w_{rec}x_{t-1}] + Min[b] \quad (90)$$

$$+ n_{in}Min[w_{in}z_t] \quad (91)$$

$$= Max[\alpha_{decay}]Min[y_{t-1}] + (n_{rec} - 1)Min[w_{rec}] + Min[b] \quad (92)$$

$$+ n_{in}Min[w_{in}] \quad (93)$$

$$Min[y_t] = \frac{1}{1 - Max[\alpha_{decay}]} \left( (n_{rec} - 1)Min[w_{rec}] + Min[b] + n_{in}Min[w_{in}] \right) \quad (94)$$

Second, we consider another definition of  $Max$  and  $Min$ , where we consider the maximum value achievable by the current sample from the weight distribution. The real maximum value of the voltage will be achieved when the presynaptic neurons to fire are those that are connected with positive weight, we then have that our equation turns to

$$Max[y_t]_i = \alpha_{decay,i}Max[y_{t-1}]_i + \sum_j ReLU[W_{rec}]_{ij} + b_i + \sum_j ReLU[W_{in}]_{ij} \quad (95)$$

$$Max[y_t]_i = \frac{1}{1 - \alpha_{decay,i}} \left( \sum_j ReLU[W_{rec}]_{ij} + b_i + \sum_j ReLU[W_{in}]_{ij} \right) \quad (96)$$

where we refer as  $\sum_j ReLU[W_{rec}]_{ij}$  the sum over columns, where we have typically omitted the index  $i$  for the element of the vector for cleanliness in the rest of the article. The case for the minimum is analogous

$$Min[y_t]_i = \alpha_{decay,i}Min[y_{t-1}]_i + Min[W_{rec}x_{t-1}] + b_i + Min[W_{in}z_t] \quad (97)$$

$$Min[y_t]_i = \frac{1}{1 - \alpha_{decay,i}} \left( - \sum_j ReLU[-W_{rec}]_{ij} + b_i - \sum_j ReLU[-W_{in}]_{ij} \right) \quad (98)$$

With this we showed how we calculated the maximal and minimal value of the voltage, to be able to use condition (IV) to define the sharpness of the SG in section E.5.

## F APPLYING CONDITIONS I-IV TO AN ALTERNATIVE DEFINITION OF RESET

We want to show how the constraints on the weights initialization and on the SG choice change, when the neuron model definition changes. We will use the notation  $i_t = W_{rec}x_t + W_{in}z_t + b$ . The reset used by

---


$$y_t = (\alpha_{decay} y_{t-1} + i_t)(1 - x_{t-1}) \quad (99)$$

that we will call *post-reset*. Instead,

$$y_t = \alpha_{decay} y_{t-1}(1 - x_{t-1}) + i_t \quad (100)$$

that we will call *pre-reset*, since it resets before applying the new current. Another example is given by

$$y_t = \alpha_{decay} y_{t-1} + i_t - \vartheta x_{t-1} \quad (101)$$

and we will call it *minus-reset*.

The first definition performs as well one refractory period, while the second does not result in a  $y_t$  clamped to zero when  $x_t = 1$ . The factor  $(1 - x_t)$  takes the voltage exactly to zero every time the neuron has fired, zero being the equilibrium voltage. What is interesting about this form of reset is that the voltage is reset exactly to  $y = 0$  after firing, while with the subtractive reset it is not the case. We consider training without passing the gradient through the reset, since

**Post-reset:**

$$y_t = (\alpha_{decay} y_{t-1} + i_t)(1 - x_{t-1})$$

$$\bar{w}_{rec} = \frac{2}{n_{rec} - 1} (1 - \alpha_{decay}) \vartheta \quad \text{I}$$

$$Var[w_{rec}] = 2(Var[z_t] + \bar{z}_t^2) \frac{n_{in}}{n_{rec} - 1} Var[w_{in}] - \frac{1}{2} \bar{w}_{rec}^2 \quad \text{II}$$

$$\gamma = \frac{1}{(n_{rec} - 1) \hat{w}_{rec}} (1 - \alpha_{decay} - \xi n_{in} \hat{w}_{in} \gamma_{in}) \quad \text{III}$$

$$\overline{\sigma'^2} = \frac{2 - \alpha_{decay}^2 - \xi n_{in} \overline{w_{in}^2} \overline{\sigma_{in}^2}}{(n_{rec} - 1) \bar{w}_{rec}^2} \quad \text{IV}$$

**Pre-reset:**

$$y_t = \alpha_{decay} y_{t-1}(1 - x_{t-1}) + i_t$$

$$\bar{w}_{rec} = \frac{1}{n_{rec} - 1} (2 - \alpha_{decay}) \vartheta \quad \text{I}$$

$$Var[w_{rec}] = 2(Var[z_t] + \bar{z}_t^2) \frac{n_{in}}{n_{rec} - 1} Var[w_{in}] - \frac{1}{2} \bar{w}_{rec}^2 \quad \text{II}$$

$$\gamma = \frac{1}{(n_{rec} - 1) \hat{w}_{rec}} (1 - \alpha_{decay} - \xi n_{in} \hat{w}_{in} \gamma_{in}) \quad \text{III}$$

$$\overline{\sigma'^2} = \frac{1 - \frac{1}{2} \alpha_{decay}^2 - \xi n_{in} \overline{w_{in}^2} \overline{\sigma_{in}^2}}{(n_{rec} - 1) \bar{w}_{rec}^2} \quad \text{IV}$$

**Minus-reset:**

---


$$y_t = \alpha_{decay} y_{t-1} + i_t - \vartheta x_{t-1}$$

$$\bar{w}_{rec} = \frac{1}{n_{rec} - 1} (3 - 2\alpha_{decay}) \vartheta \quad \text{I}$$

$$Var[w_{rec}] = 2(Var[z_t] + \bar{z}_t^2) \frac{n_{in}}{n_{rec} - 1} Var[w_{in}] - \frac{1}{2} \bar{w}_{rec}^2 \quad \text{II}$$

$$\gamma = \frac{1}{(n_{rec} - 1) \bar{w}_{rec} - \vartheta \bar{w}_{rec}} \left( 1 - \alpha_{decay} - \xi n_{in} \bar{w}_{in} \gamma_{in} \right) \quad \text{III}$$

$$\frac{1}{\sigma'^2} = \frac{1 - \alpha_{decay}^2 - \xi n_{in} \bar{w}_{in}^2 \sigma_{in}'^2}{(n_{rec} - 1) \bar{w}_{rec}^2 + \vartheta^2} \quad \text{IV}$$

To have the conditions when the gradient does not pass through the reset, put  $\vartheta = 0$  in (III) and (IV), but not in (I).

## G ALIF AND sLSTM MODELS

To study the variability of SG training with architecture choice, we tested different SG shapes on the ALIF and sLSTM networks. We used the following ALIF implementation

$$\begin{aligned} \mathbf{y}_{t,l} = & \alpha_{decay,l}^y \mathbf{y}_{t-1,l} \\ & + W_{rec,l} \mathbf{x}_{t-1,l} + W_{in,l} \mathbf{x}_{t-1,l-1} + \mathbf{b}_l \\ & - \vartheta_{t-1,l} \mathbf{x}_{t-1,l} \end{aligned} \quad (102)$$

$$\vartheta_{t,l} = \alpha_{decay,l}^\vartheta \vartheta_{t-1,l} + \mathbf{b}_l^\vartheta + \beta_l \mathbf{x}_{t-1,l} \quad (103)$$

where we initialized  $W_{rec}, W_{in}$  as Glorot Uniform,  $b_l = 0$ ,  $\alpha_{decay,l}^y = 4 \cdot 10^{-5}$ ,  $\alpha_{decay,l}^\vartheta = 0.992$  for the SHD task and  $\alpha_{decay,l}^\vartheta = 0.98$  for the sl-MNIST task,  $b_l^\vartheta = 0.01$ , and  $\beta_l = 1.8$ .

The LSTM implementation that we used is the following

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \quad (104)$$

$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \quad (105)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \quad (106)$$

$$\tilde{c}_t = \sigma_c(W_c x_t + U_c h_{t-1} + b_c) \quad (107)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \quad (108)$$

$$h_t = o_t \circ \sigma_h(c_t) \quad (109)$$

The dynamical variables  $i_t, f_t, o_t$  represent the input, forget and output gates, that prevent representations and gradients from exploding, while  $c_t, h_t$  represent the two hidden layers of the LSTM, that work as the working memory and are maintained and updated through data time  $t$ . To construct the spiking version of the LSTM (sLSTM) we turned the activations into  $\sigma_g(x) = H(x)$  and  $\sigma_c = \sigma_h = 2H(x) - 1$ . The matrices  $W_j, U_j$  are initialized with Glorot Uniform initialization, and the biases  $b_j$  as zeros, with  $j \in \{i, f, o, c\}$ .

## H MORE ON SPARSITY

We investigate if the role of sparsity remains consistent across SG shapes in Fig. 7, and across tasks in Fig. 8. Notice that Fig. 3 is repeated in Fig. 7 and 8 to ease the comparison.

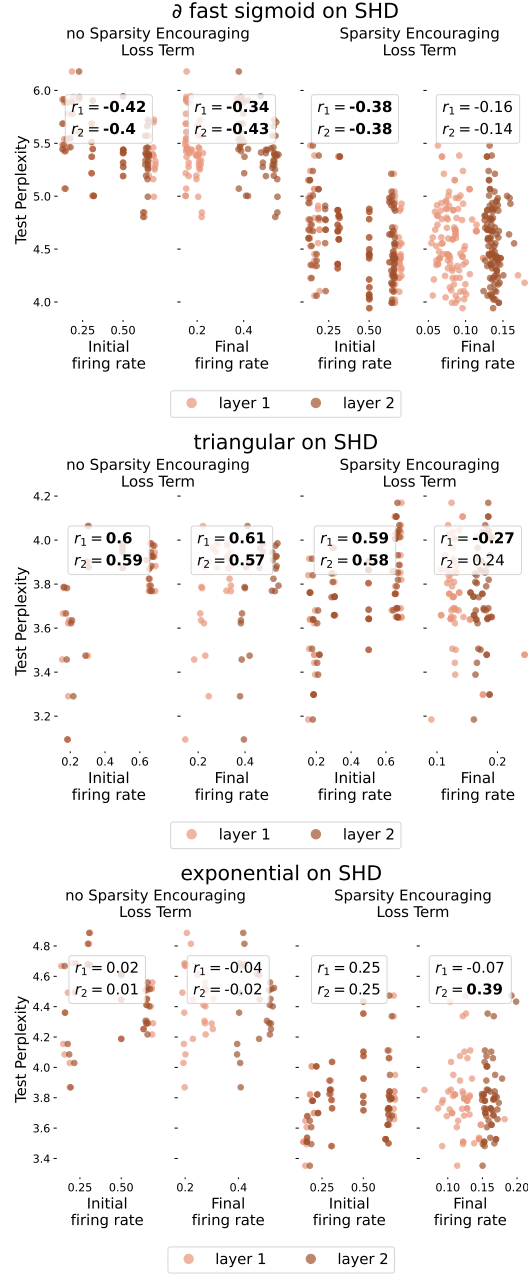


Figure 7: **Sparsity role is not consistent across SG shapes.** When we fix the task to be the SHD task, we see that the derivative of the fast sigmoid has preference for high  $p_i$ , the triangular SG has preference for low  $p_i$ , while for the exponential,  $p_i$  does not seem to correlate with final performance.



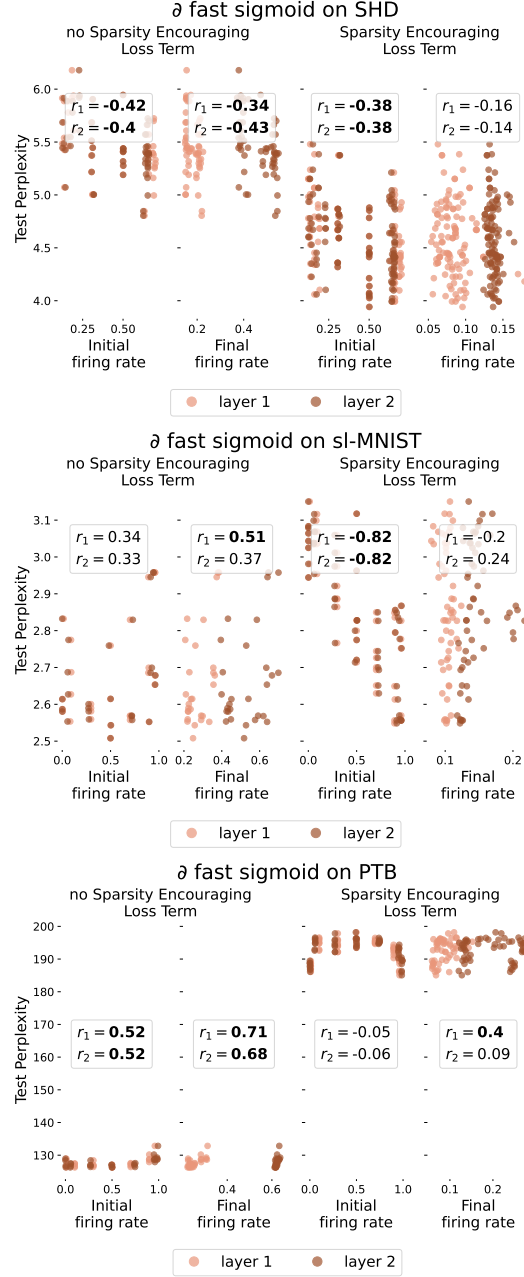


Figure 8: **Sparsity role is consistent across tasks.** Here we fix the SG shape to the derivative of the fast sigmoid and we change the task. On sl-MNIST, we see a similar trend than on SHD, where high initial firing rate is preferred for better performance when sparsity is encouraged. Encouraging sparsity has a negative effect on learning language modeling on the PTB task. However, when no sparsity is encouraged, best performance on PTB is still at  $p_i = 0.5$ .