# Supplementary information "EGRU: Event-based GRU for activity-sparse inference and learning"

**Anonymous Author(s)**
Affiliation
Address
email

## A    Derivation of continuous time version of GRU

In this section we derive the continuous-time version of the GRU model. Note that our definition of the GRU differs from the original version, presented in [3], by inverting the role of the $\mathbf{u}^{\langle t \rangle}$ and $1 - \mathbf{u}^{\langle t \rangle}$ terms in the updates equations for $\mathbf{y}^{\langle t \rangle}$ (a change in sign). This substitution does not change the behavior of the model but simplifies the notation in the continuous-time version of the model.

We first rewrite the dynamics of a layer of GRU units at time step $t$ from Eq. (1) of the main text, separating out the input and recurrent weights:

$$\mathbf{u}^{\langle t \rangle} = \sigma\Big(\mathbf{U}_u \mathbf{x}^{\langle t \rangle} + \mathbf{V}_u \mathbf{y}^{\langle t-1 \rangle} + \mathbf{b}_u\Big), \quad \mathbf{r}^{\langle t \rangle} = \sigma\Big(\mathbf{U}_r \mathbf{x}^{\langle t \rangle} + \mathbf{V}_r \mathbf{y}^{\langle t-1 \rangle} + \mathbf{b}_r\Big),$$

$$\mathbf{z}^{\langle t \rangle} = g\Big(\mathbf{U}_z \mathbf{x}^{\langle t \rangle} + \mathbf{V}_z\Big(\mathbf{r}^{\langle t \rangle} \odot \mathbf{y}^{\langle t-1 \rangle}\Big) + \mathbf{b}_z\Big), \quad \mathbf{y}^{\langle t \rangle} = \mathbf{u}^{\langle t \rangle} \odot \mathbf{z}^{\langle t \rangle} + (1 - \mathbf{u}^{\langle t \rangle}) \odot \mathbf{y}^{\langle t-1 \rangle}, \tag{S1}$$

we can write this as

$$\mathbf{y}^{\langle t \rangle} - \mathbf{y}^{\langle t-1 \rangle} = -\mathbf{u}^{\langle t \rangle} \odot \mathbf{y}^{\langle t-1 \rangle} + \mathbf{u}^{\langle t \rangle} \odot \mathbf{z}^{\langle t \rangle}. \tag{S2}$$

Note that $\mathbf{u}$ here is equivalent to $\tilde{\mathbf{u}} = 1 - \mathbf{u}$ used in the standard GRU model. Eq. (S2) is in the form of a forward Euler discretization of a continuous time dynamical system. Defining $\mathbf{y}(t) \equiv \mathbf{y}^{\langle t-1 \rangle}$, we get $\mathbf{r}(t) \equiv \mathbf{r}^{\langle t \rangle}, \mathbf{u}(t) \equiv \mathbf{u}^{\langle t \rangle}, \mathbf{z}(t) \equiv \mathbf{z}^{\langle t \rangle}$. Let $\Delta t$ define an arbitrary time step. Then Eq. (S2) becomes:

$$\mathbf{y}(t + \Delta t) - \mathbf{y}(t) = -\mathbf{u}(t) \odot (\mathbf{y}(t) + \mathbf{z}(t)) \Delta t \tag{S3}$$

Dividing by $\Delta t$ and taking limit $\Delta t \to 0$, we get:

$$\dot{\mathbf{y}}(t) = -\mathbf{u}(t) \odot (\mathbf{y}(t) - \mathbf{z}(t)), \tag{S4}$$

where $\dot{\mathbf{y}}(t) \equiv \frac{d\mathbf{y}(t)}{dt}$ is the time derivative of $\mathbf{y}(t)$.

## B    Full details of the continuous time EGRU

In this section we establish the continuous time version of the EGRU model. To describe the event generating mechanism and state dynamics it is convenient to express the dynamical system equations in therms of the activations $\mathbf{a}_\mathrm{x}$.

We first rewrite Eqs. (3) & (4) of the main text, as:

$$f_{a_\mathrm{x}} \quad \equiv \quad \tau_s \dot{\mathbf{a}}_\mathrm{x} + \mathbf{a}_\mathrm{x} + \mathbf{b}_\mathrm{x} = \mathbf{0}, \quad \mathrm{x} \in \{u, r, z\} \tag{S5}$$

$$f_c \quad \equiv \quad \tau_m \dot{\mathbf{c}}(t) + \mathbf{u}(t) \odot (\mathbf{c}(t) - \mathbf{z}(t)) = \tau_m \dot{\mathbf{c}}(t) - F(t, \mathbf{a}_u, \mathbf{a}_r, \mathbf{a}_z, \mathbf{c}) = 0. \tag{S6}$$

We write the event transitions for $\mathbf{c}$ at network event $e_k \in \mathbf{e}$, $e_k = (s_k, n_k)$, where $s_k$ are the continuous (real-valued) event times, and $n_k$ denotes which unit got activated, and using the superscript $.^-$ $(.^+)$ to the quantity just before (after) the event, as:

$$c_{n_k}^-(s_k) = \vartheta_{n_k}, \qquad c_{n_k}^+(s_k) = 0, \qquad c_m^+(s_k) = c_m^-(s_k). \tag{S7}$$

22  where $m \neq n_k$ denotes all the units connected to unit $n_k$ that are not activated. At the time of this event,
23  the activations $a_{\mathrm{x},m}$ ($\mathrm{X} \in \{u,r,x\}$) experiences a jump in its state value, given by:

$$a_{u,m}^+(s_k) = a_{u,m}^-(s_k) + v_{u,mn_k} \times c_{n_k}^-(s_k),$$ (S8)

$$a_{r,m}^+(s_k) = a_{r,m}^-(s_k) + v_{r,mn_k} \times c_{n_k}^-(s_k),$$ (S9)

$$a_{z,m}^+(s_k) = a_{z,m}^-(s_k) + v_{z,mn_k} \times r_{n_k} \times c_{n_k}^-(s_k),$$ (S10)

$$a_{\mathrm{X},n_k}^+(s_k) = a_{\mathrm{X},n_k}^-(s_k).$$ (S11)

24  External inputs also come in as events $\tilde{e}_k \in \tilde{\mathbf{e}}$, $\tilde{e}_k = (s_k, i_k)$, where $s_k$ are the continuous (real-valued)
25  event times, and $i_k$ denotes the index of the input component that got activated. Only the activations
26  $a_{\mathrm{X},l}$ for the $l$-th unit experience a transition/jump on incoming external input events, as follows:

$$a_{\mathrm{X},l}^+(s_k) = a_{\mathrm{X},l}^-(s_k) + u_{\mathrm{X},ln_k} \times x_{i_k}(s_k),$$ (S12)

27  where $x_{i_k}(s_k) = (\mathbf{x}(s_k))_{i_k}$ is the $i_k$-th component of the input $\mathbf{x}$ at time $s_k$. The internal state $\mathbf{c}$ remains
28  the same on the external input event. That is, $c_l^+ = c_l^-$.

## C   Derivation of event-based learning rule for EGRU

30  In this section we derive the event-based updates for the network weights. The update questions yield
31  different results for the recurrent weights ($\mathbf{V}_\mathrm{X}$), biases ($\mathbf{b}_\mathrm{X}$) and input weights ($\mathbf{U}_\mathrm{X}$), which are derived
32  in the remainder of this section. To increase readability important terms are highlighted in color.

### C.1   Gradient updates for the recurrent weights $\mathbf{V_x}$

34  We first split the integral Eq. (7) across events as:

$$\mathcal{L} = \sum_{k=0}^N \int_{s_k}^{s_{k+1}} \left[ \ell_c(\mathbf{c}(t),t) + \boldsymbol{\lambda}_c \cdot f_c + \sum_{\mathrm{X} \in \{u,r,z\}} \boldsymbol{\lambda}_{a_\mathrm{x}} \cdot f_{a_\mathrm{x}} \right] dt.$$ (S13)

35  Then taking the derivative of the full loss function, we get:

$$\frac{d\mathcal{L}}{dv_{ji}} = \frac{d}{dv_{ji}} \left\{ \sum_{k=0}^N \int_{s_k}^{s_{k+1}} \left[ \ell_c(\mathbf{c}(t),t) + \boldsymbol{\lambda}_c \cdot f_c + \sum_{\mathrm{X} \in \{u,r,z\}} \boldsymbol{\lambda}_{a_\mathrm{x}} \cdot f_{a_\mathrm{x}} \right] dt \right\}.$$ (S14)

36  By application of Leibniz integral rule we get,

$$\frac{d}{dv_{ji}} \int_{s_k}^{s_{k+1}} \ell_c(\mathbf{c}(t),t) dt = \ell_c(\mathbf{c},s_{k+1}) \frac{ds_{k+1}}{dv_{ji}} - \ell_c(\mathbf{c},s_k) \frac{ds_k}{dv_{ji}} + \int_{s_k}^{s_{k+1}} \frac{\partial \ell_c}{\partial \mathbf{c}} \cdot \frac{\partial \mathbf{c}}{\partial v_{ji}} dt.$$ (S15)

37  and

$$\frac{d}{dv_{ji}} \int_{s_k}^{s_{k+1}} \boldsymbol{\lambda}_c \cdot f_c dt$$ (S16)

$$= \int_{s_k}^{s_{k+1}} \boldsymbol{\lambda}_c \cdot \frac{df_c}{dv_{ji}} dt = \int_{s_k}^{s_{k+1}} \boldsymbol{\lambda}_c \cdot \left\{ \tau_m \frac{d}{dt} \frac{\partial \mathbf{c}}{\partial v_{ji}} + \frac{\partial F}{\partial v_{ji}} \right\} dt$$ (S17)

$$= \tau_m \left[ \boldsymbol{\lambda}_c \cdot \frac{\partial \mathbf{c}}{\partial v_{ji}} \right]_{s_k}^{s_{k+1}}$$ (S18)

$$- \tau_m \int_{s_k}^{s_{k+1}} \left\{ \dot{\boldsymbol{\lambda}}_c \cdot \frac{\partial \mathbf{c}}{\partial v_{ji}} + \boldsymbol{\lambda}_c \cdot \left( \left( \frac{\partial F}{\partial \mathbf{c}} \right)^T \frac{\partial \mathbf{c}}{\partial v_{ji}} + \sum_{\mathrm{X} \in \{u,r,z\}} \left( \frac{\partial F}{\partial \mathbf{a}_\mathrm{X}} \right)^T \frac{\partial \mathbf{a}_\mathrm{X}}{\partial v_{ji}} \right) \right\} dt,$$ (S19)

38  where we first apply Gronwall's theorem [5], then integration by parts, and $M^T$ denotes the transpose
39  of matrix $M$. $\ell_c(\mathbf{c}(t),t)$ is the instantaneous loss evaluated at time $t$. Similarly,

$$\frac{d}{dv_{ji}} \int_{s_k}^{s_{k+1}} \sum_{\mathrm{X} \in \{u,r,z\}} \boldsymbol{\lambda}_{a_\mathrm{x}} \cdot f_{a_\mathrm{x}} dt = \sum_{\mathrm{X} \in \{u,r,z\}} \int_{s_k}^{s_{k+1}} \boldsymbol{\lambda}_{a_\mathrm{x}} \cdot \left\{ \tau_s \frac{d}{dt} \frac{\partial \mathbf{a}_\mathrm{X}}{\partial v_{ji}} + \frac{\partial \mathbf{a}_\mathrm{X}}{\partial v_{ji}} \right\} dt$$ (S20)

$$= \tau_s \left[ \boldsymbol{\lambda}_{a_\mathrm{x}} \cdot \frac{\partial \mathbf{a}_\mathrm{X}}{\partial v_{ji}} \right]_{s_k}^{s_{k+1}} - \tau_s \int_{s_k}^{s_{k+1}} \left\{ \dot{\boldsymbol{\lambda}}_{a_\mathrm{x}} \cdot \frac{\partial \mathbf{a}_\mathrm{X}}{\partial v_{ji}} + \boldsymbol{\lambda}_{a_\mathrm{x}} \cdot \frac{\partial \mathbf{a}_\mathrm{X}}{\partial v_{ji}} \right\} dt,$$ (S21)

2

40    since $\frac{\partial \mathbf{b}}{\partial v_{ji}} = 0$.

41    Substituting these values into Eq. (S14), and setting the coefficients of terms with $\frac{\partial \mathbf{c}}{\partial v_{ji}}$ and $\frac{\partial \mathbf{a_x}}{\partial v_{ji}}$ to zero

42    (using the fact that we can choose the adjoint variables freely due to $f_c$ and $f_{a_x}$ being everywhere zero

43    by definition), we get the dynamics of the adjoint variable described in Eq. (8). The adjoint variable

44    is usually integrated backwards in time starting from $t = T$, also due to its dependence on the loss

45    values (See Fig. S2). The initial conditions for the adjoint variables is defined as $\boldsymbol{\lambda}_c = \boldsymbol{\lambda}_{a_x} = \mathbf{0}$.

46    Setting the coefficients of terms with $\frac{\partial \mathbf{c}}{\partial v_{ji}}$ and $\frac{\partial \mathbf{a_x}}{\partial v_{ji}}$ to zero allows us to write the parameter updates as:

$$\frac{d\mathcal{L}}{dv_{ji}} = \sum_{k=0}^{N} \left\{ \left(l_c^- - l_c^+\right) \frac{ds}{dv_{ji}} + \tau_s \sum_{\mathrm{X}} \left( \boldsymbol{\lambda}_{a_x}^- \cdot \frac{\partial \mathbf{a_x^-}}{\partial v_{ji}} - \boldsymbol{\lambda}_{a_x}^+ \cdot \frac{\partial \mathbf{a_x^+}}{\partial v_{ji}} \right) + \tau_m \left( \boldsymbol{\lambda}_c^- \cdot \frac{\partial \mathbf{c^-}}{\partial v_{ji}} - \boldsymbol{\lambda}_c^+ \cdot \frac{\partial \mathbf{c^+}}{\partial v_{ji}} \right) \right\} \tag{S22}$$

$$= \sum_{k=0}^{N} \xi_{\mathrm{X},ijk} \tag{S23}$$

47    To define the required jumps at event times for the adjoint variables, we start with finding the

48    relationship between $\frac{\partial \mathbf{c^-}}{\partial v_{ji}}$ and $\frac{\partial \mathbf{c^+}}{\partial v_{ji}}$. Eqs. (S7) define $s_k$ as a differentiable function of $v_{ji}$ under the

49    condition $\dot{c}_{n_k}^- \neq 0$ and $\dot{c}_{n_k}^+ \neq 0$ due to the implicit function theorem [16, 17].

$$c_{n_k}^- - \vartheta_{n_k} = 0 \tag{S24}$$

$$\frac{\partial c_{n_k}^-}{\partial v_{ji}} + \frac{dc_{n_k}^-}{ds} \frac{\partial s}{\partial v_{ji}} = 0 \tag{S25}$$

$$\frac{\partial c_{n_k}^-}{\partial v_{ji}} + \dot{c}_{n_k}^- \frac{\partial s}{\partial v_{ji}} = 0 \tag{S26}$$

$$\frac{\partial s}{\partial v_{ji}} = \frac{-1}{\dot{c}_{n_k}^-} \frac{\partial c_{n_k}^-}{\partial v_{ji}}, \tag{S27}$$

50    where we write $\frac{dc_{n_k}^-}{ds} \equiv \dot{c}_{n_k}^-$ and $\dot{c}_{n_k}^- \neq 0$. Similarly,

$$c_{n_k}^+ = 0 \tag{S28}$$

$$\frac{\partial c_{n_k}^+}{\partial v_{ji}} + \dot{c}_{n_k}^+ \frac{\partial s}{\partial v_{ji}} = 0 \tag{S29}$$

51    which allows us to write

$$\frac{\partial c_{n_k}^+}{\partial v_{ji}} = \frac{\dot{c}_{n_k}^+}{\dot{c}_{n_k}^-} \frac{\partial c_{n_k}^-}{\partial v_{ji}} \tag{S30}$$

52    Similarly, starting from $c_m^+ = c_m^-$, we can derive

$$\frac{\partial c_m^+}{\partial v_{ji}} = \frac{\partial c_m^-}{\partial v_{ji}} - \frac{1}{\dot{c}_{n_k}^-} \frac{\partial c_{n_k}^-}{\partial v_{ji}} \left( \dot{c}_m^- - \dot{c}_m^+ \right) \tag{S31}$$

53    For the activations $\mathbf{a_x}$, we use Eqs. (S8)–(S11) to derive the relationships between $\frac{\partial a_x}{\partial v_{ji}}^+$ and $\frac{\partial a_x}{\partial v_{ji}}^-$.

54    Thus, we have:

$$\frac{\partial a_{\mathrm{X},m}^+}{\partial v_{ji}} = \frac{\partial a_{\mathrm{X},m}^-}{\partial v_{ji}} - \frac{1}{\tau_s} \frac{v_{mn_k} r_{\mathrm{X},n_k}^- c_{n_k}^-}{\dot{c}_{n_k}^-} \frac{\partial c_{n_k}^-}{\partial v_{ji}} + \delta_{in_k} \delta_{jm} c_{n_k}^- + c_{n_k}^- v_{mn_k} \frac{\partial r_{\mathrm{X},n_k}^-}{\partial v_{ji}} - c_{n_k}^- v_{mn} \frac{\dot{r}_{\mathrm{X},n_k}^-}{\dot{c}_{n_k}^-} \frac{\partial c_{n_k}^-}{\partial v_{ji}} \tag{S32}$$

$$\frac{\partial a_{\mathrm{X},n_k}^+}{\partial v_{ji}} = \frac{\partial a_{\mathrm{X},n_k}^-}{\partial v_{ji}} \tag{S33}$$

55    where $\mathbf{r_X} = \mathbf{0}$ if $\mathrm{X} \in \{u, r\}$ and $\mathbf{r_X} = \mathbf{r}$ if $\mathrm{X} = \{z\}$.

3

56  Substituting Eqs. (S30),(S31),(S33), (S32) into Eq. (S22), we get:

$$\xi_{\text{x},ijk} = \left\{ \frac{\partial c_{n_k}^-}{\partial v_{ji}} \left( \frac{-1}{\dot{c}_{n_k}^-}(\ell_c^+ - \ell_c^-) + \tau_m \left( \lambda_{c,n_k}^- - \frac{\dot{c}_{n_k}^+}{\dot{c}_{n_k}^-}\lambda_{c,n_k}^+ \right) + \tau_m \frac{1}{\dot{c}_{n_k}^-}\sum_{m \neq n_k} \lambda_{c,m}^+ (\dot{c}_m^- - \dot{c}_m^+) \right. \right. \tag{S34}$$

$$\left. + \sum_{\text{x}} \frac{r_{\text{x},n_k}^- c_{n_k}^-}{\dot{c}_{n_k}^-} \sum_{m \neq n_k} v_{mn}\lambda_{a_{\text{x}},m}^+ + \tau_s \sum_{\text{x}} \frac{\dot{r}_{\text{x},n_k}^- c_{n_k}^-}{\dot{c}_{n_k}^-} \sum_{m \neq n_k} v_{mn}\lambda_{a_{\text{x}},m}^+ \right) \tag{S35}$$

$$+ \tau_m \sum_{m \neq n_k} \frac{\partial c_m^-}{\partial v_{ji}} \left( \lambda_{c,m}^- - \lambda_{c,m}^+ \right) \tag{S36}$$

$$\tau_s \sum_{\text{x}} \frac{\partial a_{\text{x},n_k}^-}{\partial v_{ji}} \left( (\lambda_{a_{\text{x}},n_k}^- - \lambda_{a_{\text{x}},n_k}^+) - c_{n_k}^- G'(a_{\text{x},n_k}^-) \sum_{m \neq n_k} v_{mn}\lambda_{a_{\text{x}},m}^+ \right) \tag{S37}$$

$$\tau_s \sum_{\text{x}} \sum_{m \neq n_k} \frac{\partial a_{\text{x},m}^-}{\partial v_{ji}} \left( \lambda_{a_{\text{x}},m}^- - \lambda_{a_{\text{x}},m}^+ \right) \tag{S38}$$

$$\left. - \tau_s \delta_{in_k} r_{\text{x},n_k}^- c_{n_k}^- \sum_{m \neq n_k} \delta_{jm}\lambda_{a_{\text{x}},m}^+ \right\} \tag{S39}$$

57  where we use $\mathbf{r}_{\text{x}} = G(\mathbf{a}_{\text{x}})$ to denote $G(\mathbf{a}_r) = \mathbf{r}$ and $G(\mathbf{a}_z) = G(\mathbf{a}_u) = \mathbf{1}$, $\delta_{ab}$ is the kronecker delta
58  defined as:

$$\delta_{ab} = \left\{ \begin{array}{ll} 1 & \text{if} \quad a = b, \\ 0 & \text{otherwise} \end{array} \right. \tag{S40}$$

59  Setting the coefficients of $\frac{\partial c^-}{\partial v_{ji}}$ and $\frac{\partial a_{\text{x}}^-}{\partial v_{ji}}$ to 0 (again, using our ability to choose the adjoint variables
60  freely), we can get both $\xi_{\text{x},ijk}$ and the transitions for the adjoint variables.

61  For the parameter updates we get:

$$\xi_{ijk} = -\tau_s \delta_{in_k} r_{\text{x},n_k}^- c_{n_k}^- \sum_{m \neq n_k} \delta_{jm}\lambda_{a_{\text{x}},m}^+ \tag{S41}$$

$$= -\tau_s r_{\text{x},i}^- c_i^- \lambda_{a_{\text{x}},j}^+ . \tag{S42}$$

62  The jumps/transitions of the adjoint variables are:

$$\lambda_{a_{\text{x}},m}^+ = \lambda_{a_{\text{x}},m}^- \tag{S43}$$

$$\lambda_{a_{\text{x}},n_k}^+ = \lambda_{a_{\text{x}},n_k}^- - c_{n_k}^- G'(a_{\text{x},n_k}) \sum_{m \neq n_k} v_{mn}\lambda_{a_{\text{x}},m}^+ \tag{S44}$$

$$\lambda_{c,m}^+ = \lambda_{c,m}^- \tag{S45}$$

$$\tau_m \dot{c}_{n_k}^+ \lambda_{c,n_k}^+ = -(\ell_c^+ - \ell_c^-) + \tau_m \dot{c}_{n_k}^- \lambda_{c,n_k}^- + \tau_m \sum_{m \neq n_k} \lambda_{c,m}^+ (\dot{c}_m^- - \dot{c}_m^+)$$

$$+ \tau_s c_{n_k}^- \sum_{\text{x}} \left( \dot{r}_{\text{x},n_k}^- + \frac{r_{\text{x},n_k}^-}{\tau_s} \right) \sum_{m \neq n_k} v_{mn}\lambda_{a_{\text{x}},m}^+ , \tag{S46}$$

63  where $(\ell_c^+ - \ell_c^-)$ denotes the jumps in the instantaneous loss around event time $s_k$. Thus, all the
64  quantities on the right hand side of Eq. (S22) can be calculated from known quantities.

65  **C.2  Gradient updates for biases $\mathbf{b}_{\text{x}}$**

66  Proceeding similarly for the biases $\mathbf{b}_{\text{x}}$ for each of $\text{x} \in \{u,r,z\}$ (dropping the subscript x for simplicity):

$$\frac{d\mathcal{L}}{db_i} = \frac{d}{db_i} \left\{ \sum_{k=0}^N \int_{s_k}^{s_{k+1}} \left[ \ell_c(\mathbf{c}(t),t) + \boldsymbol{\lambda}_c \cdot f_c + \sum_{\text{x} \in \{u,r,z\}} \boldsymbol{\lambda}_{a_{\text{x}}} \cdot f_{a_{\text{x}}} \right] dt \right\}. \tag{S47}$$

the $\xi_{\mathrm{x},ik}^{\mathrm{bias}}$ term can be shown to be:

$$\xi_{\mathrm{x},ik}^{\mathrm{bias}} = \int_{s_k}^{s_{k+1}} \lambda_{a_{\mathrm{x}},i} dt \tag{S48}$$

with

$$\frac{d\mathcal{L}}{db_i} = \sum_{k=0}^{N} \xi_{\mathrm{x},ik}^{\mathrm{bias}} \tag{S49}$$

### C.3 Gradient updates for input weights $\mathbf{U_x}$

Proceeding similarly for the input weights $\mathbf{U_x}$ for each of $\mathrm{x} \in \{u,r,z\}$ (dropping the subscript x for simplicity):

$$\frac{d\mathcal{L}}{du_{jx}} = \frac{d}{du_{jx}} \left\{ \sum_{k=0}^{N} \int_{s_k}^{s_{k+1}} \left[ \ell_c(\mathbf{c}(t),t) + \boldsymbol{\lambda}_c \cdot f_c + \sum_{\mathrm{x} \in \{u,r,z\}} \boldsymbol{\lambda}_{a_\mathrm{x}} \cdot f_{a_\mathrm{x}} \right] dt \right\}. \tag{S50}$$

the $\xi_{\mathrm{x},jxk}^{\mathrm{input}}$ term can be shown to be:

$$\xi_{\mathrm{x},jxk}^{\mathrm{input}} = -\tau_s \lambda_{a_\mathrm{x},j}^{+} x_x \tag{S51}$$

with

$$\frac{d\mathcal{L}}{du_{jx}} = \sum_{k=0}^{N} \xi_{\mathrm{x},jxk}^{\mathrm{input}} \tag{S52}$$

## D   Details of experiments

### D.1   DVS128 Gesture recognition

In this experiment we use Tonic library [10] to prepare the dataset. The recordings in the dataset are sliced by time without any overlap to produce samples of length 1.7 seconds. The data is denoised with a filter time of 10ms and normalised to [0;1] before being fed to the model. The positive and negative polarity events are represented by 2 separate channels. Our model consists of a preprocessing layer which performs downscaling and flattening transformations, followed by two RNN layers. Both RNN layers have the same number of hidden dimensions. Finally, a fully connected layer of size 11 performs the classification. We use cross-entropy loss and Adam optimizer with default parameters (0.001 learning rate, $\beta_1 = 0.9$, $\beta_2 = 0.999$). The learning rate is scaled by 80% every 100 epochs.

We use additional loss to regularize the output and increase sparsity of the network. The applied regularization losses are shown in Equation S55. $L_{reg}$ is applied indirectly to the active outputs and $L_{act}$ is applied on the auxiliary internal state $c_i^{\langle t \rangle}$, the threshold parameter $\vartheta_i$ is detached from the graph in the second equation so the loss only affects the internal state. We set the regularization weights $w_{reg}$ and $w_v$ to 0.01 and 0.05 respectively.

Fig. S3(a) shows comparison of training curves for LSTM, GRU and EGRU, mean activity of the EGRU network is also shown, the network achieved 80%+ sparsity without significant drop in accuracy. The activities of LSTM and GRU are not shown in Fig S3(a) since they are always 100%. In our experiments we calculate sparsity of these networks as average number of activations close to zero with an absolute tolerance of $1 \times 10^{-8}$, however in Fig. S3(b) we show that even if we increase the absolute tolerance to $1 \times 10^{-3}$, the sparsity of these networks is still an order of magnitude lower than EGRU.

$$L_{reg} = w_{reg} \left( \frac{1}{N} \frac{1}{n_{\mathrm{units}}} \sum_{n=1}^{N} \sum_{1}^{n_{\mathrm{units}}} H\left( c_i^{\langle t \rangle} - \vartheta_i \right) \quad -0.05 \right) \tag{S53}$$

$$L_{act} = w_v \left( \frac{1}{N} \frac{1}{n_{\mathrm{units}}} \sum_{n=1}^{N} \sum_{i=1}^{n_{\mathrm{units}}} c_i - (\vartheta_i - 0.05) \right) \tag{S54}$$

$$\tag{S55}$$

where $N$ spans mini-batch.

### D.1.1 Ablation study

We performed ablation studies, showing the performance of the EGRU models with variation of the gating mechanism. All models in this study are a variation of our **EGRU**(1024) model. The results of these experiments are presented in Table S2. By using a scalar threshold $\vartheta$ where all units share a same threshold parameter we find that the accuracy drops by 2% but the the activity sparsity is increased to 90%.

Next, we evaluate a model with 'hard reset' where the auxiliary internal state $c_i^{\langle t \rangle}$ is set to 0 every time an event is emitted by an unit. We observe a drop in accuracy since this hard reset essentially ignores the input between $t-1$ and $t$. This drop in performance might be significant for applications which require high temporal resolution, which necessitates the term $-y_i^{\langle t-1 \rangle}$ in Eq. (2).

### D.2 Sequential MNIST

All the weights were initialised using Xavier uniform distribution, while the biases were initialised using a uniform distribution. The unit thresholds were initialised using a normal distribution with mean 0 and standard deviation of $\sqrt{2}$, but was transformed to be between 0 and 1 by passing through a standard sigmoid/logistic function after every update. In all the experiments, we used a batch size of 500, and trained the network with Adam with default parameters (0.001 learning rate, $\beta_1 = 0.9$, $\beta_2 = 0.999$) on a cross-entropy loss function. We used gradient clipping with a max gradient norm of 0.25. All models were trained for 200 epochs. The outputs of all the units were convolved with an exponential filter with time constant of 10 time units i.e. with $e^{\frac{1}{10}}$ to calculate an output trace. The value of this trace at the last time step was used to predict the class through a softmax function.

### D.3 PTB Language modeling

Our experimental setup largely follows [13]. In particular, we download and preprocess PennTreebank with their published code [1]. Words are projected to a 400-dimensional dense vector by a linear transformation, followed by three RNN layers. The first two RNN layers feature the same hidden dimension, while the hidden dimension of the last RNN layer always equals the word vector embedding dimension. As common in language modeling, we apply cross entropy loss and use weight tying [6, 14].

We apply the regularization strategies of [13]. Backpropagation through time is conducted with a variable sequence length. With 95% probability, the sequence length is drawn from $\mathcal{N}(70,5)$, and with 5% probability the sequence length is drawn from $\mathcal{N}(35,5)$. We apply variational dropout [4] to the vocabulary with probability $p = 0.1$, to the word embedding vectors with probability $p = 0.4$ as well as to each layer output with probability $p_l$. DropConnect [15] was applied to the hidden-to-hidden weight matricies with probability $p_h$. While vocabulary dropout and embedding dropout where fixed for all models, we tuned layer-to-layer and hidden-to-hidden dropout rates for each model individually. We experimented with both Adam [7] and NT-AvSGD [13] optimization procedures. While Adam lead to competitive results for all models, only LSTM models converged using NT-AvSGD. When optimized with SGD based optimizers, both GRU and EGRU fell behind Adam optimized models. Thus, we optimized all models with Adam, and set momentum to 0 as reported in [12]. Weight decay was set to $1.2 \times 10^{-6}$ as reported in [13]. Gradient clipping was applied to all models, where the magnitude of clipped gradients only made very small differences in results. We trained our models for 800 total epochs. The first 400 epochs were trained at constant learning rate $\lambda$. Then a cosine decay from $\lambda$ to $0.1 \cdot \lambda$ was applied for the remaining 400 epochs.

Parameter search was conducted exclusively on the PTB training and validation set. See table S4 for detailed hyperparameters of the best models.

Our experiments exhibit both forward and backward sparsity also for the challenging task of language modeling. For this task, we don't apply any explicit loss terms to improve sparsity. Figure S4 shows how sparsity evolves during the training process. It is evident that forward sparsity naturally evolves from EGRU, even without explicit sparsity regularization. The sparsity of the backward pass is governed by the width parameter $\epsilon$ of the pseudo-derivative. We experimented with different values of $\epsilon$ on PTB. As shown in figure S5 backward sparsity increases with smaller values of $\epsilon$ as expected. At the same, time convergence speed is slightly limited. For the sake of demonstrating competitive performance, we thus choose $\epsilon = 1$ for the language modeling results reported in this work.

---

[1] https://github.com/salesforce/awd-lstm-lm

| reference | architecture (# units) | para-meters | effective MAC | accu-racy | activity sparsity | backward sparsity |
|---|---|---|---|---|---|---|
| **ours** | LSTM (867) | 16.28M | 20.97M | 87.89% | 0% | - |
| **ours** | GRU (1024) | 15.75M | 15.73M | 88.07% | 0% | - |
| **ours** | **EGRU** (512) | 5.51M | 4.19M | 88.02% | 83.79% | 53.55% |
| **ours** | **EGRU** (1024) | 15.75M | 10.54M | 90.22% | 82.53% | 56.63% |
| **ours** | **EGRU+DA** (1024) | 15.75M | 10.77M | 97.13% | 78.77% | 58.20% |

**Table S1:** Model performance over multiple runs for the DVS Gesture recognition task. Effective number of MAC operations as described in section 3.5.

| model | accuracy | activity sparsity |
|---|---|---|
| **EGRU**(1024) | 90.2 | 82.5 |
| without regularization | 89.3 | 76.5 |
| scalar $\vartheta$ | 88.3 | 90.8 |
| hard reset | 87.2 | 90 |

**Table S2**

| architecture (# units) | parameters | effective MAC (mean±std) | test accuracy (%) (mean±std) | activity sparsity (%) (mean±std) | backwards sparsity (%) at epochs 20/50/100 (mean±std) |
|---|---|---|---|---|---|
| GRU (512)* | 791K | 795K | 98.7±0.2 | - | - |
| GRU (590) | 1.049M | 1.054M | 98.7±0.1 | - | - |
| **EGRU** (590) | 1.048M | 210K±51K | 95.5±1.6 | 80.5±4.9 | 24.9±6.8 / 26.1±5.9 / 26.0±4.3 |

**Table S3:** Model performance over 4 runs for sequential MNIST task. Test scores are given as percentage accuracy, where higher is better. (*) mean over 3 runs. The fourth run was unsuccessful due to network instability.

| RNN cell | val ppl best | val ppl mean±std | hidden units | learning rate | batch size | dropout $p_h$ | dropout $p_l$ | gradient clip |
|---|---|---|---|---|---|---|---|---|
| LSTM | 61.0 | 61.2±0.2 | 1150 | 0.003 | 40 | 0.5 | 0.4 | 0.25 |
| GRU | 75.1 | 75.2±0.1 | 1350 | 0.001 | 80 | 0.5 | 0.4 | 0.25 |
| **EGRU** | | | 1350 | 0.0003 | 80 | 0.5 | 0.3 | 2.0 |
| **EGRU** | | | 2700 | 0.0003 | 80 | 0.5 | 0.3 | 2.0 |

**Table S4:** Best parameters for the trained models. Mean and uncertainty are calculated from 10 runs with different random seeds.

| reference | architecture (# units) | para-meters | effective MAC | validation | test | activity sparsity |
|---|---|---|---|---|---|---|
| **ours** | LSTM (1150) | 24M | 24M | 63.6 | 59.7 | - |
| **ours** | GRU (1350) | 24M | 24M | 74.8 | 70.1 | - |
| **ours** | **EGRU** (1350) | 24M | 5.4M | 73.1 | 68.48 | 85.7% |
| **ours** | **EGRU** (2700) | 76M | 8.4M | 72.0 | 67.8 | 91.3% |

**Table S5:** Model performance over multiple runs for PennTreebank. Validation and test scores are given as perplexities, where lower is better. Sparsity refers to activity-sparsity of the EGRU output, and effecive MAC operations consider the layer-wise sparsity in the forward pass.

## E Dataset licenses

Penn Treebank [11] is subject to the Linguistic Data Consortium User Agreement for Non-Members [2]. Following [13], we download Penn Treebank data from http://www.fit.vutbr.cz/~imikolov/rnnlm/simple-examples.tgz.

> LDC Not-For-Profit members, government members and nonmember licensees may use LDC data for noncommercial linguistic research and education only. For-profit organizations who are or were LDC members may conduct commercial technology development with LDC data received when the organization was an LDC for-profit member unless use of that data is otherwise restricted by a corpus-specific license agreement. Not-for-profit members, government members and nonmembers, including nonmember for-profit organizations, cannot use LDC data to develop or test products for commercialization, nor can they use LDC data in any commercial product or for any commercial purpose.

The DVS128 Gesture Dataset [1] is released under the Creative Commons Attribution 4.0 license and can be retrieved from: https://research.ibm.com/interactive/dvsgesture/. We used Tonic library [10] for Pytorch to preprocess data and to apply transformations.

The sequential MNIST task [8] is based on the MNIST dataset first introduced in [9].
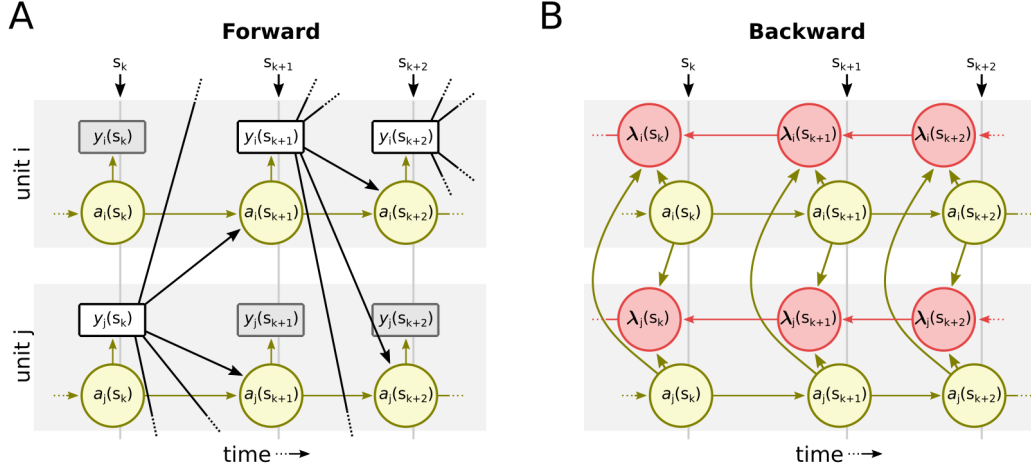
## F Hardware and software details

Most of our experiments were run on NVIDIA A100 GPUs. Some initial hyper-parameter searches were conducted on NVIDIA V100 and Quadro RTX 5000 GPUs. We used about 12,000 computational hours in total for training and hyper-parameter searches. All models and experiments were implemented in PyTorch. For the continuous time EGRU model, we also used the torchdiffeq [2] library.

The machines used for the DVS128 gesture recognition task and for the PTB language modeling task feature 8x NVIDIA A100-SXM4 (40GB) GPUs, 2x AMD EPYC CPUs 7352 with 24 cores each, and 1TB RAM on each compute node. For each run, we only use a single GPU, and a fraction of the cores and memory available on the node to run multiple experiments in parallel. The nodes operate Red Hat Enterprise Linux Server (release 7.9).
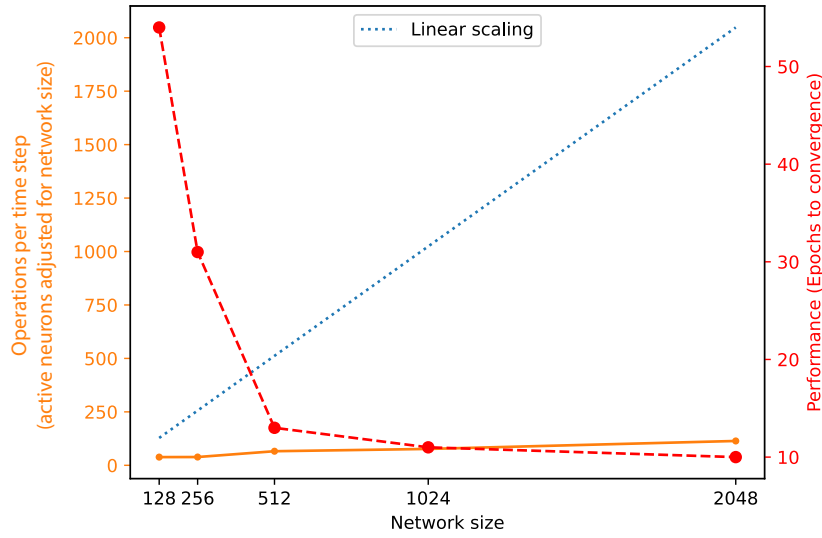
---

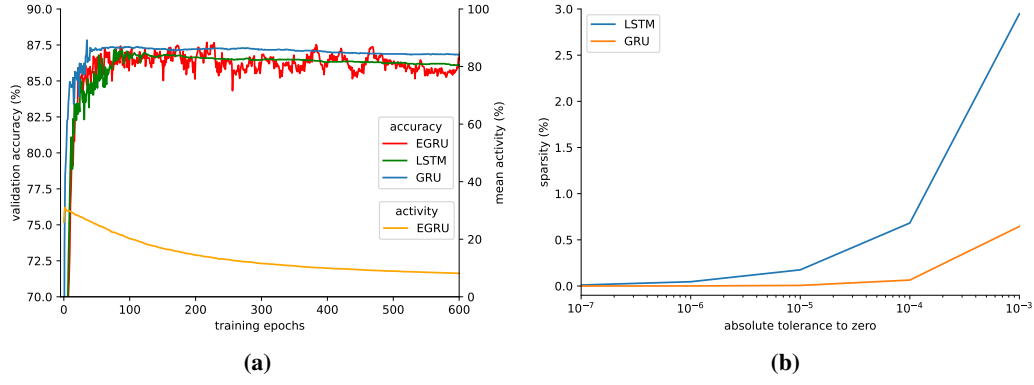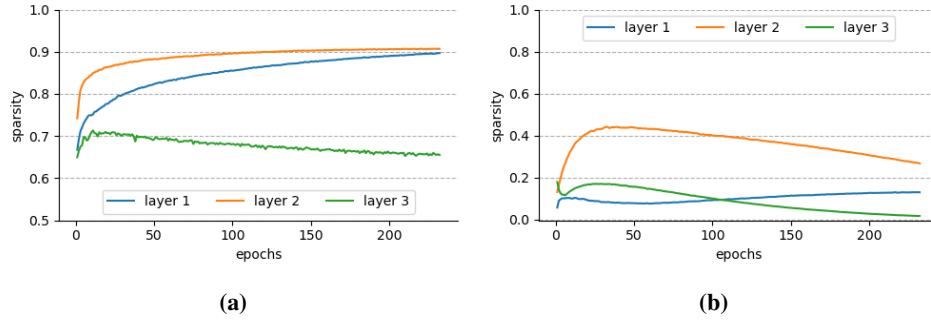[2] https://www.ldc.upenn.edu/data-management/using/licensing

**Figure S1:** Illustrate the event-based state dynamics for two EGRU units ($i$ and $j$) (A) Forward dynamics. Information only propagates from units that generate an event. (B) Backward dynamics.
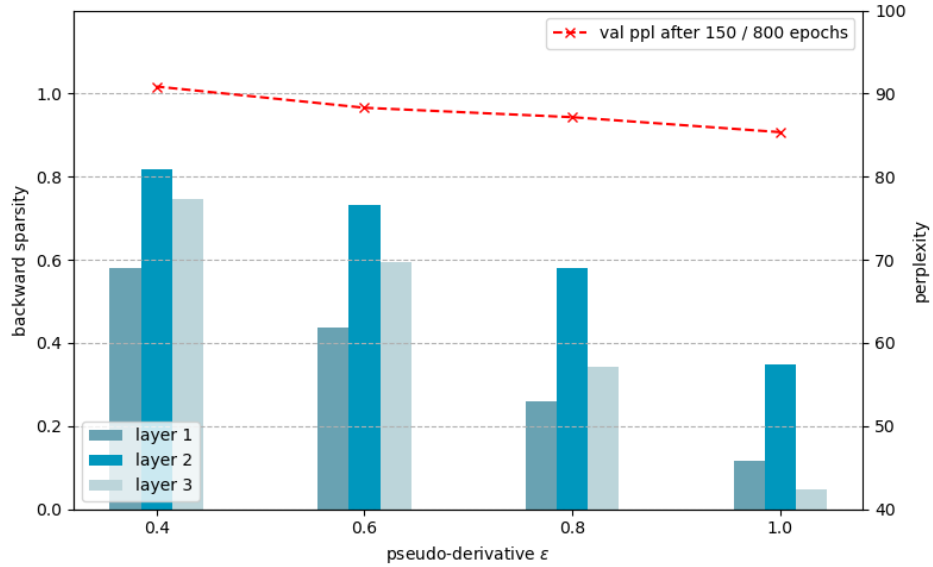


**Figure S2:** Illustration of the scaling properties of the EGRU on a $14 \times 14$ sequential MNIST task. As the size of the network increases, the network converges faster. Increasing the network size 10x increases the speed of convergence 5x, while increasing the total amount of computation per sample only 2x. The total amount of computation is adjusted for network size. The smaller subsampled 14x14 sMNIST task was chosen here for reasons of computational limitations.

**Figure S3:** **(a)** mean training curves over 5 runs for DVS gesture task. **(b)** activity sparsity of LSTM and GRU for DVS gesture task across various values of absolute tolerance to zero.



**Figure S4:** EGRU with 1350 hidden units on the Penn Treebank language modeling task with pseudo-derivative $\epsilon = 1$ . **(a)** layer-wise forward sparsity **(b)** layer-wise backward sparsity



**Figure S5:** Backward sparsity for EGRU with 1350 hidden units on the Penn Treebank language modeling task with varying pseudo-derivative support $\epsilon$

10

# References

[1] A. Amir, B. Taba, D. Berg, T. Melano, J. McKinstry, C. Di Nolfo, T. Nayak, A. Andreopoulos, G. Garreau, M. Mendoza, et al. A low power, fully event-based gesture recognition system. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7243–7252, 2017.

[2] R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural Ordinary Differential Equations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 6572–6583. Curran Associates, Inc., 2018.

[3] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.

[4] Y. Gal and Z. Ghahramani. A theoretically grounded application of dropout in recurrent neural networks. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper/2016/file/076a0c97d09cf1a0ec3e19c7f2529f2b-Paper.pdf.

[5] T. H. Gronwall. Note on the derivatives with respect to a parameter of the solutions of a system of differential equations. *Annals of Mathematics*, pages 292–296, 1919.

[6] H. Inan, K. Khosravi, and R. Socher. Tying word vectors and word classifiers: A loss framework for language modeling. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL https://openreview.net/forum?id=r1aPbsFle.

[7] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Y. Bengio and Y. LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL http://arxiv.org/abs/1412.6980.

[8] Q. V. Le, N. Jaitly, and G. E. Hinton. A simple way to initialize recurrent networks of rectified linear units. *arXiv preprint arXiv:1504.00941*, 2015.

[9] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

[10] G. Lenz, K. Chaney, S. B. Shrestha, O. Oubari, S. Picaud, and G. Zarrella. Tonic: event-based datasets and transformations., July 2021. URL https://doi.org/10.5281/zenodo.5079802. Documentation available under https://tonic.readthedocs.io.

[11] M. P. Marcus, M. A. Marcinkiewicz, and B. Santorini. Building a large annotated corpus of english: The penn treebank. *Comput. Linguist.*, 19(2):313–330, jun 1993. ISSN 0891-2017.

[12] G. Melis, C. Dyer, and P. Blunsom. On the state of the art of evaluation in neural language models. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL https://openreview.net/forum?id=ByJHuTgA-.

[13] S. Merity, N. S. Keskar, and R. Socher. Regularizing and Optimizing LSTM Language Models. *arXiv:1708.02182 [cs]*, Aug. 2017.

[14] O. Press and L. Wolf. Using the output embedding to improve language models. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 157–163, Valencia, Spain, Apr. 2017. Association for Computational Linguistics. URL https://aclanthology.org/E17-2025.

[15] L. Wan, M. Zeiler, S. Zhang, Y. Le Cun, and R. Fergus. Regularization of neural networks using dropconnect. In S. Dasgupta and D. McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1058–1066, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL https://proceedings.mlr.press/v28/wan13.html.

[16] T. C. Wunderlich and C. Pehle. Event-based backpropagation can compute exact gradients for spiking neural networks. *Scientific Reports*, 11(1):12829, June 2021. ISSN 2045-2322. doi: 10.1038/s41598-021-91786-z. URL https://www.nature.com/articles/s41598-021-91786-z.

[17] W. Yang, D. Yang, and Y. Fan. A proof of a key formula in the error-backpropagation learning algorithm for multiple spiking neural networks. In *International Symposium on Neural Networks*, pages 19–26. Springer, 2014.