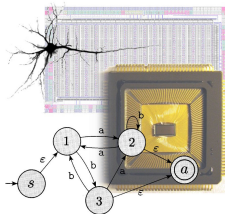


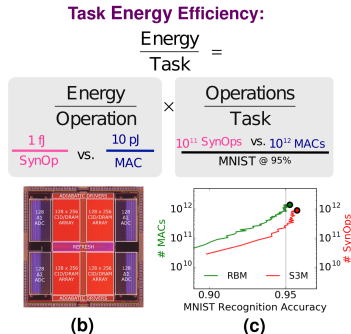
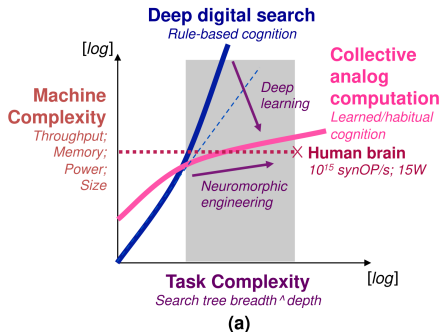
# Event-Driven Random Backpropagation: Enabling Neuromorphic Deep Learning Machines

Emre Neftci

Department of Cognitive Sciences, UC Irvine,  
Department of Computer Science, UC Irvine,

March 7, 2017





Cauwenberghs, *Proceedings of the National Academy of Sciences*, 2013

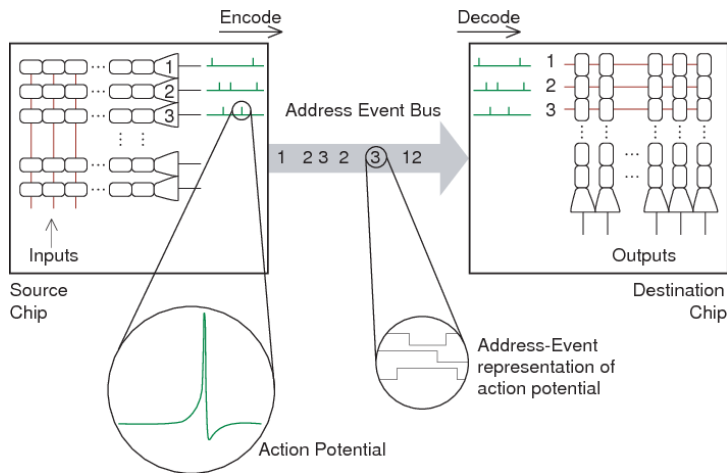
Karakiewicz, Genov, and Cauwenberghs, *IEEE Sensors Journal*, 2012

Neftci, Augustine, Paul, and Detorakis, *arXiv preprint arXiv:1612.05596*, 2016

## 1000x power improvements compared to future GPU technology through two factors:

- Architecture and device level optimization in event-based computing
- Algorithmic optimization in neurally inspired learning and inference

# Neuromorphic Computing Can Enable Low-power, Massively Parallel Computing



- Only spikes are communicated & routed between neurons (weights, internal states are local)
- To use this architecture for practical workloads, we need algorithms that operate on local information

For many industrial applications involving controlled environments, where existing data is readily available, off-chip/off-line learning is often sufficient.

So why do embedded learning?

### **Two main use cases:**

- Mobile, low-power platform in uncontrolled environments, where adaptive behavior is required.
- Working around device mismatch/non-idealities.

### **Potentially rules out:**

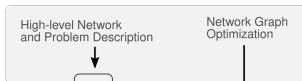
- Self-driving cars
- Data mining
- Fraud Detection

## Neuromorphic Learning Machines: Online learning for data-driven autonomy and algorithmic efficiency

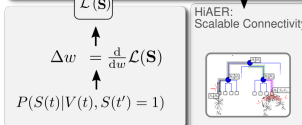
VARIABLE COMPLEXITY TASKS



EVENT-BASED MACHINE LEARNING FRAMEWORK

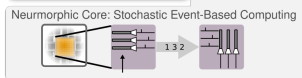


STOCHASTIC LEARNING



THEORY OF EVENT-BASED STOCHASTIC COMPUTING

NEUROMORPHIC HARDWARE DESIGN



- **Hardware & Architecture:** Scalable Neuromorphic Learning Hardware Design
- **Programmability:** Neuromorphic supervised, unsupervised and reinforcement learning framework



theano



---

neon\_mlp\_extract.py

---

# setup model layers

```
layers = [Affine(nout=100, init=init_norm, activation=Rectlin()),  
          Affine(nout=10, init=init_norm, activation=Logistic(shortcut=True))]
```

# setup cost function as CrossEntropy

```
cost = GeneralizedCost(costfunc=CrossEntropyBinary())
```

# setup optimizer

```
optimizer = GradientDescentMomentum(  
    0.1, momentum_coef=0.9, stochastic_round=args.rounding)
```

---

Can we design a digital neuromorphic learning machine  
that is flexible and efficient?

- **Leaky Stochastic I&F Neuron (LIF)**

$$V[t+1] = -\alpha V[t] + \sum_{j=1}^n \xi_j w_j(t) s_j(t) \quad (1a)$$

$$V[t+1] \geq T : V[t+1] \leftarrow V_{reset} \quad (1b)$$



- LIF with first order kinetic synapse

$$V[t + 1] = -\alpha V[t] + I_{syn} \quad (2a)$$

$$I_{syn}[t + 1] = -a_1 I_{syn}[t] + \sum_{j=1}^n w_j(t) s_j(t) \quad (2b)$$

$$V[t + 1] \geq T : V[t + 1] \leftarrow V_{reset} \quad (2c)$$

- **LIF with second order kinetic synapse**

$$V[t + 1] = -\alpha V[t] + I_{syn} + I_{syn}, \quad (3a)$$

$$I_{syn}[t + 1] = -a_1 I_{syn}[t] + c_1 I_s[t] + \eta[t] + b \quad (3b)$$

$$I_s[t + 1] = -a_2 I_s[t] + \sum_{j=1}^n w_j s_j[t] \quad (3c)$$

$$V[t + 1] \geq T : V[t + 1] \leftarrow V_{reset} \quad (3d)$$

- **Dual-Compartment LIF with synapses**

$$V_1[t + 1] = -\alpha V_1[t] + \alpha_{21} V_2[t] \quad (4a)$$

$$V_2[t + 1] = -\alpha V_2[t] + \alpha_{12} V_1[t] + I_{syn} \quad (4b)$$

$$I_{syn}[t + 1] = -a_1 I_{syn}[t] + \sum_{j=1}^n w_j^1(t) s_j(t) + \eta[t] + b \quad (4c)$$

$$V_1[t + 1] \geq T : V_1[t + 1] \leftarrow V_{reset} \quad (4d)$$

- **Mihalas Niebur Neuron (MNN)**

$$V[t + 1] = \alpha V[t] + I_e - G \cdot E_L + \sum_{i=1}^n I_i[t] \quad (5a)$$

$$\Theta[t + 1] = (1 - b)\Theta[t] + aV[t] - aE_L + b \quad (5b)$$

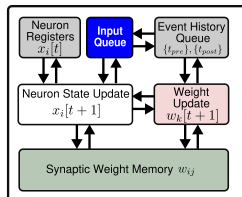
$$I_1[t + 1] = -\alpha_1 I_1[t] \quad (5c)$$

$$I_2[t + 1] = -\alpha_2 I_2[t] \quad (5d)$$

$$V[t + 1] \geq \Theta[t + 1] : \text{Reset}(V[t + 1], I_1, I_2, \Theta) \quad (5e)$$

MNN can produce a wide variety of spiking behaviors

NSAT Core (2048 Neurons)



Neuron and Synapse Model

$$\begin{aligned} \mathbf{x}[t+1] = & \mathbf{A}\mathbf{x}[t] && \text{(Leak \& Coupling)} \\ & + \mathbf{\Xi}[t] \otimes \mathbf{W}[t]\mathbf{S}[t] && \text{(Synaptic inputs)} \\ & + \boldsymbol{\eta}[t] && \text{(Noise)} \end{aligned}$$

$$\mathbf{x}[t+1] \geq \boldsymbol{\theta}, \mathbf{x}[t+1] \leftarrow \mathbf{X}_r \quad \text{(Thresholds \& Reset)}$$

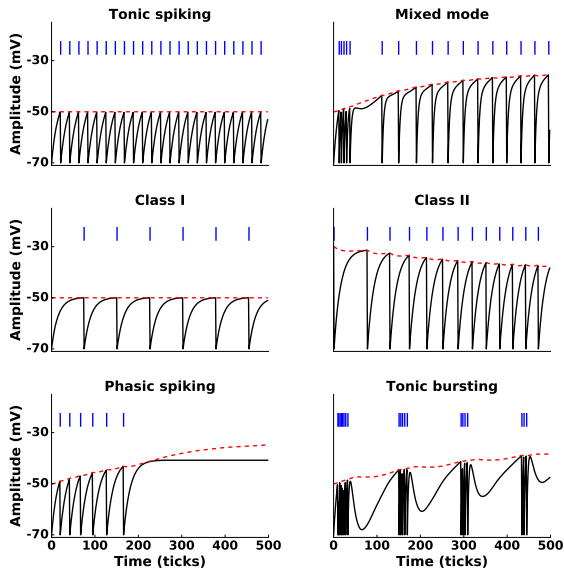
$$x_0[t+1] \geq \theta_0, s_i[t+1] \leftarrow 1 \quad \text{(Spiking Output)}$$

$$e_k = x_m[t] (K[t - t_k] + K[t_k - t_{last}]) \quad \text{(Eligibility)}$$

$$w_k[t+1] = w_k[t] + s_k[t+1]e_k \quad \text{(Weight update)}$$

- Multicompartment generalized integrate-and-fire neurons
- Multiplierless design
- Weight sharing (convnets) at the level of the core

Equivalent software simulations for analyzing fault tolerance, precision, performance, and efficiency trade-offs (available publicly soon!)



$$w_k[t + 1] = w_k[t] + s_k[t + 1]e_k \quad \text{(Weight update)}$$

$$e_k = x_m \underbrace{(K[t - t_k] + K[t_k - t_{last}])}_{STDP} \quad \text{(Eligibility)}$$

$$x_m = \sum_i \gamma_i x_i \quad \text{(Modulation)}$$

Detorakis, Augustine, Paul, Pedroni, Sheik, Cauwenberghs, and Neftci (in preparation)

$$w_k[t + 1] = w_k[t] + s_k[t + 1]e_k \quad (\text{Weight update})$$

$$e_k = x_m \underbrace{(K[t - t_k] + K[t_k - t_{last}])}_{STDP} \quad (\text{Eligibility})$$

$$x_m = \sum_i \gamma_i x_i \quad (\text{Modulation})$$

Detorakis, Augustine, Paul, Pedroni, Sheik, Cauwenberghs, and Nefci (in preparation)

## Based on two insights:

*Causal and acausal STDP weight updates on pre-synaptic spikes only, using only forward lookup access of the synaptic connectivity table*

Pedroni et al., 2016

*“Plasticity involves as a third factor a local dendritic potential, besides pre- and postsynaptic firing times”*

Urbanczik and Senn, *Neuron*, 2014

Clopath, Büsing, Vasilaki, and Gerstner, *Nature Neuroscience*, 2010



- **Reinforcement Learning**

$$\Delta w_{ij} = \eta r STDP_{ij}$$

Florian, *Neural Computation*, 2007

- **Unsupervised Representation Learning**

$$\Delta w_{ij} = \eta g(t) STDP_{ij}$$

Neftci, Das, Pedroni, Kreutz-Delgado, and Cauwenberghs, *Frontiers in Neuroscience*, 2014

- **Unsupervised Sequence Learning**

$$\Delta w_{ij} = \eta (\Theta(V) - \alpha(\nu_i - C)) \nu_j$$

Sheik et al. 2016

- **Supervised Deep Learning**

$$\Delta w_{ij} = \eta (\nu_{tgt} - \nu_i) \phi'(V) \nu_j$$

Neftci, Augustine, Paul, and Detorakis, *arXiv preprint arXiv:1612.05596*, 2016

- **Reinforcement Learning**

$$\Delta w_{ij} = \eta r STDP_{ij}$$

Florian, *Neural Computation*, 2007

- **Unsupervised Representation Learning**

$$\Delta w_{ij} = \eta g(t) STDP_{ij}$$

Neftci, Das, Pedroni, Kreutz-Delgado, and Cauwenberghs, *Frontiers in Neuroscience*, 2014

- **Unsupervised Sequence Learning**

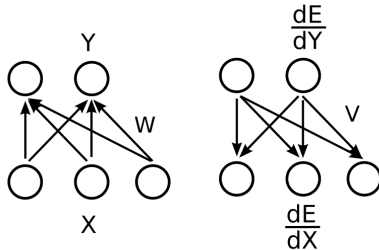
$$\Delta w_{ij} = \eta (\Theta(V) - \alpha(\nu_i - C)) \nu_j$$

Sheik et al. 2016

- **Supervised Deep Learning**

$$\Delta w_{ij} = \eta (\nu_{tgt} - \nu_i) \phi'(V) \nu_j$$

Neftci, Augustine, Paul, and Detorakis, *arXiv preprint arXiv:1612.05596*, 2016

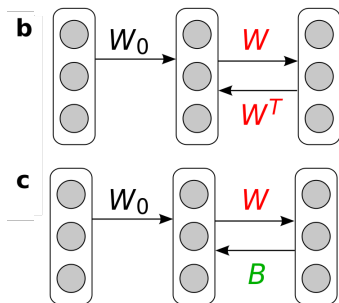


(A) ForwardProp

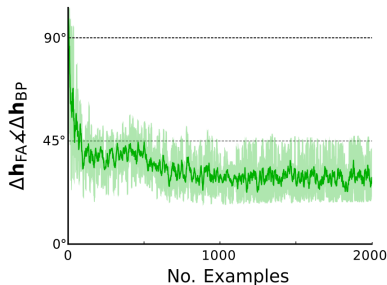
(B) BackProp

## Potential incompatibilities of BP on a neural (neuromorphic) substrate:

- 1 Symmetric Weights
- 2 Computing Multiplications and Derivatives
- 3 Propagating error signals with high precision
- 4 Precise alternation between forward and backward passes
- 5 Synaptic weights can change sign
- 6 Availability of targets



**e**



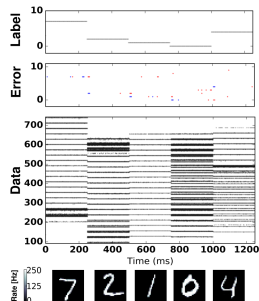
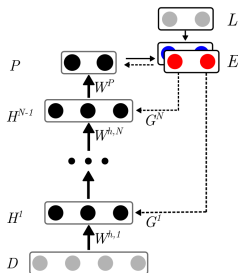
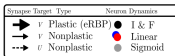
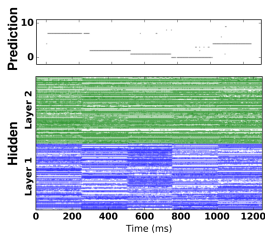
Replace weight matrices in backprop phase with (fixed) random weights

Lillicrap, Cownden, Tweed, and Akerman, *arXiv preprint arXiv:1411.0247*, 2014

Baldi, Sadowski, and Lu, *arXiv preprint arXiv:1612.02734*, 2016

- Event-driven Random Backpropagation Learning Rule: Error-modulated, membrane voltage-gated, event-driven, supervised.

$$\Delta w_{ik} \propto \underbrace{\phi'(I_{syn,i}[t])}_{\text{Derivative}} S_k[t] \sum_j \underbrace{G_{ij} (L_j[t] - P_j[t])}_{\text{Error}} \quad (\text{eRBP})$$



- Event-driven Random Backpropagation Learning Rule: Error-modulated, membrane voltage-gated, event-driven, supervised.

$$\Delta w_{ik} \propto \underbrace{\phi'(I_{syn,i}[t])}_{\text{Derivative}} S_k[t] \underbrace{\sum_j G_{ij} \underbrace{(L_j[t] - P_j[t])}_{\text{Error}}}_{T_i} \quad (\text{eRBP})$$

Approximate derivative with a boxcar function:

**function** ERBP

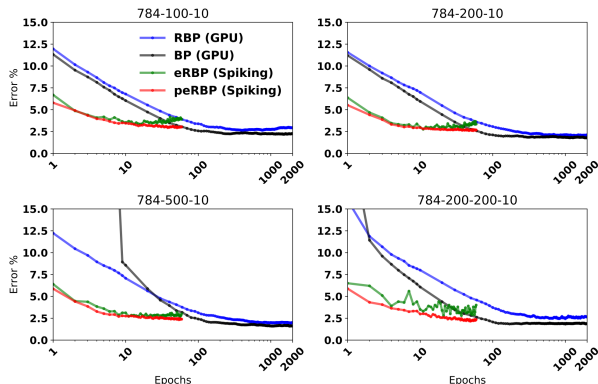
**for**  $k \in \{\text{presynaptic spike addresses } \mathbf{S}^{pre}\}$  **do**

**if**  $b_{min} < I < b_{max}$  **then**  $w_k \leftarrow w_k + T$ ,

**end if**

**end for**

**end function**



## Network

Dataset

PI MNIST 784-100-10

PI MNIST 784-200-10

PI MNIST 784-500-10

PI MNIST 784-200-200-10

PI MNIST 784-500-500-10

## Classification Error

eRBP

peRBP

RBP (GPU)

BP (GPU)

3.94%

3.53%

2.76%

3.48%

2.02%

2.74%

2.15%

2.08%

2.42%

2.20%

2.19%

1.81%

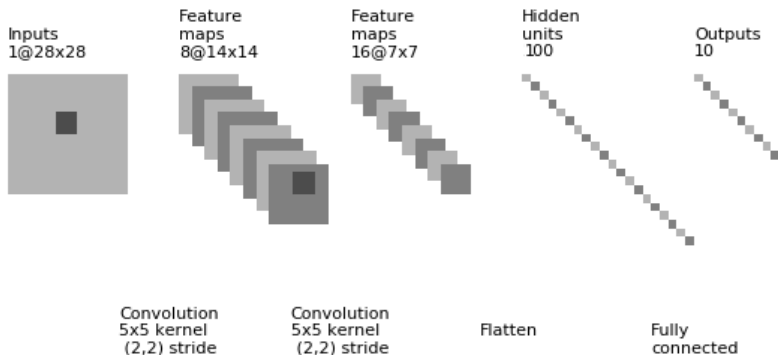
1.8%

1.91%

1.90%

peRBP = eRBP with stochastic synapses

## peRBP MNIST Benchmarks (Convolutional Neural Net)



### Network

Dataset  
MNIST

### Classification Error

peRBP

**3.8 (5 epochs)%**

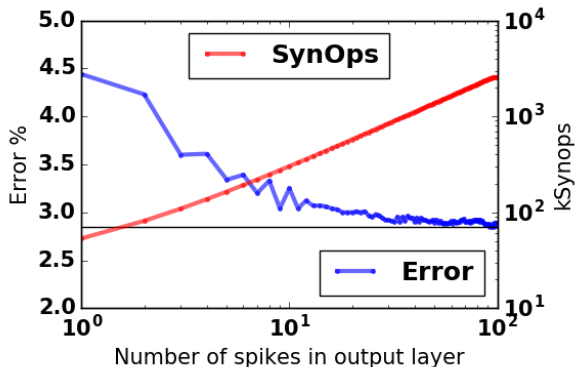
RBP (GPU)

1.95%

BP (GPU)

1.23%

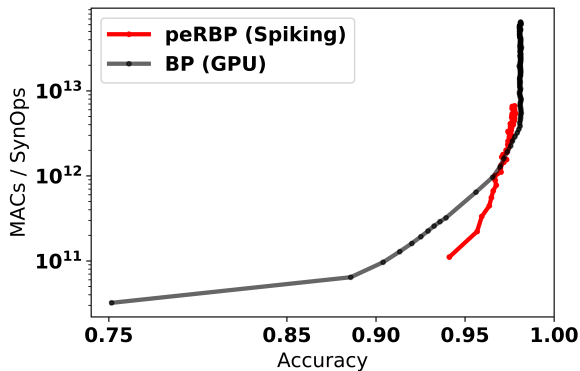




### Energy Efficiency During Inference:

- Inference:  $\cong 100k$  Synops until first spike: <5% error, 100,000 SynOps per classification

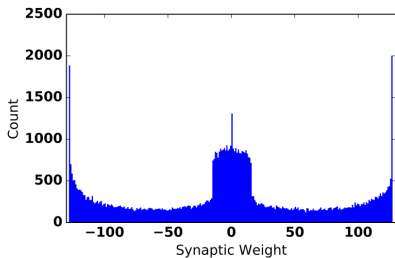
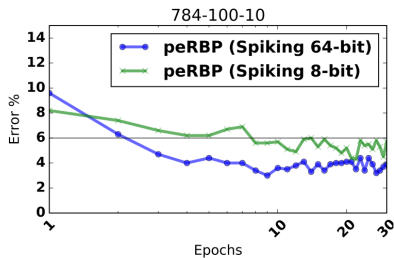
	eRBP	DropConnect (GPU)	Spinnaker	True North
Implementation	(20 pJ/Synop)	CPU/GPU	ASIC	ASIC
Accuracy	95%	99.79%	95%	95%
Energy/classify	2 $\mu J$	1265 $\mu J$	6000 $\mu J$	4 $\mu J$
Technology		28 nm	Unknown	28 nm



## Energy Efficiency During Training:

- Training: SynOp-MAC parity

Embedded local plasticity dynamics for continuous (life-long) learning



- 16 bits neural states
- 8 bits synaptic weights
- $\cong$  1Mbit Synaptic Weight Memory

All-digital implementation for exploring scalable event-based learning

### Summary:

- 1 NSAT: Flexible and efficient neural learning machines
- 2 Supervised deep learning with event-driven random back-propagation can achieve good learning results at  $>100\times$  energy improvements

### Challenges:

- 1 Catastrophic Forgetting: Need for Hippocampus, Intrinsic Replay and Neurogenesis
- 2 Build a neuromorphic library of “deep learning tricks” (Batch normalization, Adam, ...)

## Collaborators:



Georgios Detorakis  
(UCI)



Somnath Paul (Intel)



Charles Augustine  
(Intel)

## Support:





P. Baldi, P. Sadowski, and Zhiqin Lu. “Learning in the Machine: Random Backpropagation and the Learning Channel”. In: *arXiv preprint arXiv:1612.02734* (2016).



Gert Cauwenberghs. “Reverse engineering the cognitive brain”. In: *Proceedings of the National Academy of Sciences* 110.39 (2013), pp. 15512–15513.



C. Clopath, L. BÜsing, E. Vasilaki, and W. Gerstner. “Connectivity reflects coding: a model of voltage-based STDP with homeostasis”. In: *Nature Neuroscience* 13.3 (2010), pp. 344–352.



R.V. Florian. “Reinforcement learning through modulation of spike-timing-dependent synaptic plasticity”. In: *Neural Computation* 19.6 (2007), pp. 1468–1502.



R. Karakiewicz, R. Genov, and G. Cauwenberghs. “1.1 TMACS/mW Fine-Grained Stochastic Resonant Charge-Recycling Array Processor”. In: *IEEE Sensors Journal* 12.4 (Apr. 2012), pp. 785–792.



Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Akerman. “Random feedback weights support learning in deep neural networks”. In: *arXiv preprint arXiv:1411.0247* (2014).



S. Mihalas and E. Niebur. “A generalized linear integrate-and-fire neural model produces diverse spiking behavior”. In: *Neural Computation* 21 (2009), pp. 704–718.



E. Neftci, S. Das, B. Pedroni, K. Kreutz-Delgado, and G. Cauwenberghs. “Event-Driven Contrastive Divergence for Spiking Neuromorphic Systems”. In: *Frontiers in Neuroscience* 7.272 (Jan. 2014). ISSN: 1662-453X. DOI: 10.3389/fnins.2013.00272. URL: [http://www.frontiersin.org/neuromorphic\\_engineering/10.3389/fnins.2013.00272/abstract](http://www.frontiersin.org/neuromorphic_engineering/10.3389/fnins.2013.00272/abstract).



Emre Neftci, Charles Augustine, Somnath Paul, and Georgios Detorakis. “Event-driven Random Back-Propagation: Enabling Neuromorphic Deep Learning Machines”. In: *arXiv preprint arXiv:1612.05596* (2016).



Bruno U Pedroni, Sadique Sheik, Siddharth Joshi, Georgios Detorakis, Somnath Paul, Charles Augustine, Emre Neftci, and Gert Cauwenberghs. “Forward Table-Based Presynaptic Event-Triggered Spike-Timing-Dependent Plasticity”. In: Oct. 2016. URL: <https://arxiv.org/abs/1607.03070>.  
URL: %7BIEEE%20Biomedical%20Circuits%20and%20Systems%20Conference%20(BioCAS),%20https://arxiv.org/abs/1607.03070%7D.



Robert Urbanczik and Walter Senn. “Learning by the dendritic prediction of somatic spiking”. In: *Neuron* 81.3 (2014), pp. 521–528.