

Training deep resistive networks with equilibrium propagation

Aspen Winter Conference
Computing with Physical Systems

11 January 2024

Benjamin Scellier

Research scientist at Rain AI



Acknowledgements



Jack Kendall



Maxence Ernout



Suhas Kumar



Ross Pantone



Yoshua Bengio



Axel Laborieux



Julie Grollier



Damien Querlioz



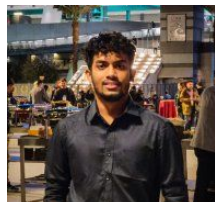
Yann Ollivier



Siddhartha
Mishra



Kalpana
Manickavasagam



Vidyesh Anisetti



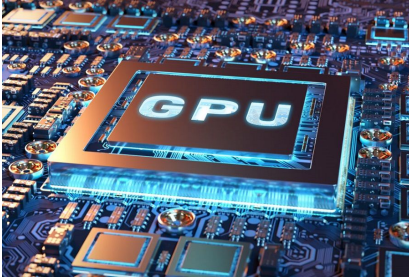
Ananth Kandala



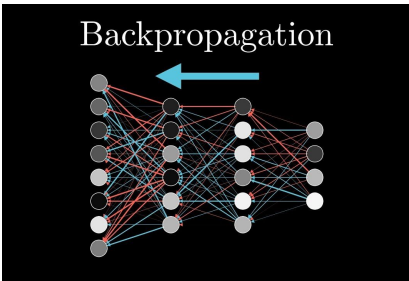
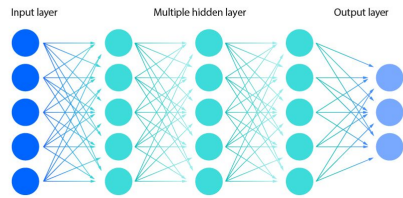
Jennifer Schwarz

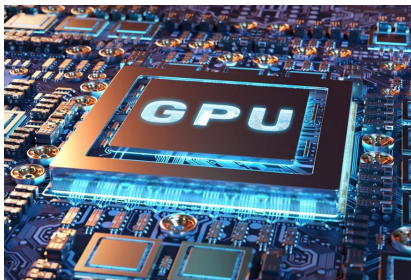


Deep Learning Computing Paradigm

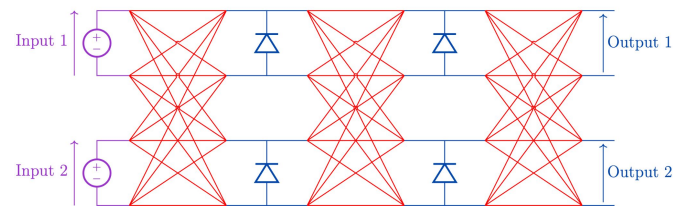
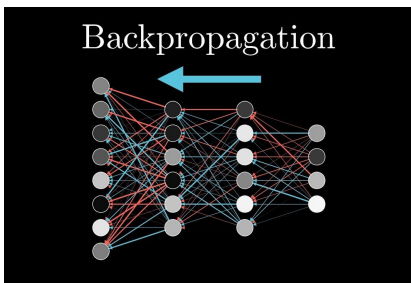
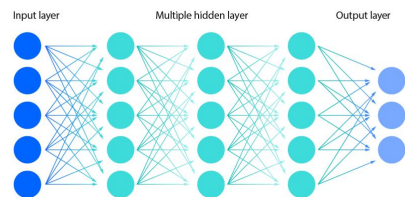


Deep neural network

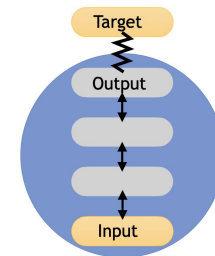




Deep neural network



Deep resistive network



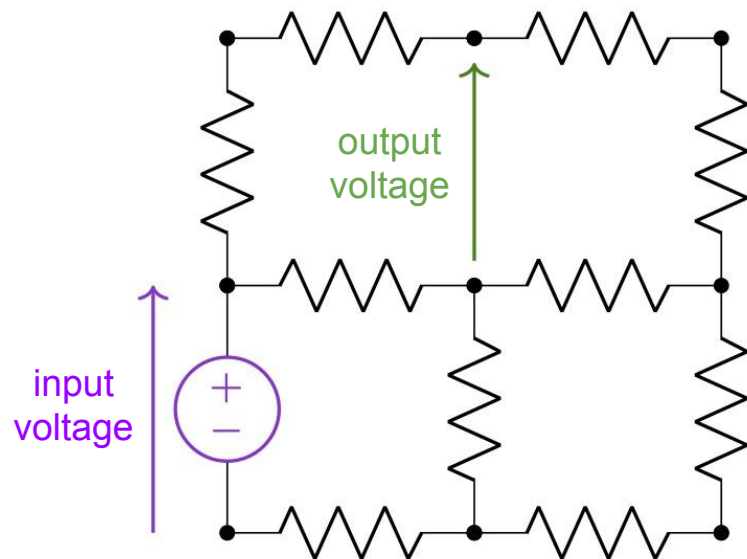
Equilibrium propagation

What is a deep resistive network?

Computing with a Resistor Network

Circuit elements:

- voltage sources (inputs)

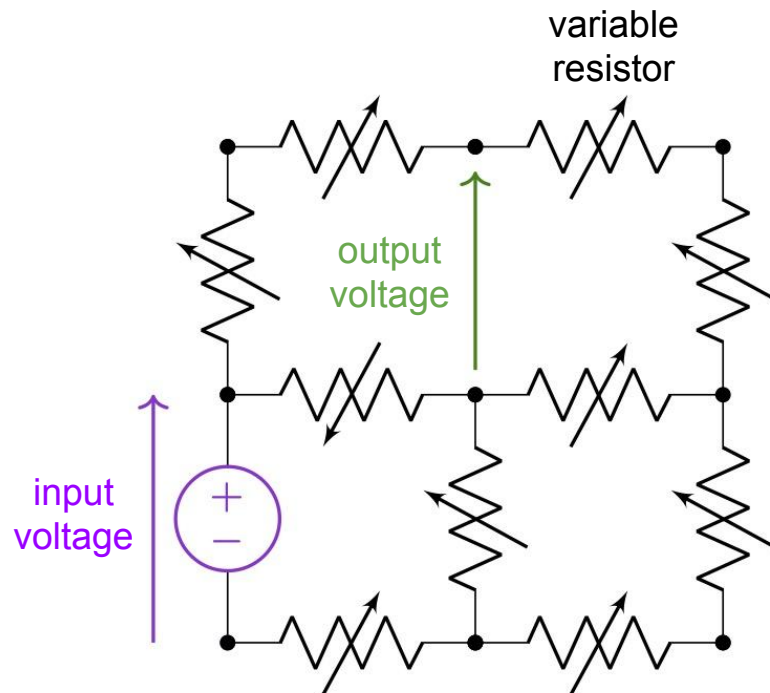


$$\text{output} = F(\text{input})$$

Computing with a Resistor Network

Circuit elements:

- voltage sources (inputs)
- variable resistors (trainable weights)

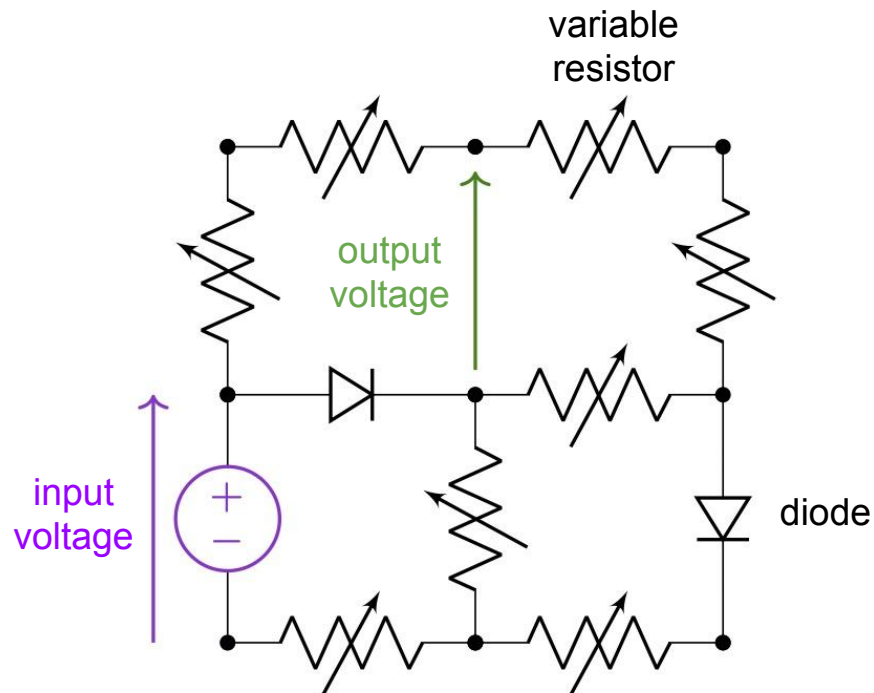


$$\text{output} = F_{\text{conductances}}(\text{input})$$

Computing with a Resistor Network

Circuit elements:

- voltage sources (inputs)
- variable resistors (trainable weights)
- diodes (nonlinearities)



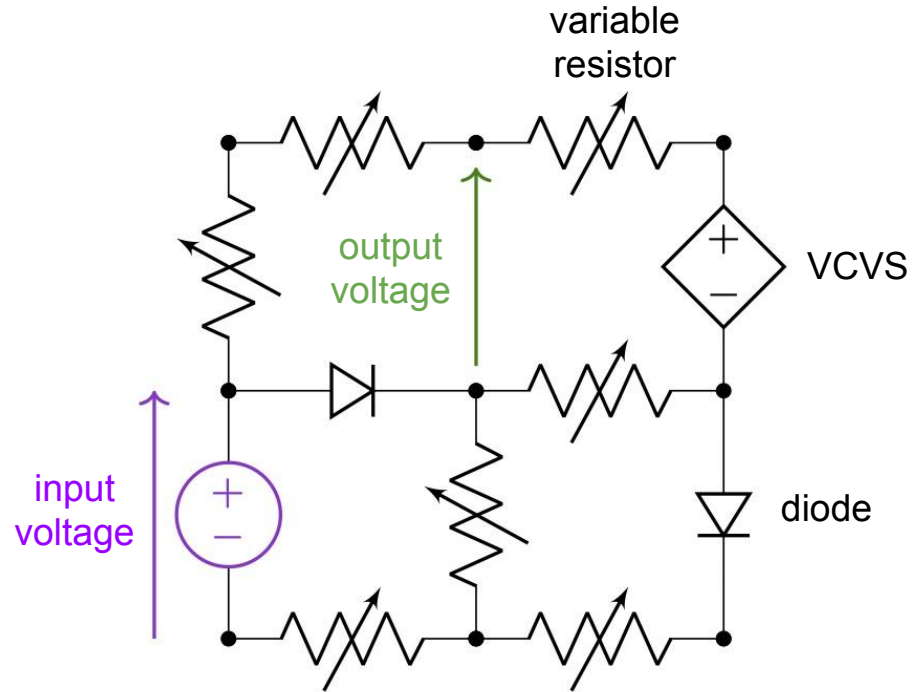
$$\text{output} = F_{\text{conductances}}(\text{input})$$

Computing with a Resistor Network

Circuit elements:

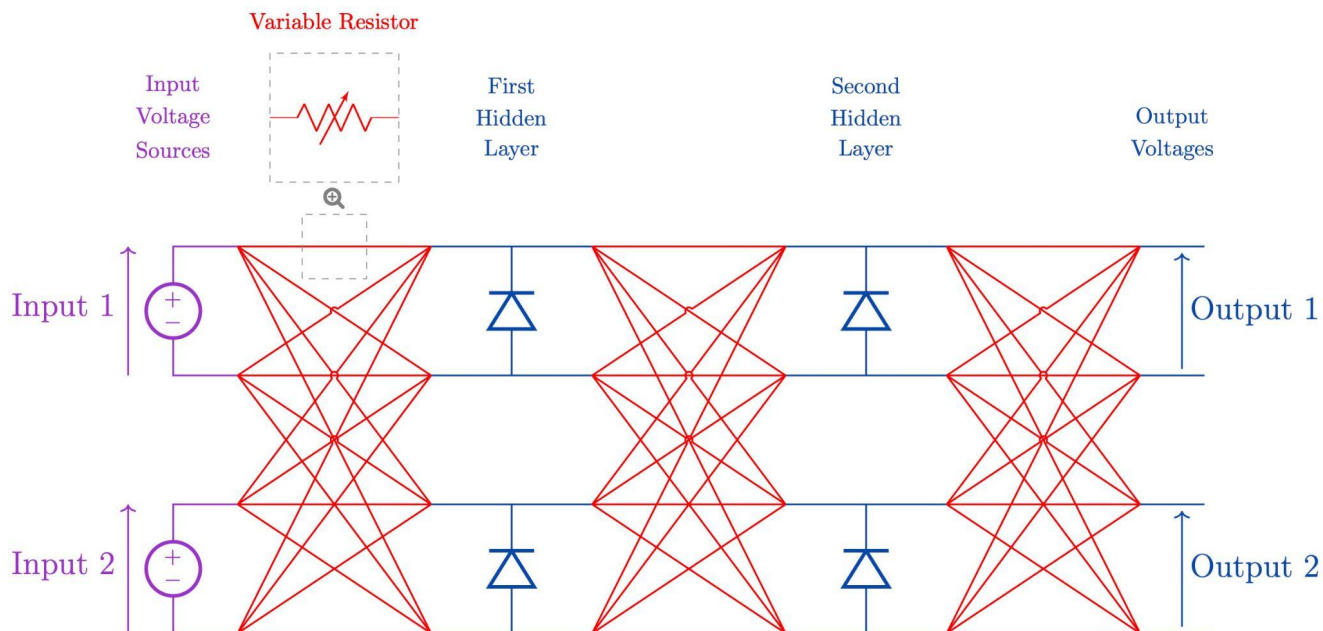
- voltage sources (inputs)
- variable resistors (trainable weights)
- diodes (nonlinearities)
- VCVS (amplification)

Such electrical circuits are
universal function approximators



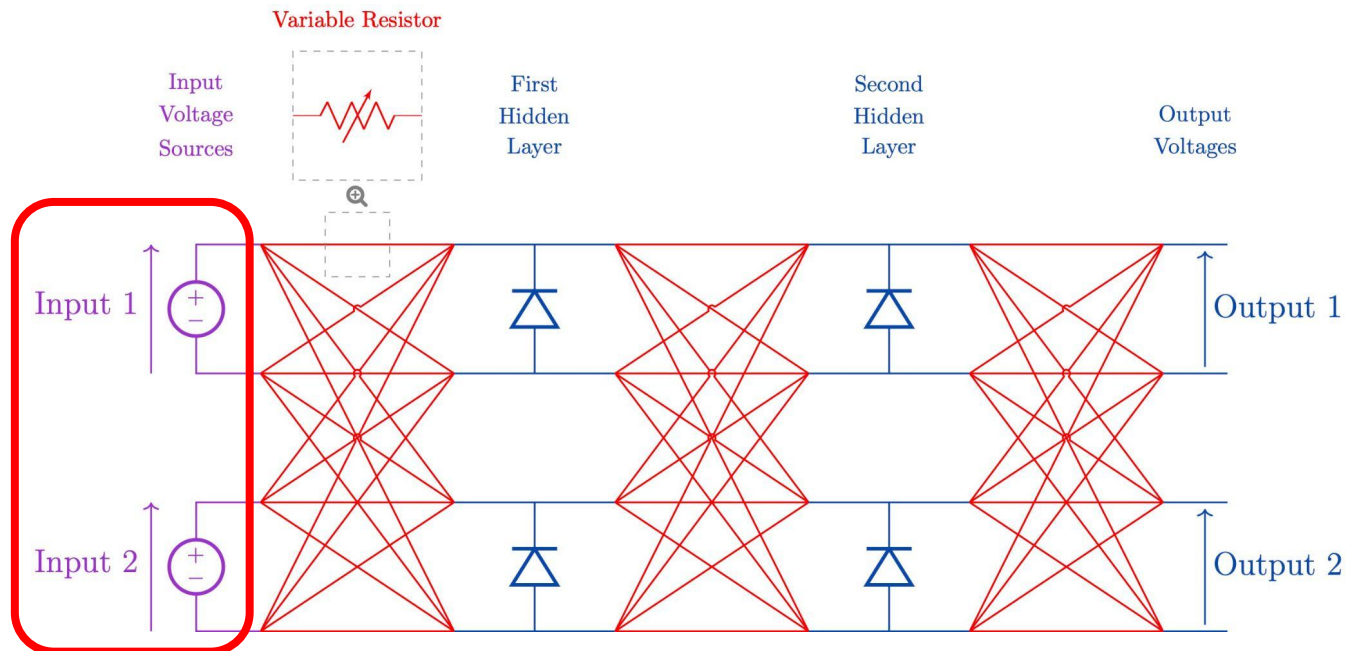
$$\text{output} = F_{\text{conductances}}(\text{input})$$

Deep Resistive Network



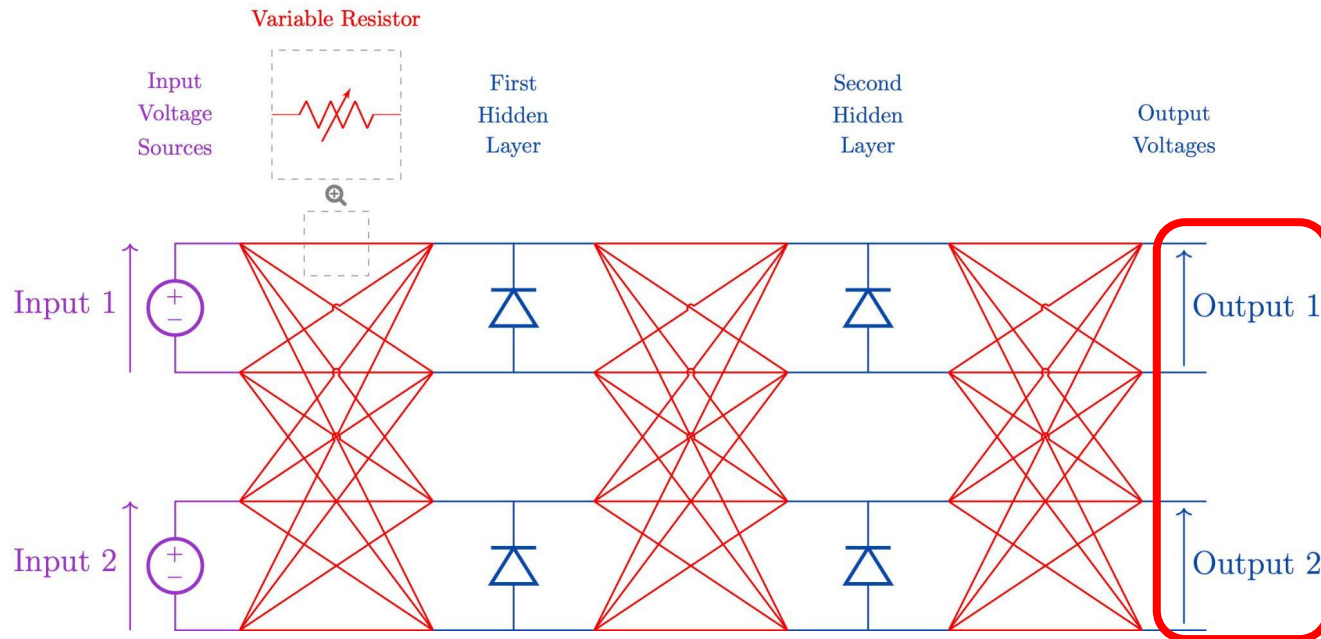
Deep Resistive Network

- input voltage sources



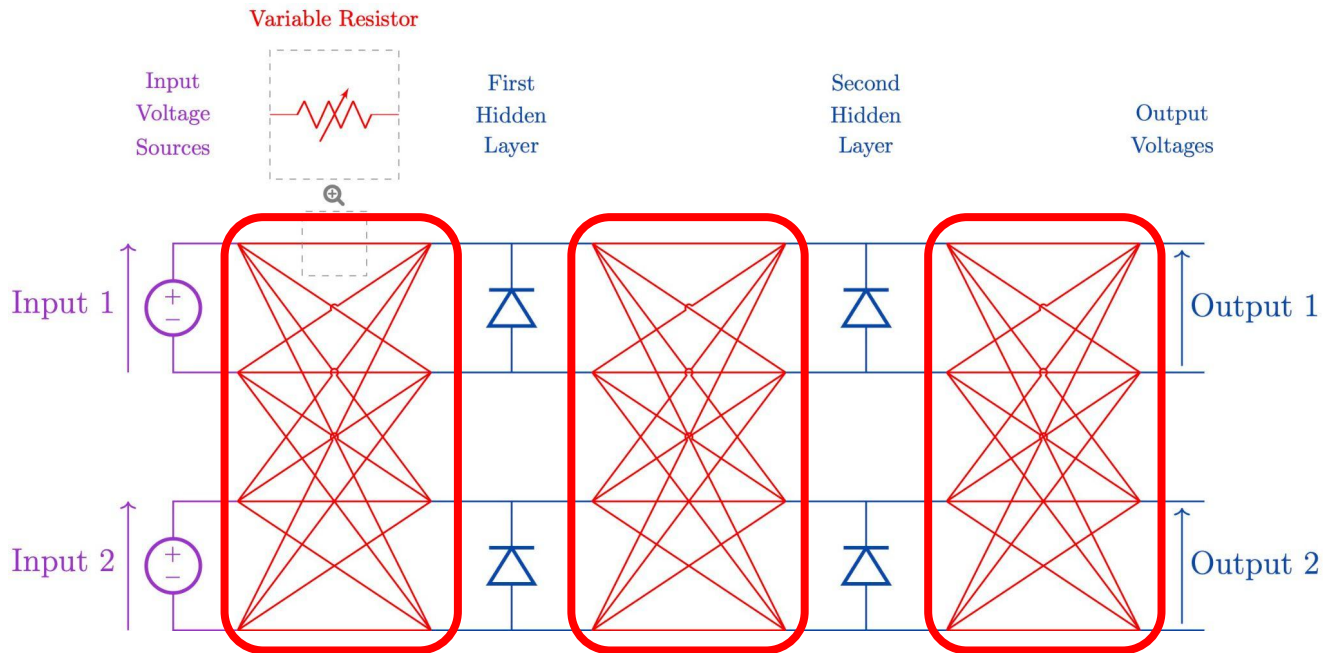
Deep Resistive Network

- input voltage sources
- output voltages



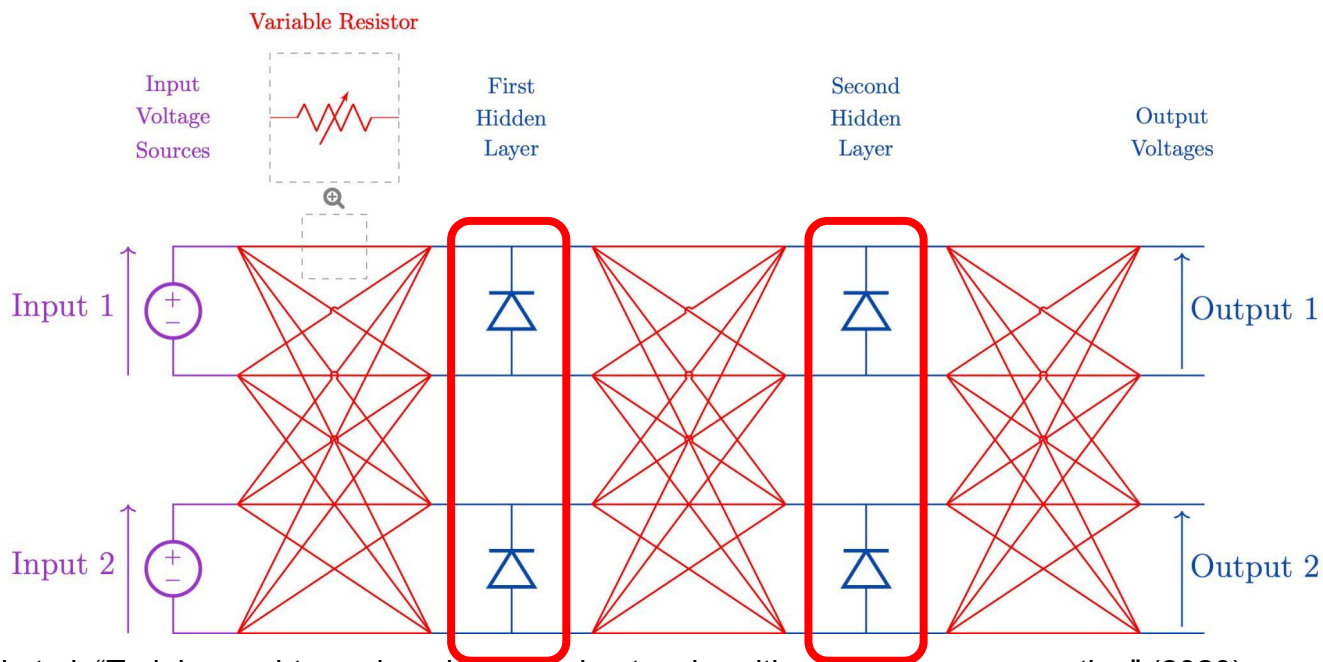
Deep Resistive Network

- input voltage sources
- output voltages
- crossbar arrays of variable resistors (linear trainable weights)



Deep Resistive Network

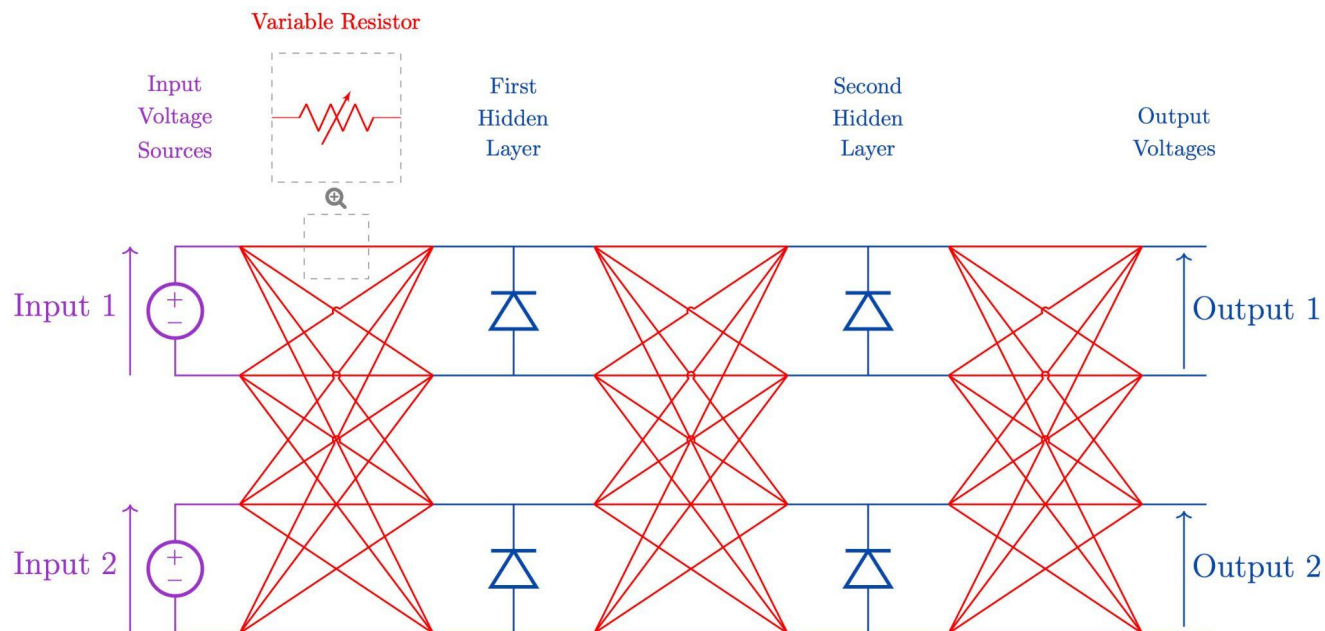
- input voltage sources
- output voltages
- crossbar arrays of variable resistors (linear trainable weights)
- diodes (nonlinearities)



Deep Resistive Network

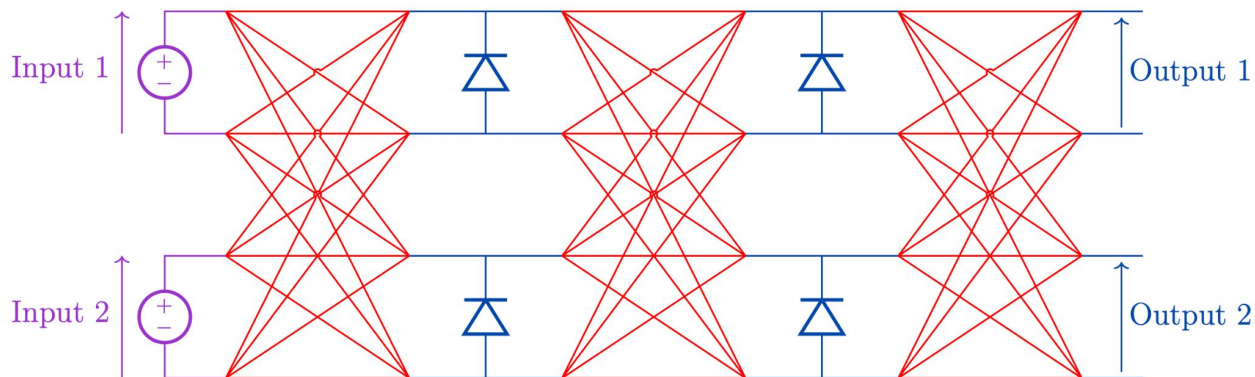
Properties:

- universal function approximator
- can be trained via equilibrium propagation



What is equilibrium propagation?

Equilibrium Propagation in a Deep Resistive Network



References:

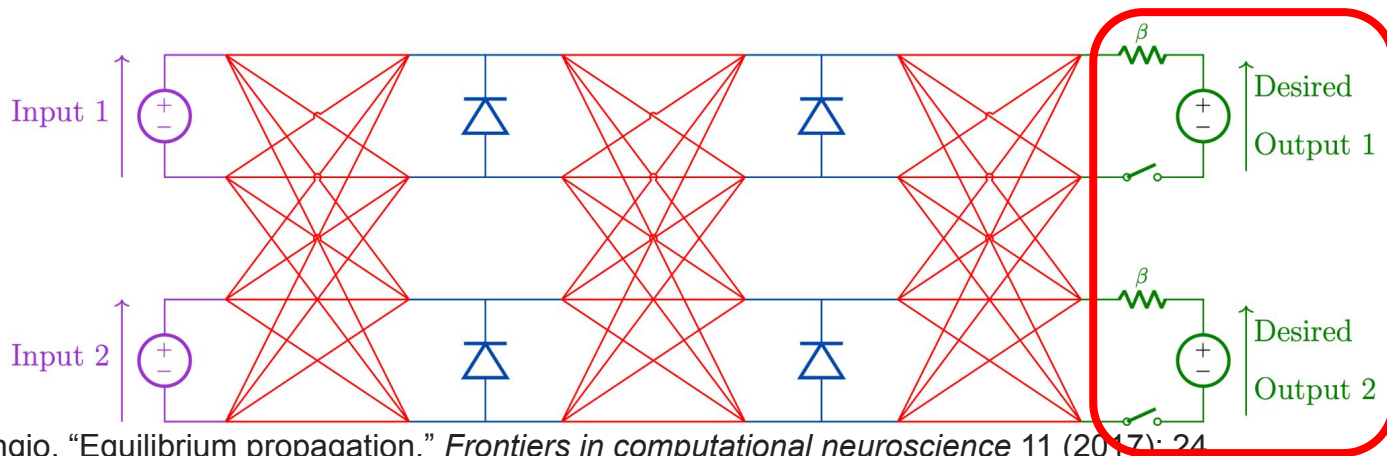
Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. "Training end-to-end analog neural networks with equilibrium propagation" (2020)

Equilibrium Propagation in a Deep Resistive Network

EP learning requires augmenting the network:

- for each pair of output nodes, add switch + voltage source + resistor



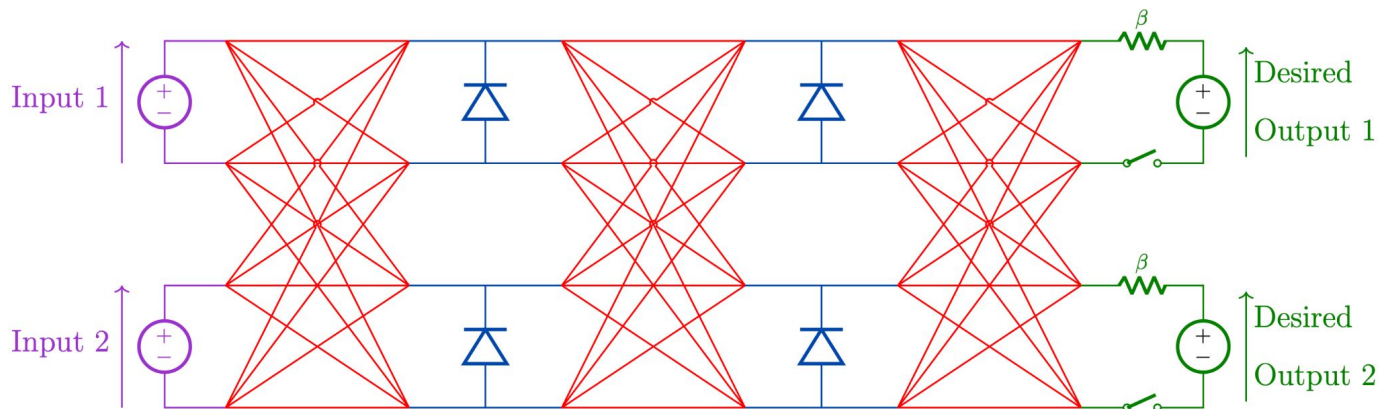
References:

Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24

Kendall et al. "Training end-to-end analog neural networks with equilibrium propagation" (2020)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:



References:

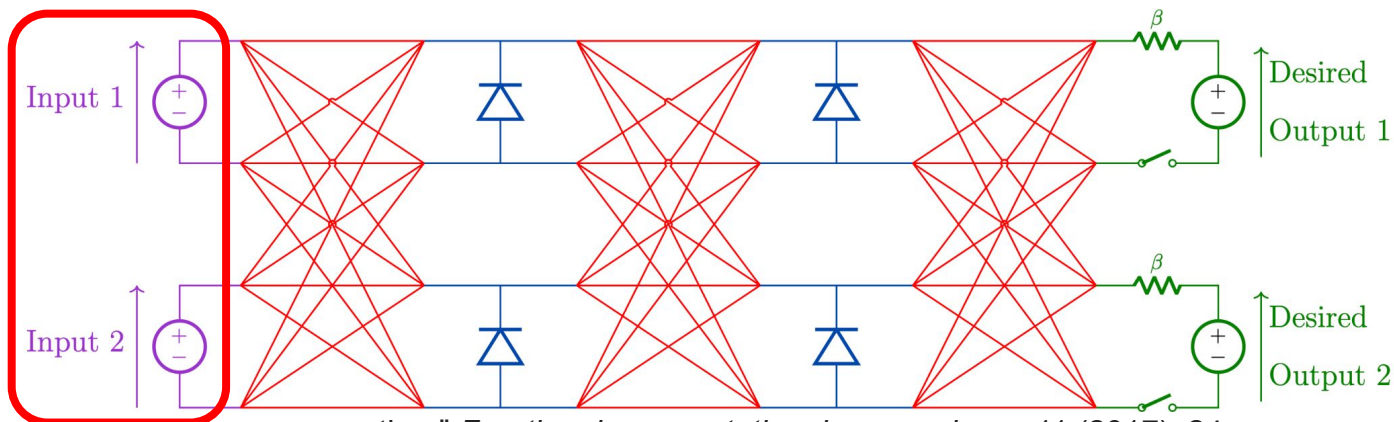
Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. "Training end-to-end analog neural networks with equilibrium propagation" (2020)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Set input voltages.



References:

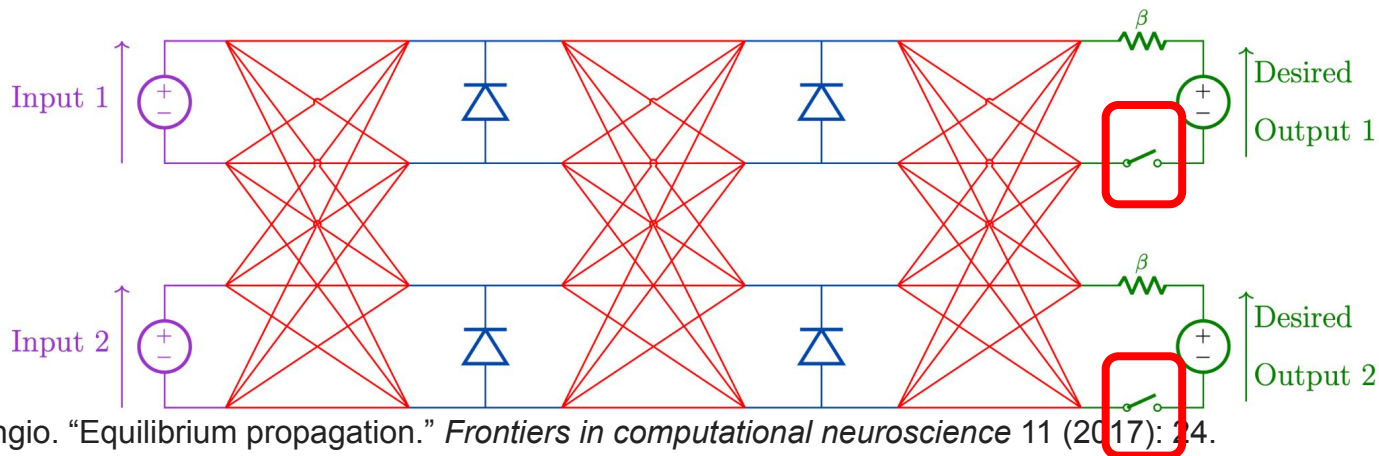
Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. "Training end-to-end analog neural networks with equilibrium propagation" (2020)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches.



References:

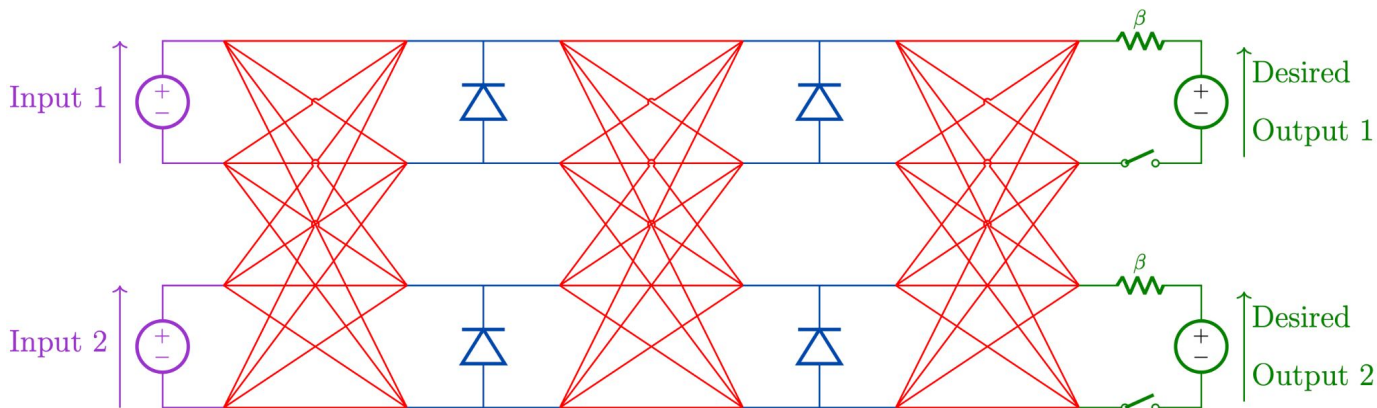
Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. "Training end-to-end analog neural networks with equilibrium propagation" (2020)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches. Observe the steady state (“free state”).



References:

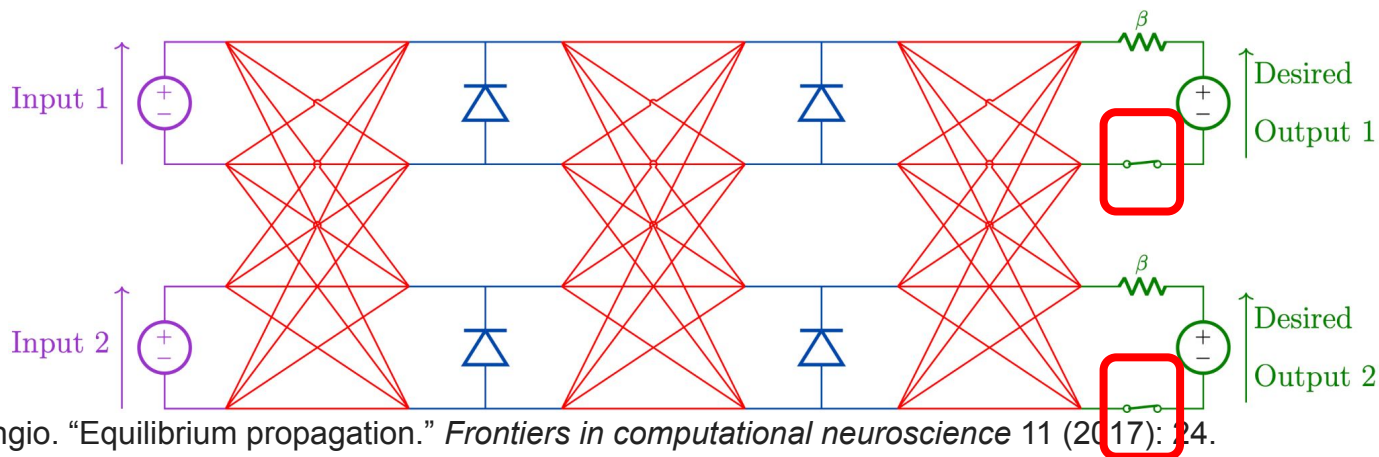
Scellier and Bengio. “Equilibrium propagation.” *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. “Training end-to-end analog neural networks with equilibrium propagation” (2020)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches: “free state”
2. Close output switches.



References:

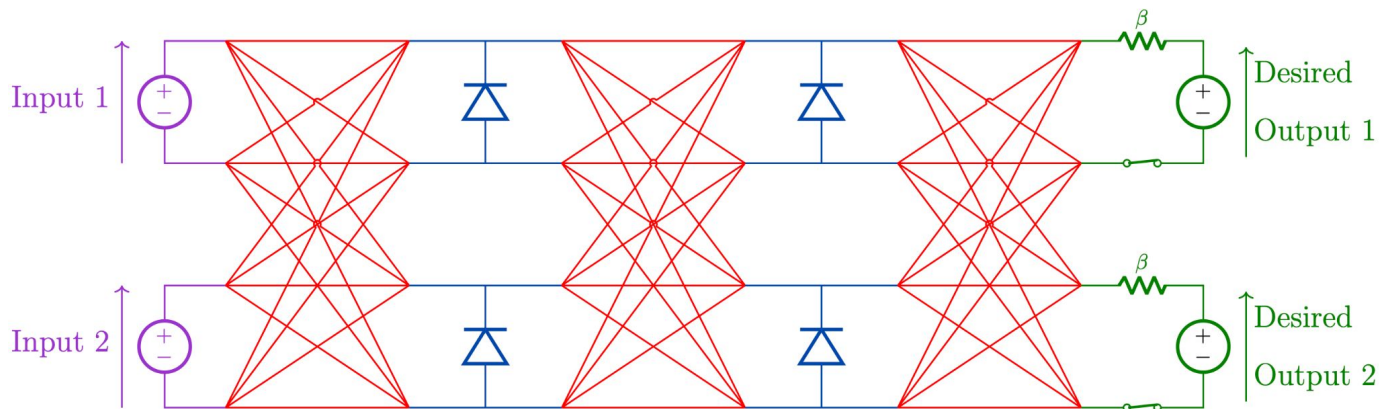
Scellier and Bengio. “Equilibrium propagation.” *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. “Training end-to-end analog neural networks with equilibrium propagation” (2020)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches: “free state”
2. Close output switches. Observe the new steady state (“nudged state”).



References:

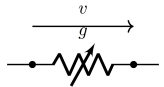
Scellier and Bengio. “Equilibrium propagation.” *Frontiers in computational neuroscience* 11 (2017): 24.

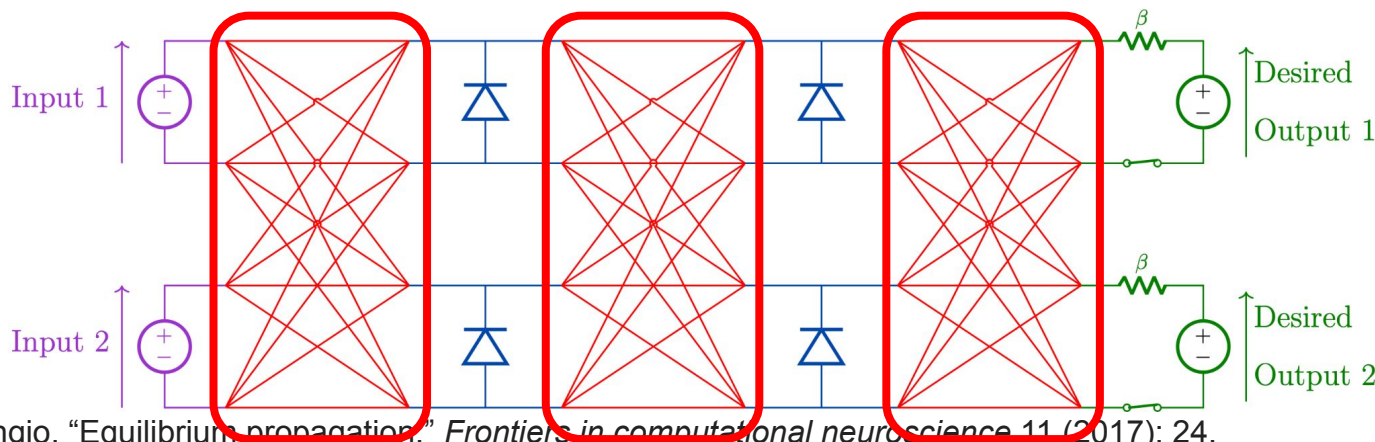
Kendall et al. “Training end-to-end analog neural networks with equilibrium propagation” (2020)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches: “free state”
2. Close output switches: “nudged state”

3. For each resistor  update the conductance g $\Delta g = \frac{\eta}{2\beta} \left(v_{\text{free}}^2 - v_{\text{nudged}}^2 \right)$



References:

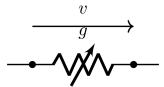
Scellier and Bengio. “Equilibrium propagation.” *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. “Training end-to-end analog neural networks with equilibrium propagation” (2020)

Equilibrium Propagation in a Deep Resistive Network

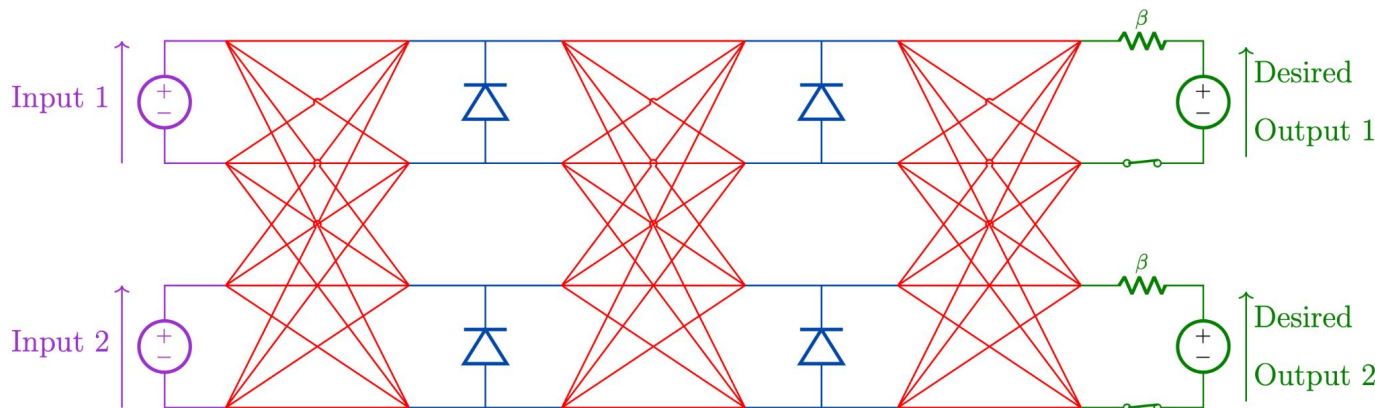
Training procedure:

1. Open output switches: “free state”
2. Close output switches: “nudged state”

3. For each resistor  update the conductance g

$$\Delta g = \frac{\eta}{2\beta} \left(v_{\text{free}}^2 - v_{\text{nudged}}^2 \right)$$

local in space



References:

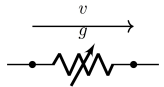
Scellier and Bengio. “Equilibrium propagation.” *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. “Training end-to-end analog neural networks with equilibrium propagation” (2020)

Equilibrium Propagation in a Deep Resistive Network

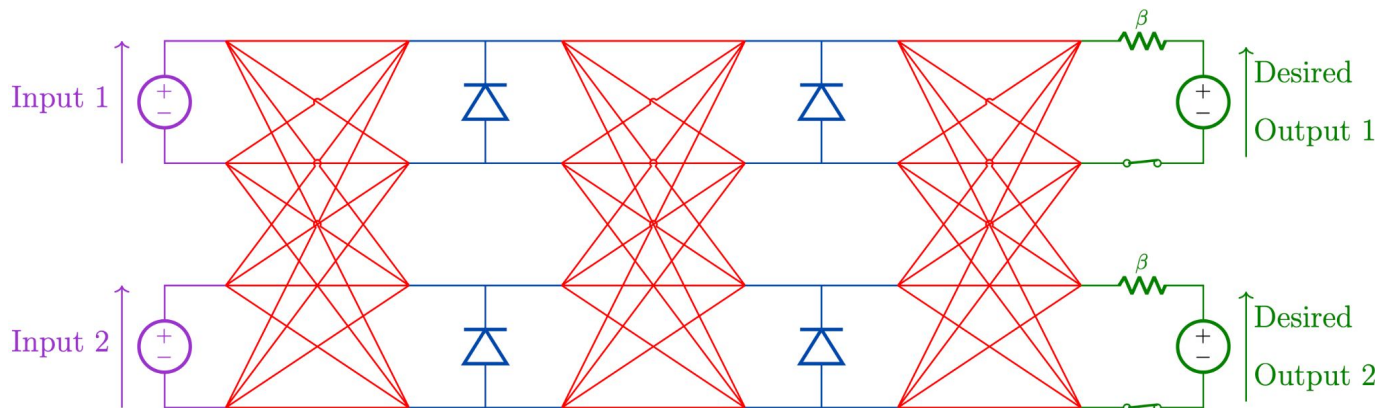
Training procedure:

1. Open output switches: “free state”
2. Close output switches: “nudged state”

3. For each resistor  update the conductance g $\Delta g = \frac{\eta}{2\beta} \left(v_{\text{free}}^2 - v_{\text{nudged}}^2 \right)$

Theorem: weight updates approximate gradient descent on the MSE

$$\Delta g = -\eta \nabla_g \text{MSE} + O(\beta)$$



References:

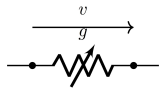
Scellier and Bengio. “Equilibrium propagation.” *Frontiers in computational neuroscience* 11 (2017): 24.

Kendall et al. “Training end-to-end analog neural networks with equilibrium propagation” (2020)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches: “free state”
2. Close output switches: “nudged state”

3. For each resistor  update the conductance g $\Delta g = \frac{\eta}{2\beta} \left(v_{\text{free}}^2 - v_{\text{nudged}}^2 \right)$

Theorem: weight updates approximate gradient descent on the MSE

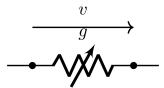
$$\Delta g = -\eta \nabla_g \text{MSE} + O(\beta)$$

Bottom line: **local learning rule** performs **stochastic gradient descent (SGD)**

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches: “free state”
2. Close output switches: “nudged state”

3. For each resistor  update the conductance g $\Delta g = \frac{\eta}{2\beta} \left(v_{\text{free}}^2 - v_{\text{nudged}}^2 \right)$

Theorem: weight updates approximate gradient descent on the cost function

$$\Delta g = -\eta \nabla_g C + O(\beta)$$

Bottom line: **local learning rule** performs **stochastic gradient descent (SGD)**

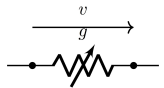
Remarks:

1. applicable to any cost function C (not just the MSE)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches: “free state”
2. Close output switches: “nudged state”

3. For each resistor  update the conductance g $\Delta g = \frac{\eta}{2\beta} \left(v_{\text{free}}^2 - v_{\text{nudged}}^2 \right)$

Theorem: weight updates approximate gradient descent on the cost function

$$\Delta g = -\eta \nabla_g C + O(\beta)$$

Bottom line: **local learning rule** performs **stochastic gradient descent (SGD)**

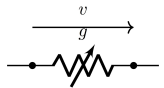
Remarks:

1. applicable to any cost function C (not just the MSE)
2. applicable to any network topology (not just a DRN)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches: “free state”
2. Close output switches: “nudged state”

3. For each resistor  update the conductance g $\Delta g = \frac{\eta}{2\beta} \left(v_{\text{free}}^2 - v_{\text{nudged}}^2 \right)$

Theorem: weight updates approximate gradient descent on the cost function

$$\Delta g = -\eta \nabla_g C + \mathcal{O}(\beta)$$

Bottom line: **local learning rule** performs **stochastic gradient descent (SGD)**

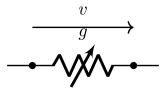
Remarks:

1. applicable to any cost function C (not just the MSE)
2. applicable to any network topology (not just a DRN)

Equilibrium Propagation in a Deep Resistive Network

Training procedure:

1. Open output switches: “free state”
2. Close output switches: “nudged state”

3. For each resistor  update the conductance g $\Delta g = \frac{\eta}{2\beta} \left(v_{\text{free}}^2 - v_{\text{nudged}}^2 \right)$

Theorem: weight updates perform gradient descent on a surrogate function

$$\Delta g = -\eta \nabla_g \mathcal{L}_\beta$$

where

$$\mathcal{L}_\beta = C + O(\beta)$$

Bottom line: **local learning rule** performs **stochastic gradient descent (SGD)**

Remarks:

1. applicable to any cost function C (not just the MSE)
2. applicable to any network topology (not just a DRN)

Equilibrium Propagation in a Deep Resistive Network

Sketch of the proof:

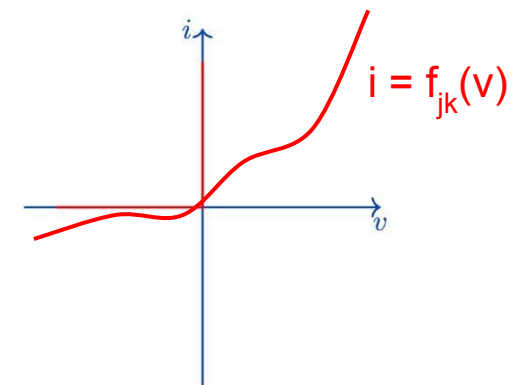
- Equilibrium propagation generally applies to systems whose steady state minimizes a functional (variational principle) [1,2]
- Nonlinear resistive networks minimize a functional called *co-content* [3]

References

- [1] Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.
- [2] Scellier. "A deep learning theory for neural networks grounded in physics." *PhD thesis, Université de Montréal* (2021).
- [3] Millar. "CXVI. Some general theorems for non-linear systems possessing resistance." *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* 42.333 (1951): 1150-1160.

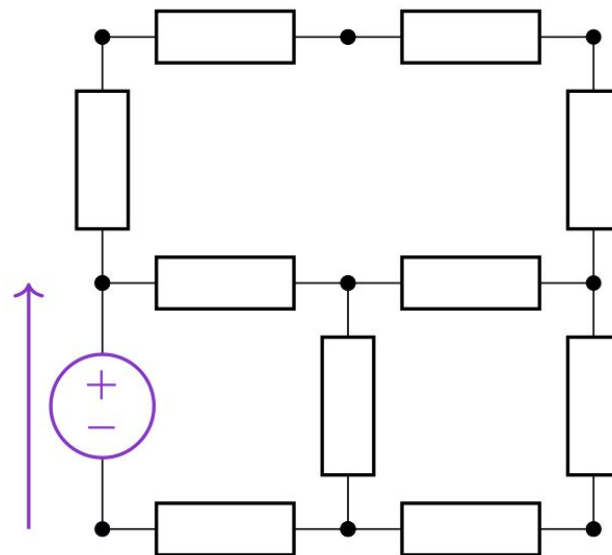
What is the co-content?

i-v curve of branch (j,k)



Co-content of branch (j,k)

$$E_{jk} = \int_0^{v_{jk}} f_{jk}(u) du$$



Total co-content of the circuit

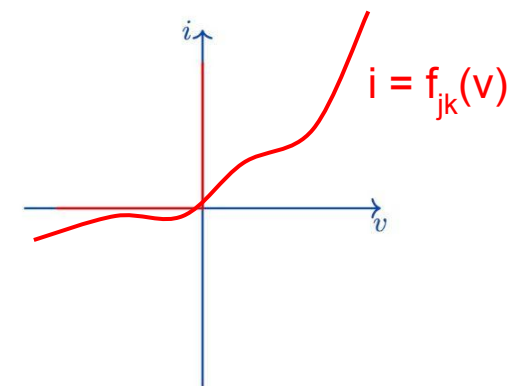
$$E_{\text{total}} = \sum_{j,k} E_{jk}$$

References

Millar. "CXVI. Some general theorems for non-linear systems possessing resistance." The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science 42.333 (1951): 1150-1160.

Kendall et al. "Training end-to-end analog neural networks with equilibrium propagation" (2020)

i-v curve of branch (j,k)

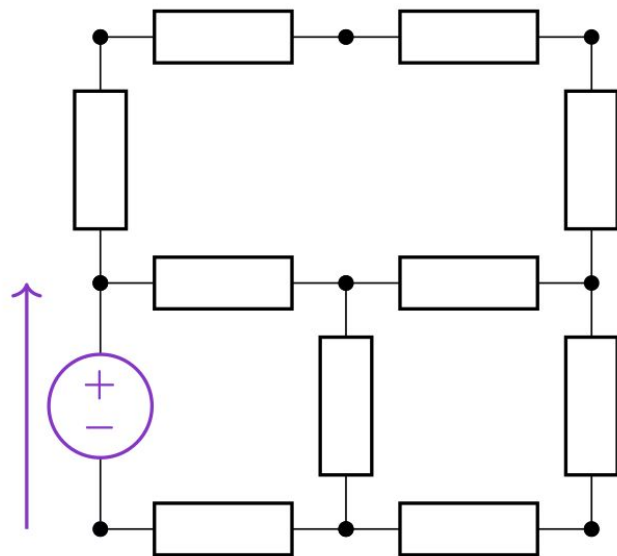


Co-content of branch (j,k)

$$E_{jk} = \int_0^{v_{jk}} f_{jk}(u) du$$

branch equation + KVL

$$\delta E_{\text{total}} = \sum_{j,k} \delta E_{jk} = \sum_{j,k} f_{jk} \delta v_{jk} = \sum_{j,k} i_{jk} (\delta v_j - \delta v_k) = \sum_j \delta v_j \left(\sum_k i_{jk} \right) \stackrel{\text{KCL}}{=} 0$$

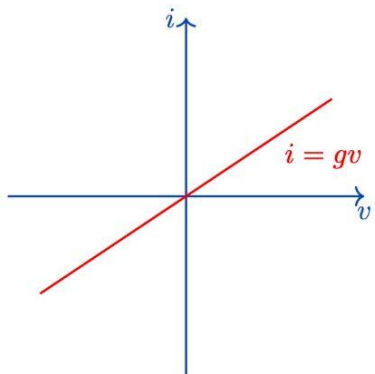


Total co-content of the circuit

$$E_{\text{total}} = \sum_{j,k} E_{jk}$$

KCL

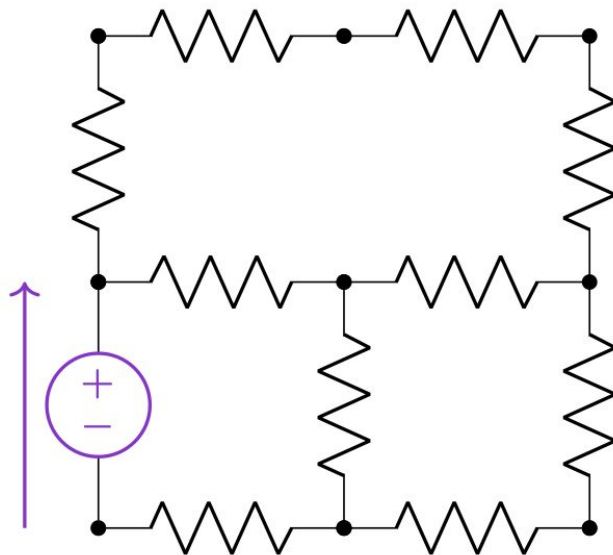
Linear resistor



Co-content of linear resistor

$$E = \frac{1}{2}gv^2$$

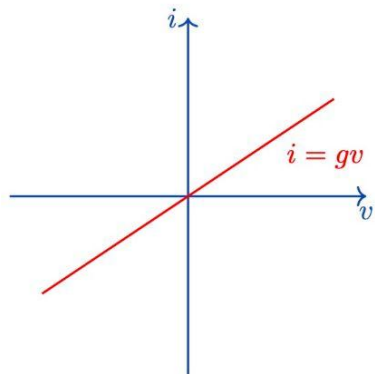
(half the power dissipation)



Total co-content of the circuit

$$E_{\text{total}} = \frac{1}{2} \sum_{j,k} g_{jk} v_{jk}^2$$

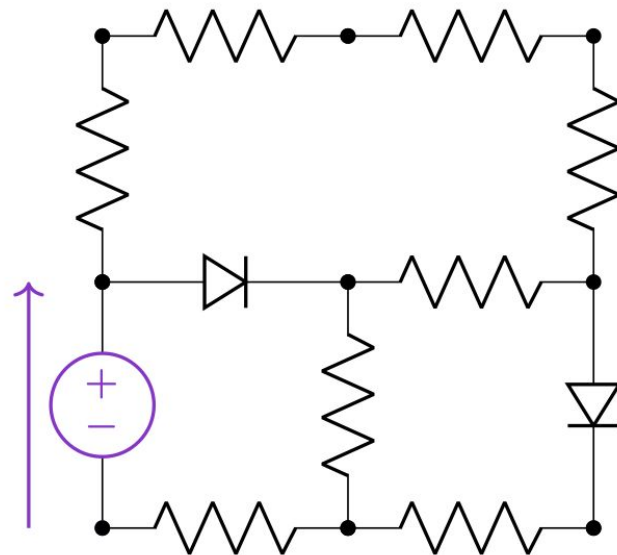
Linear resistor



Co-content of linear resistor

$$E = \frac{1}{2}gv^2$$

(half the power dissipation)



Total co-content of the circuit

$$E_{\text{total}} = \frac{1}{2} \sum_{j,k} g_{jk} v_{jk}^2 + E_{\text{diodes}}$$

Equilibrium Propagation Learning Rule

Learning rule

$$\Delta g = \frac{\eta}{2\beta} \left(\underbrace{v_{\text{free}}^2}_{\partial_g E(\text{free})} - \underbrace{v_{\text{nudged}}^2}_{\partial_g E(\text{nudged})} \right)$$

Co-content

$$E = \frac{1}{2} \sum_{\text{resistor}} gv^2 + E_{\text{diodes}}$$

Equilibrium Propagation Learning Rule

Learning rule

$$\Delta g = \frac{\eta}{\beta} (\partial_g E(\text{free}) - \partial_g E(\text{nudged}))$$

Co-content

$$E = \frac{1}{2} \sum_{\text{resistor}} gv^2 + E_{\text{diodes}}$$

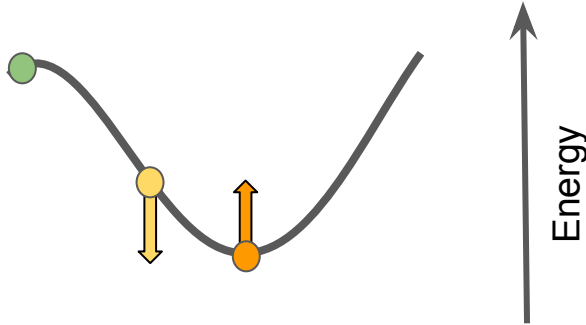
Equilibrium Propagation Learning Rule

Learning rule

$$\Delta w = \frac{\eta}{\beta} (\partial_w E(\text{free}) - \partial_w E(\text{nudged}))$$

Equilibrium Propagation Learning Rule

- Desired output
- Nudged output
- Free output

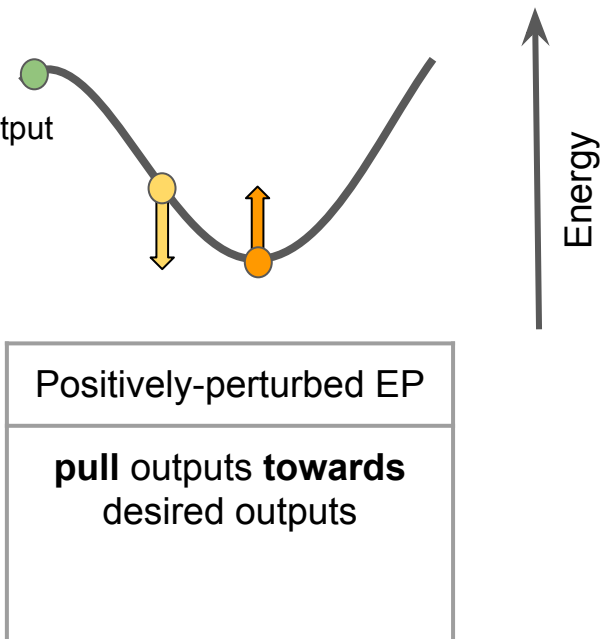


Learning rule

$$\Delta w = \frac{\eta}{\beta} (\partial_w E(\text{free}) - \partial_w E(\text{nudged}))$$

Variants of Equilibrium Propagation (EP)

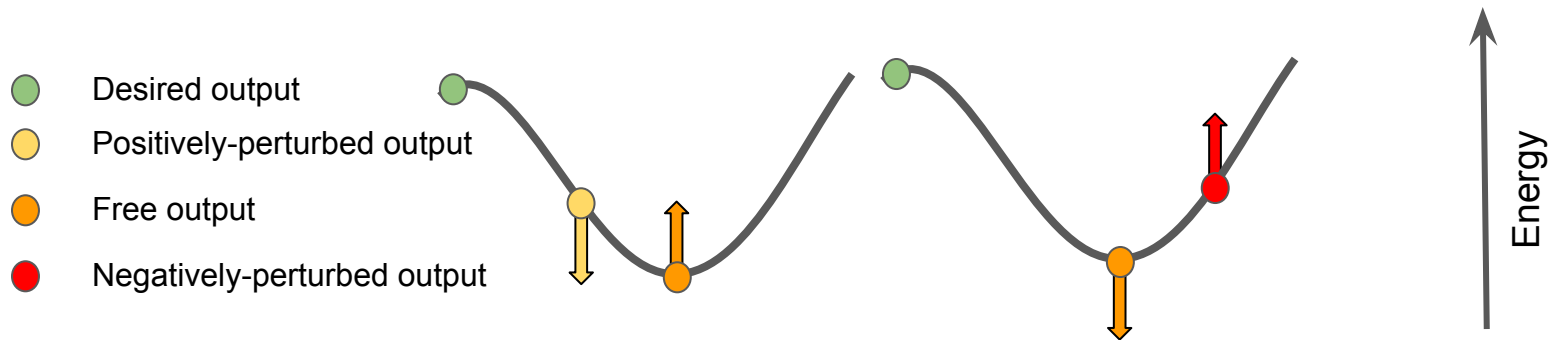
- Desired output
- Positively-perturbed output
- Free output



References:

- Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.
- Laborieux et al. "Scaling equilibrium propagation to deep ConvNets by reducing its gradient estimator bias" (2021)
- Scellier et al. "Energy-based learning algorithms: a comparative study" NeurIPS (2023)

Variants of Equilibrium Propagation (EP)



Positively-perturbed EP	Negatively-perturbed EP
pull outputs towards desired outputs	push outputs away from desired outputs

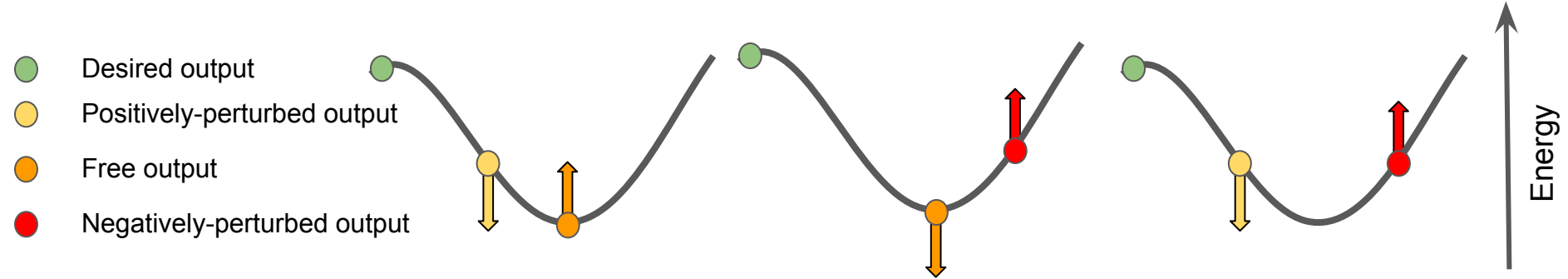
References:

Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.

Laborieux et al. "Scaling equilibrium propagation to deep ConvNets by reducing its gradient estimator bias" (2021)

Scellier et al. "Energy-based learning algorithms: a comparative study" NeurIPS (2023)

Variants of Equilibrium Propagation (EP)

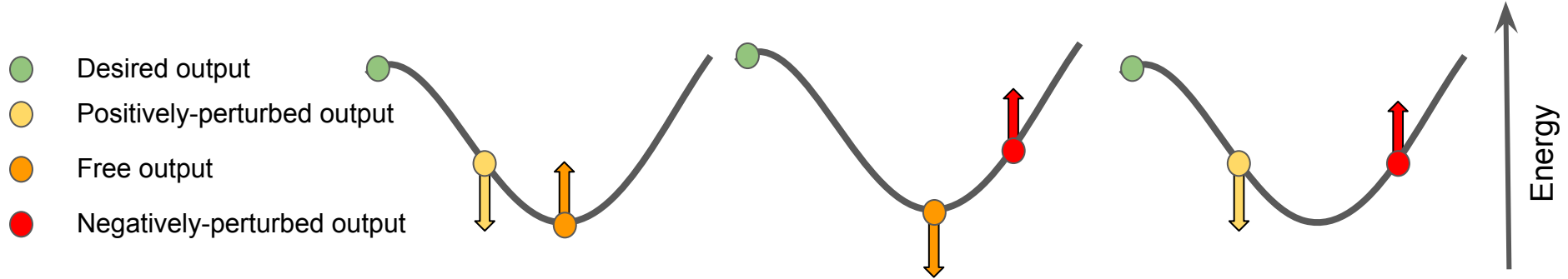


Positively-perturbed EP	Negatively-perturbed EP	Centered EP
pull outputs towards desired outputs	push outputs away from desired outputs	→ one positive perturbation → one negative perturbation

References:

- Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.
- Laborieux et al. "Scaling equilibrium propagation to deep ConvNets by reducing its gradient estimator bias" (2021)
- Scellier et al. "Energy-based learning algorithms: a comparative study" NeurIPS (2023)

Variants of Equilibrium Propagation (EP)

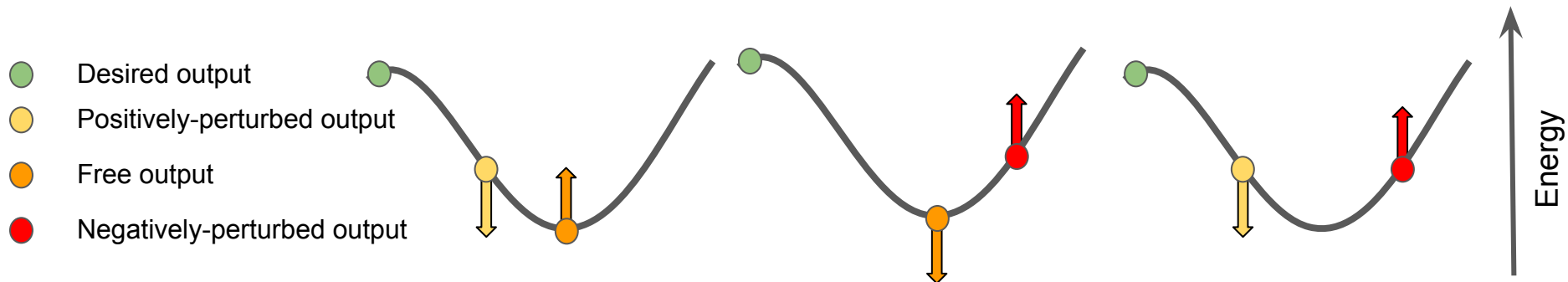


Positively-perturbed EP	Negatively-perturbed EP	Centered EP
pull outputs towards desired outputs	push outputs away from desired outputs	→ one positive perturbation → one negative perturbation

Performance in practice?



Variants of Equilibrium Propagation (EP)



Positively-perturbed EP

Negatively-perturbed EP

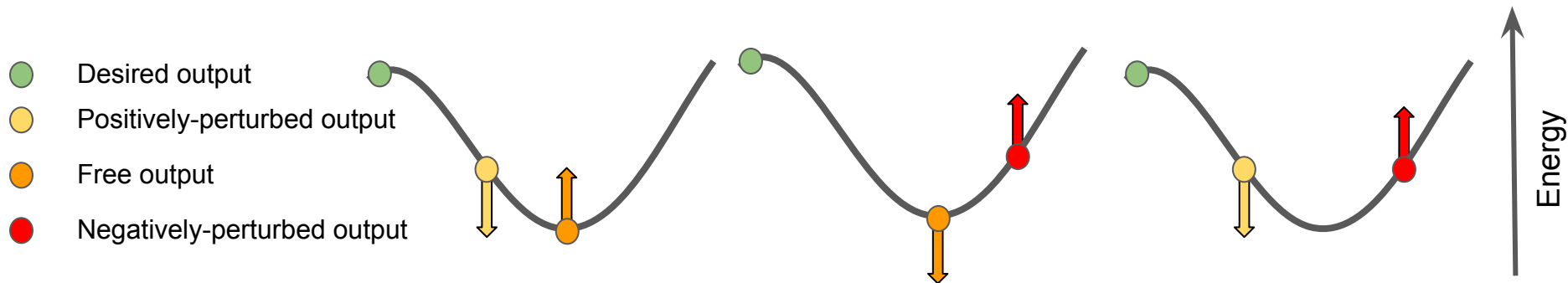
Centered EP

Theorem: EP computes the gradient of a surrogate function that approximates the cost function

$$\Delta w = -\eta \nabla_w \mathcal{L}_\beta$$

$$\mathcal{L}_\beta = C + O(\beta)$$

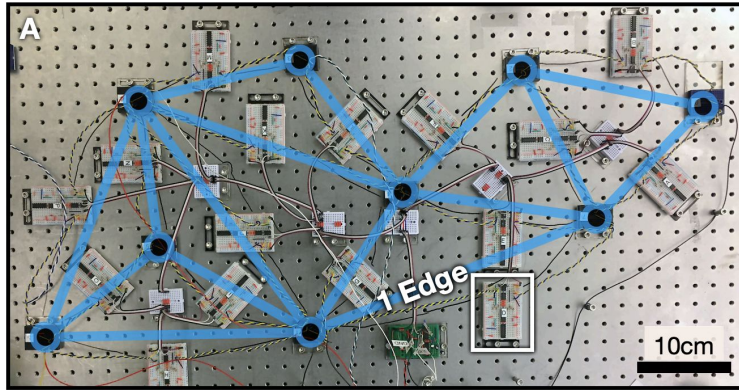
Variants of Equilibrium Propagation (EP)



	Positively-perturbed EP	Negatively-perturbed EP	Centered EP
Properties of the surrogate function	$L_{P-EP} = C + O(\beta)$ $L_{P-EP} \leq C$ <p>(lower bound)</p>	$L_{N-EP} = C + O(\beta)$ $L_{N-EP} \geq C$ <p>(upper bound)</p>	$L_{C-EP} = C + O(\beta^2)$

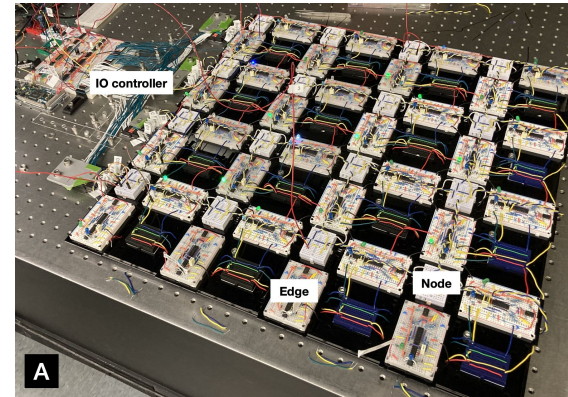
Hardware Experiments

Linear resistor network [1]



- Network of 9 nodes and 16 edges
- Task: classify the Iris dataset

Nonlinear resistive network [2]



- Network of 32 (transistor-based) twin edges
- Tasks: XOR and nonlinear regression

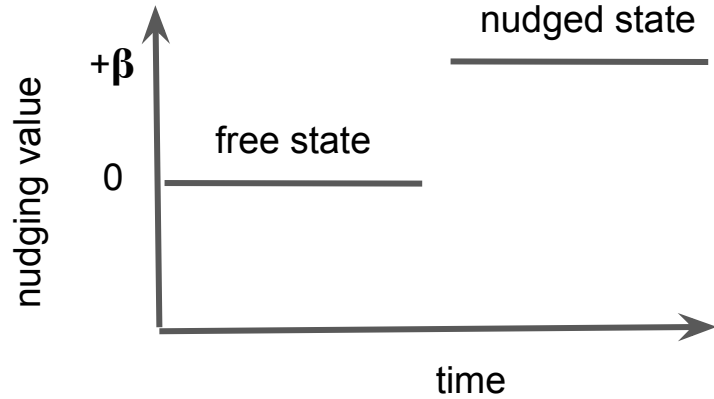
Caveat: use two copies of the same network
(one for the free state, one for the nudge state)

References

[1] Dillavou et al. "Demonstration of decentralized physics-driven learning." Physical Review Applied (2022)

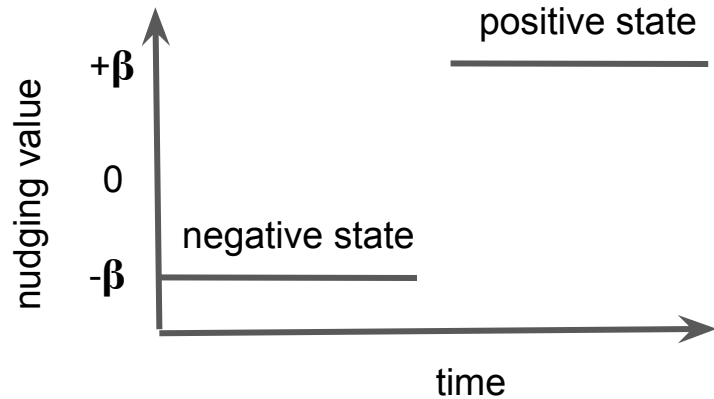
[2] Dillavou et al. "Circuits that train themselves: decentralized, physics-driven learning." (2023)

Equilibrium propagation



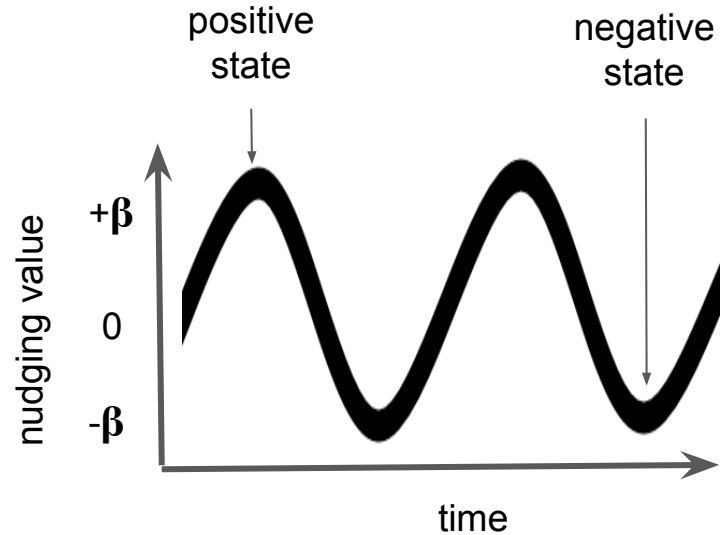
Learning rule $\Delta g = \frac{\eta}{2\beta} (v_{\text{free}}^2 - v_{\text{nudged}}^2)$

Equilibrium propagation



Learning rule
$$\Delta g = \frac{\eta}{4\beta} (v_{\text{neg}}^2 - v_{\text{pos}}^2)$$
$$= \frac{\eta}{4\beta} (v_{\text{neg}} + v_{\text{pos}}) (v_{\text{neg}} - v_{\text{pos}})$$

Frequency propagation [1]



$$\text{state}(t) = \text{mean} + \text{amplitude} * \sin(\omega t) + O(\beta^2)$$

Reference

[1] Anisetti et al. "Frequency propagation: Multi-mechanism learning in nonlinear physical networks." (2022)

Simulations of deep resistive networks

Task: train a DRN to classify MNIST digits

	Network size (number of parameters)	Number of epochs	Total duration	Test error rate
SPICE simulations [1]	0.16M	10	1 week	3.43%

References

[1] Kendall et al. "Training end-to-end analog neural networks with equilibrium propagation" (2020)

Simulations of deep resistive networks

Task: train a DRN to classify MNIST digits

	Network size (number of parameters)	Number of epochs	Total duration	Test error rate
SPICE simulations [1]	0.16M	10	1 week	3.43%
DRN simulator [2,3]	51M	50	6 hours	1.40%

320x larger

5x more

28x shorter

The “network size” to “epoch duration” ratio is 45000 times larger

References

[1] Kendall et al. “Training end-to-end analog neural networks with equilibrium propagation” (2020)

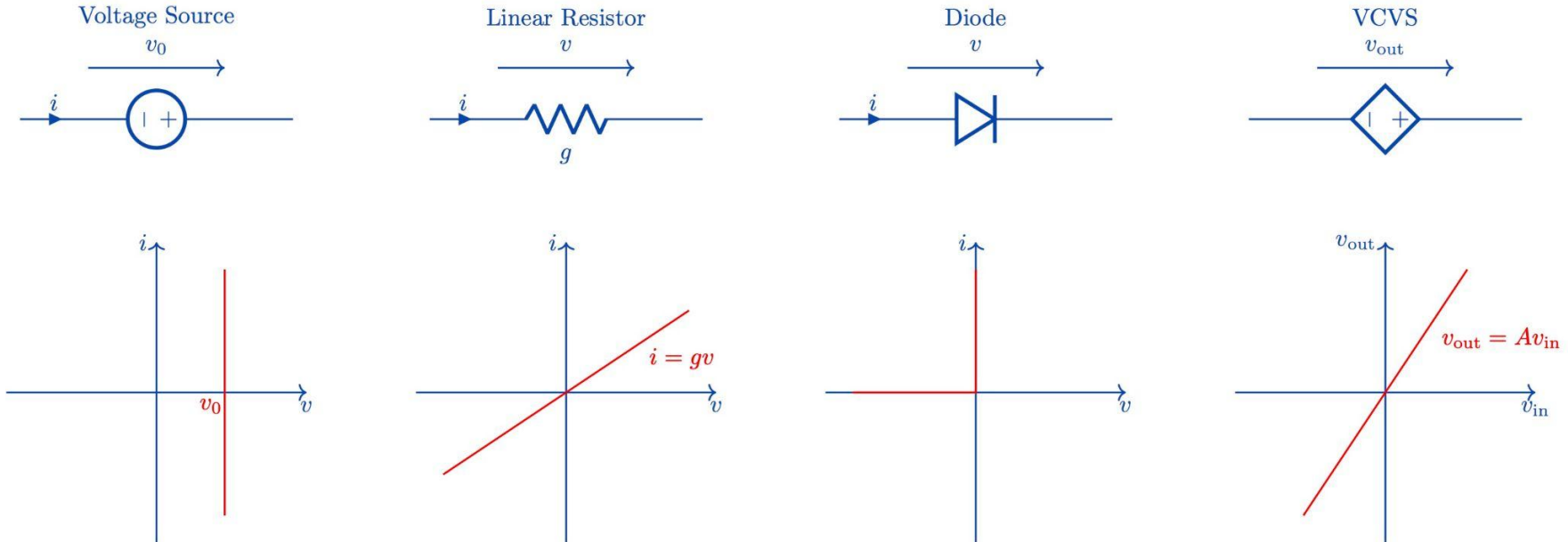
[2] “A fast algorithm to simulate nonlinear resistive networks”. To appear.

[3] Code will be available at <https://github.com/rain-neuromorphics/energy-based-learning>

Simulations of deep resistive networks

Main idea:

- assumption: the circuit elements are **ideal**



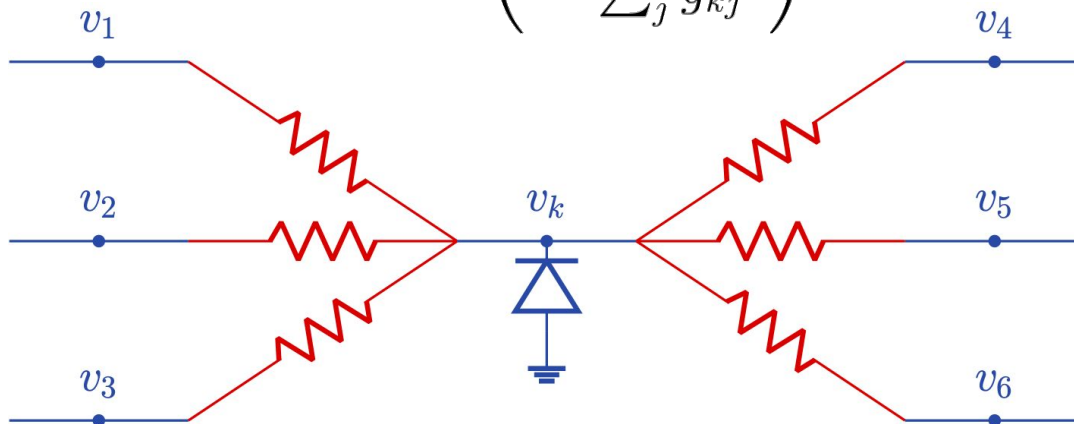
Simulations of deep resistive networks

Main idea:

- assumption: the circuit elements are **ideal**
- the math considerably simplifies

Equation of a node voltage

$$v_k = \max \left(0, \frac{\sum_j g_{kj} v_j}{\sum_j g_{kj}} \right)$$



Simulator for energy-based algorithms

Three key abstractions

Energy function

- deep resistive network (DRN)
- ...

Energy minimizer

(in the space of network configurations)

- block coordinate descent
- ...

Learning algorithm

(in the weight space)

- equilibrium propagation (positive, negative, centered, ...)
- ...

Simulator for energy-based algorithms

Three key abstractions

Energy function

- deep resistive network (DRN)
- deep Hopfield network (DHN)
- ...

Energy minimizer

(in the space of network configurations)

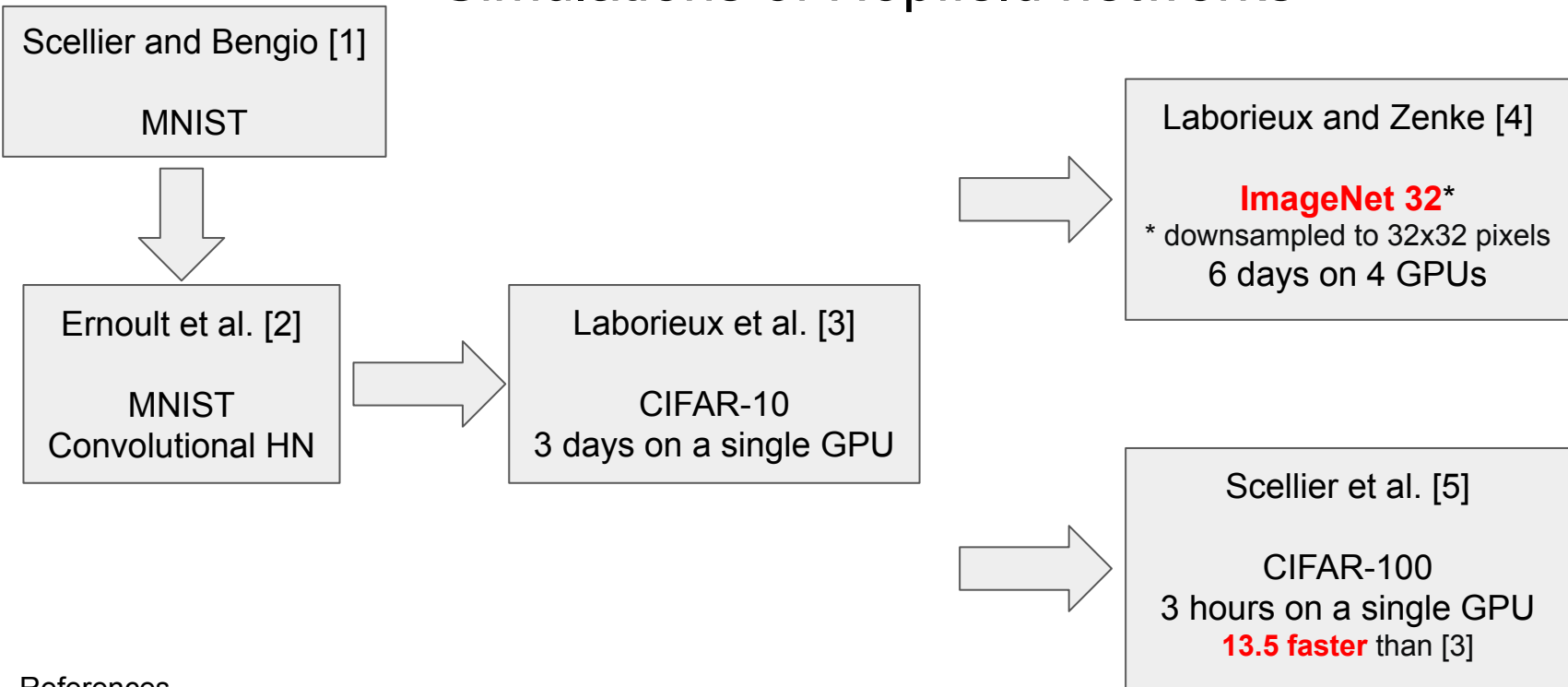
- block coordinate descent
- gradient descent
- ...

Learning algorithm

(in the weight space)

- equilibrium propagation (positive, negative, centered, ...)
- contrastive (Hebbian) learning
- coupled learning
- ...
- truncated backpropagation (baseline)
- recurrent backpropagation (baseline)

Simulations of Hopfield networks



References

- [1] Scellier and Bengio. "Equilibrium propagation." *Frontiers in computational neuroscience* 11 (2017): 24.
- [2] Ernoul et al. "Updates of equilibrium prop match gradients of backprop through time." *NeurIPS* (2019).
- [3] Laborieux et al. "Scaling equilibrium propagation to deep ConvNets." *Frontiers in neuroscience* (2021).
- [4] Laborieux and Zenke. "Holomorphic equilibrium propagation." *NeurIPS* (2022).
- [5] Scellier et al. "Energy-based learning algorithms: a comparative study." *NeurIPS* (2023)

Hardware experiments of Hopfield networks

nature electronics

Article

<https://doi.org/10.1038/s41928-022-00869-w>

Activity-difference training of deep neural networks using memristor crossbars

Received: 10 March 2022

Su-in Yi¹, Jack D. Kendall², R. Stanley Williams¹ & Suhas Kumar³✉

Accepted: 13 October 2022

- 64 × 64 memristor crossbar array
- Task: classify Braille words
- demonstrate **10,000x improvement** (over GPUs) in energy efficiency for training

Caveat: use additional (SRAM) memory to store the free and nudged states

Summary

Training deep resistive networks (DRNs) with equilibrium propagation (EP):

1. DRNs are universal function approximators
2. EP performs gradient descent on a cost function
3. DRN simulator is 45,000x faster than SPICE simulations
4. Small-scale experiments demonstrate 10,000x energy efficiency gains over GPUs

Thanks to all my collaborators!



Jack Kendall



Maxence Ernout



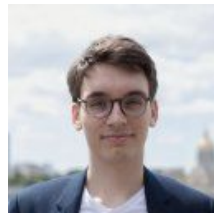
Suhas Kumar



Ross Pantone



Yoshua Bengio



Axel Laborieux



Julie Grollier



Damien Querlioz



Yann Ollivier



Siddhartha
Mishra



Kalpana
Manickavasagam



Vidyesh Anisetti



Ananth Kandala



Jennifer Schwarz



Thank you!

