
Adjoint Method: The Connection between Analog-based Equilibrium Propagation Architectures and Neural ODEs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Analog neural networks (ANNs) hold significant potential for substantial reductions
2 in power consumption in modern neural networks, particularly when employing
3 the increasingly popular Energy-Based Models (EBMs) in tandem with the local
4 Equilibrium Propagation (EP) training algorithm. This paper analyzes the relationship
5 between this family of ANNs and the concept of Neural Ordinary Differential
6 Equations (Neural ODEs). Using the adjoint method, we formally demonstrate
7 that ANN-EP can be derived from Neural ODEs by constraining the differential
8 equations to those with a steady-state response. This finding opens avenues
9 for the ANN-EP community to extend ANNs to non-steady-state scenarios.
10 Additionally, it provides an efficient setting for NN-ODEs that significantly
11 reduces the training cost.

12 1 Introduction

13 The aspiration to empower IoT devices at the edge with real-time adaptive learning capabilities has
14 been one of the primary motivators for advancing efficient neural network training methods. This
15 endeavor goes beyond merely substituting backpropagation with alternative strategies; it represents
16 a shift in the computational principles underlying model design, training, and inference. Such a
17 transformation requires the incorporation of innovative models like Energy-Based Models (EBMs),
18 which diverge significantly from conventional paradigms. This shift is pivotal for leveraging the
19 unique attributes of diverse computing paradigms, such as photonic and neuromorphic computing.

20 In light of the evolving computational landscape, the advent of Neural Ordinary Differential Equations
21 [1] (Neural ODEs) and Equilibrium Propagation [2] (EP) have emerged as promising contenders.
22 Neural ODEs extend neural network models into the continuous-time domain, allowing
23 the incorporation of differential equations to represent the evolution of the network states. EP, in
24 contrast, provides a biologically plausible learning framework that integrates the concept of a cost
25 function, a notable departure from the training methodologies of conventional EBMs.

26 The adjoint method serves as a pivotal link between EP and Neural ODEs, providing a framework
27 that allows for the efficient computation of gradients in continuous-time models. This paper posits
28 that EP can be construed as a special case of Neural ODEs, unified by the underlying adjoint method.
29 This connection is not merely theoretical but has practical implications, especially in the realm of
30 analog computation. By leveraging an analog circuit, we will demonstrate how the Lambda (λ)
31 variable in Neural ODEs corresponds precisely to the variation of the node voltages of the circuit.
32 The exploration of such connections is vital for the development of low-cost and efficient solutions
33 in neural networks, aligning with the broader goals of advancing biologically plausible and energy-
34 efficient machine learning models.

$\min_{\theta} C(\mathbf{y}, \mathbf{y}_{\text{true}})$	$\min_{\theta} C(\mathbf{s}, \mathbf{s}_{\text{true}})$	$\min_{\theta} C(\mathbf{s}, \mathbf{s}_{\text{true}})$
$\text{s.t. } \mathbf{y} = \mathbf{U}_n(\cdots(\mathbf{U}_0(\mathbf{x}, \theta)\cdots))$	$\text{s.t. } \mathbf{f}(\mathbf{s}, \mathbf{x}, \theta) = \mathbf{0}$	$\text{s.t. } \frac{d\mathbf{s}}{dt} = \mathbf{f}(\mathbf{s}(t), t, \mathbf{x}, \theta)$
(a) Output defined explicitly as a composition of functions.	(b) Output defined implicitly via a nonlinear implicit equation.	(c) Output defined implicitly via a differential equation.

Figure 1: Supervised learning optimization problem definitions

2 Formulation of Training as an Optimization Task

The goal of training any supervised machine learning model is to find the optimal parameters θ that associate low cost values C to input-output pairs $\{\mathbf{x}, \mathbf{y}_{\text{true}}\}_{i=1}^N$ from the training dataset. This objective is usually formulated in the form of an optimization problem, where the constraints define how the output is obtained.

Figure 1 shows three classes of optimization problems in machine learning. In conventional feedforward networks (Fig. 1a), the output \mathbf{y} is derived explicitly through the composition of the functions \mathbf{U}_n through \mathbf{U}_0 , where \mathbf{U}_i are usually the layers of the model. Conversely, equilibrium models represent \mathbf{y} as a function of the network’s state variable \mathbf{s} . In the case of EP (Fig. 1b), this function is expressed through an nonlinear implicit equation, whereas in Neural ODEs (Fig. 1c), it is defined through a differential equation.

Regardless of the method employed to generate the output, the supervised learning process comprises three main stages: (1) computation of \mathbf{y} or \mathbf{s} , (2) computation of the loss gradient $\frac{dC}{d\mathbf{y}}$ or $\frac{dC}{d\mathbf{s}}$, and (3) computation of the parameter gradient $\frac{dC}{d\theta}$.

Computing \mathbf{y} , often referred to as the inference/forward step, is simpler in feedforward networks due to the explicit form of \mathbf{y} . For equilibrium models, this necessitates the use of a nonlinear solver, such as the Newton-Raphson method or Euler’s method for ODEs.

In feedforward networks, the parameter gradient $\frac{dC}{d\theta}$ is typically computed via the backpropagation algorithm, which uses the chain rule to propagate gradients from the output layer towards the input layer. However, when \mathbf{y} is defined implicitly, calculating $\frac{dC}{d\theta}$ requires a different approach. This is the subject of the next two section.

3 Neural ODEs

3.1 Motivation

Neural ODEs are a family of deep neural network models that can be interpreted as a continuous equivalent of Residual Networks (ResNets). In ResNets, the hidden state $\mathbf{s}[t + 1]$ at layer $t + 1$ is derived from the hidden state $\mathbf{s}[t]$ at layer t according to the equation $\mathbf{s}[t + 1] = \mathbf{s}[t] + \mathbf{f}(\mathbf{s}[t], \mathbf{x}, \theta)$, where $\mathbf{f}(\cdot)$, for many applications, is a simple feedforward network, and θ are its parameters. This transformation can be viewed as a discretization of a derivative $\mathbf{s}'(t)$ with timestep $\Delta t = 1$. Upon taking Δt to zero, this leads to an ordinary differential equation: $\frac{d\mathbf{s}(t)}{dt} = \mathbf{f}(\mathbf{s}(t), t, \mathbf{x}, \theta)$. The output layer $\mathbf{s}[T]$ is then defined as the solution of this ODE, starting from an initial state $\mathbf{s}[0]$.

3.2 Adjoint Method for ODEs

Consider the optimization problem presented in Fig. 1c, where the goal is to find the parameters θ that minimize the cost C at the output layer $\mathbf{s}(T)$. The adjoint method for ODEs provides a way to backpropagate through a model whose output is defined implicitly through a differential equation. The algorithm is provided in Alg 1. For its derivation, see [1].

The solution of equation (2) requires knowing the value of $\mathbf{s}(t)$ along its entire trajectory. Storing these evaluations in step 1 is impractical for large systems or long time spans due to memory

Algorithm 1 Adjoint method for constraints defined via ODEs

- 1: Solve the (forward) differential equation for $0 \leq t \leq T$ to get $s(t)$.

$$\begin{cases} \frac{ds}{dt} = \mathbf{f}(s(t), t, \mathbf{x}, \boldsymbol{\theta}) \\ \mathbf{s}(t=0) = \mathbf{s}_0 \end{cases} \quad (1)$$

- 2: Using the $s(T)$ from step 1, find $\boldsymbol{\lambda}(T)$. Then solve the (backward) differential equation from $t = T$ to $t = 0$.

$$\begin{cases} \frac{d\boldsymbol{\lambda}}{dt} = -[\boldsymbol{\lambda}(t)]^T \frac{\partial \mathbf{f}}{\partial \mathbf{s}} \\ \boldsymbol{\lambda}(T) = \left[\frac{\partial C(\mathbf{s}(T), \mathbf{s}_{\text{true}})}{\partial \mathbf{s}} \right]^T \end{cases} \quad (2)$$

- 3: Plug $\boldsymbol{\lambda}(t)$ to get $\frac{dC}{d\boldsymbol{\theta}}$.

$$\frac{dC}{d\boldsymbol{\theta}} = - \int_T^0 [\boldsymbol{\lambda}(t)]^T \frac{\partial \mathbf{f}}{\partial \boldsymbol{\theta}} dt \quad (3)$$

Algorithm 2 Adjoint method for constraints defined via implicit nonlinear equations

- 1: Given $\boldsymbol{\theta}$, solve the nonlinear equation given by the constraint to get s .

$$\mathbf{f}(s, \boldsymbol{\theta}) = \mathbf{0} \quad (4)$$

- 2: Using the s from step 1, find $\boldsymbol{\lambda}$.

$$\left(\frac{\partial \mathbf{f}}{\partial \mathbf{s}} \right)^T \boldsymbol{\lambda} = - \left(\frac{\partial C}{\partial \mathbf{s}} \right)^T \quad (5)$$

- 3: Plug $\boldsymbol{\lambda}$ in the equation below to get $\frac{dC}{d\boldsymbol{\theta}}$.

$$\frac{dC}{d\boldsymbol{\theta}} = \boldsymbol{\lambda}^T(\boldsymbol{\theta}) \frac{\partial \mathbf{f}}{\partial \boldsymbol{\theta}} \quad (6)$$

72 limitations. The authors of [1] addressed this issue by forming an augmented system of ODEs that
 73 recomputes $s(t)$ concurrently with $\boldsymbol{\lambda}(t)$ during step 2.

74 While ODE solvers can be implemented in analog hardware, it is not obvious how the second and
 75 third steps could be realized. However, if we make some assumptions about the nature of the ODE,
 76 the resulting system has a direct interpretation in analog hardware. This is discussed next.

77 3.3 Bridging Neural ODEs to Physical Systems

78 Neural networks are well-known as universal function approximators, traditionally expressed in
 79 the form $\mathbf{y} = \mathbf{f}(\mathbf{x}, \boldsymbol{\theta})$. EP extends this paradigm to implicitly-defined neural networks. In such
 80 networks, the mapping from input to output is defined not directly but via an implicit equation,
 81 offering increased representational capacity.

82 Neural ODEs take this notion a step further by introducing a time component into the system. A
 83 unique feature of Neural ODEs is the ability to define a family of mappings between inputs and
 84 outputs. This is achieved by choosing different initial states for the system. After a time T , different
 85 initial states will generally lead to different final states, further enhancing the expressiveness of the
 86 model. This idea is depicted in Fig. 2.

87 The question arises: Can these Neural ODE systems be connected to stable physical systems, specif-
 88 ically analog circuits? In nonlinear resistive analog networks, for example, the system converges to
 89 the same stable state irrespective of the initial conditions (s) as long as conductances ($\boldsymbol{\theta}$) and input
 90 voltages (\mathbf{x}) are fixed. Mathematically, this is represented by $\frac{ds}{dt} = \mathbf{0}$. When a Neural ODE system
 91 reaches steady-state, at $t = \infty$, the differential equation $\frac{ds}{dt} = \mathbf{f}(s(t), t, \mathbf{x}, \boldsymbol{\theta})$ effectively reduces
 92 to an implicit nonlinear equation $\mathbf{f}(s_\infty, \mathbf{x}, \boldsymbol{\theta}) = \mathbf{0}$. Furthermore, if the differential equation has
 93 a single equilibrium point at steady-state, then it can be represented by a nonlinear analog circuit.
 94 This idea is depicted in Fig. 3, which shows that for some specific ODEs, when $\boldsymbol{\theta}$ and \mathbf{x} are fixed,
 95 the system always converges to the same steady-state irrespective of initial state of s .

96 In the next section, we develop the algorithm for solving optimization problems in which the con-
 97 straint is a nonlinear implicit equation, which is the steady-state solution of ODEs with one equilib-
 98 rium state at $t = \infty$.

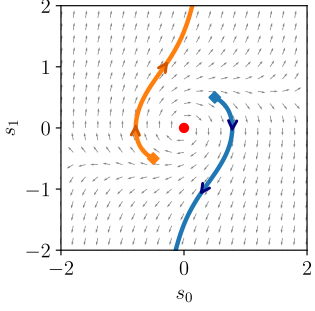


Figure 2: A example of a typical differential equation with no constraint on the trajectories.

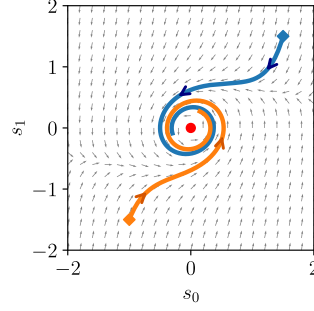


Figure 3: An example of a differential equation where all trajectories converge to the same equilibrium state.

99 4 Adjoint Method for Implicit Nonlinear Equations

100 Let $s, s_{\text{true}} \in \mathbb{R}^{n_s}$ and $\theta \in \mathbb{R}^{n_\theta}$. Suppose we have the function $C(s, s_{\text{true}}) : \mathbb{R}^{n_s} \times \mathbb{R}^{n_s} \rightarrow \mathbb{R}$, where
 101 $s(\theta)$ is the solution of the implicit function $f(s, x, \theta) = \mathbf{0}$ for a function $f : \mathbb{R}^{n_s} \times \mathbb{R}^{n_x} \times \mathbb{R}^{n_\theta} \rightarrow \mathbb{R}^{n_s}$.
 102 What is $\frac{dC}{d\theta}$?

103 C is a function of s , which is a function of θ . Therefore:

$$\frac{d\hat{C}}{d\theta} = \frac{\partial C}{\partial s} \frac{ds}{d\theta} \quad (7)$$

104 The differentiability of \hat{C} with respect to θ depends on the differentiability of s with respect to θ .
 105 By assuming that $s(\theta)$ exists and is the solution of the implicit function $f(s, x, \theta) = \mathbf{0}$, the implicit
 106 theorem guarantees the differentiability of $s(\theta)$ [3]. Therefore, at the points $s = s(\theta)$, the following
 107 relationship holds:

$$\frac{df}{d\theta} = \frac{\partial f}{\partial s} \frac{ds}{d\theta} + \frac{\partial f}{\partial \theta} \quad (8)$$

108 Since $f(s, \theta) = \mathbf{0}$ everywhere, $\frac{df(s(\theta), \theta)}{d\theta} = \mathbf{0}$. Assuming that $\frac{\partial f}{\partial s}$ is a non-singular matrix, the
 109 equation above can be expressed as:

$$\frac{ds}{d\theta} = - \left(\frac{\partial f}{\partial s} \right)^{-1} \frac{\partial f}{\partial \theta} \quad (9)$$

110 Substituting equation (9) into equation (7) yields:

$$\frac{d\hat{C}}{d\theta} = \left[- \frac{\partial C}{\partial s} \left(\frac{\partial f}{\partial s} \right)^{-1} \right] \frac{\partial f}{\partial \theta} \quad (10)$$

111 Let $\lambda(\theta) \in \mathbb{R}^{n_s \times 1}$ be the solution of term inside the square brackets.

$$\left(\frac{\partial f}{\partial s} \right)^T \lambda(\theta) = - \left(\frac{\partial C}{\partial s} \right)^T \quad (11)$$

112 Then the gradient can be written as:

$$\frac{d\hat{C}}{d\theta} = \lambda^T(\theta) \frac{\partial f}{\partial \theta} \quad (12)$$

113 The algorithm for solving optimization problems with implicitly-defined constraints is summarized
 114 in Alg. 2. It is much simpler than Alg. 1 and has a direct interpretation in analog circuits.

115 5 Interpretation of the Adjoint Variable in Analog Circuits

116 In the context of analog circuits, we can think of s as the voltages V in the nodes of the circuit,
 117 s_{true} as the desired node voltages V_{true} , θ as the conductance matrix G , $f(s, \theta) = \mathbf{0}$ as the set of
 118 equations according to Kirchoff's current law, $s(\theta)$ as the function that maps conductances to those
 119 voltages that satisfy Kirchoff's current law, and $\hat{C}(\theta, s_{\text{true}}) = C(s(\theta), s_{\text{true}})$ as the function that
 120 evaluates the difference between the measured and desired voltages of the circuit.

121 Substituting V for s , G for θ , and I for f , the rest of the terms in Alg. 2 can be interpreted as
 122 follows:

- 123 • The terms $\frac{\partial I}{\partial V}$ is the derivative of the current with respect to the voltage, and therefore has
 124 a unit of Siemens.
- 125 • The term $\frac{\partial I}{\partial G}$ is the derivative of the current with respect to the conductance, and therefore
 126 has a unit of Volts.
- 127 • In most cases, the cost function C is the mean-squared-error. C therefore has units of Watts
 128 and $\frac{\partial C}{\partial V}$ has a unit of Amperes.

129 With the link to analog circuits established, the adjoint algorithm can be interpreted as follows.
 130 The first step requires the solution of a nonlinear equation. This is exactly the steady-state (or DC)
 131 solution of an analog circuit. Let us denote these voltages as V^* . The current state of the circuit
 132 is compared against the desired state V_{true} . The term $\frac{\partial C}{\partial V^*}$ measures how the node voltages should
 133 be modified to reduce the discrepancy between the measured and the desired state. If we choose to
 134 inject a current of $\frac{\partial C}{\partial V^*}$ externally into the circuit, then since $\lambda = -\left(\frac{\partial I}{\partial V^*}\right)^{-T} \left(\frac{\partial C}{\partial V^*}\right)^T$ has units of
 135 Volts, λ is exactly the variation in the node voltages due to the injected current.

136 In a nonlinear analog circuit, the calculation of λ as presented above is impossible. Because of the
 137 presence of nonlinear elements, the injection of current into the circuit can completely change the
 138 state of a circuit. For example, a diode on the verge of conduction can switch on if the voltage across
 139 it increases as a result of the injected current. If this diode is connected across a large computational
 140 block, the block will be shorted to ground when the diode turns on, significantly altering the circuit
 141 between the two states. Since the calculation of λ requires the evaluation of the derivative $\frac{\partial I}{\partial V}$ at
 142 exactly the steady-state V^* of the circuit, we have to ensure that the injected current does not disturb
 143 the state of the circuit. One approach is to scale the injected current $\left(\frac{\partial C}{\partial V}\right)$ by some factor β to ensure
 144 that it is sufficiently small. Voltage variations in the nodes can then be measured and divided by β
 145 to obtain the unscaled version.

146 In summary, while the adjoint method can indeed be implemented on an analog circuit, the gradient
 147 calculated is an approximation since the act of injecting a current to measure the derivatives disturbs
 148 the circuit.

149 6 Connection to EP and Verification

150 To validate our analysis in the previous section, we designed a simple circuit, as an illustrative
 151 example, to learn the XOR dataset. The circuit has an input layer with dimensions 4×2 , followed
 152 by a pair of up-down diodes to introduce nonlinearity, and an output layer with dimensions 4×2 .

153 The XOR dataset has two inputs and one output. The inputs are represented by X_1 and X_2 and the
 154 output by the difference of Y_1 and Y_2 . This is to account for the constraint that resistances cannot
 155 assume negative values. X_3 and X_4 represent the bias voltages and are of opposite polarity. We
 156 considered the node voltages V_1, V_2, Y_1, Y_2 as our state variables.

157 As an example, consider the conductance G^\dagger (circled in the schematic). It is connected between
 158 nodes V_1 and Y_1 . As a result, the voltage drop across it is $V^\dagger = V_1 - Y_1$. According to equation (6),
 159 $\left[\frac{dC}{dG^\dagger}\right]_{\text{Adjoint}} = \lambda^T \frac{dI}{dG^\dagger} \cdot \frac{dI}{dG^\dagger}$, in this example, is of size 1×4 . Since G^\dagger is connected to only V_1

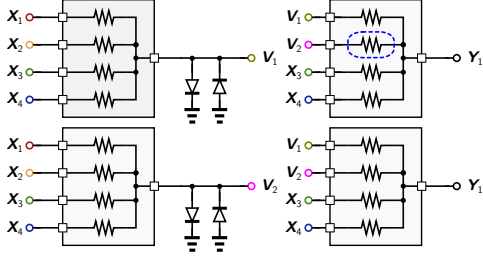


Figure 4: Schematic diagram for learning the XOR dataset.

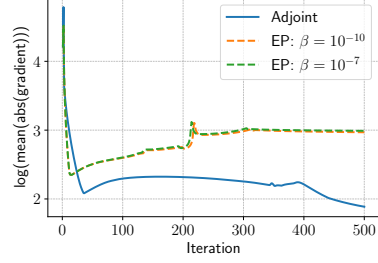


Figure 5: Plot of the gradients during training.

160 and Y_1 , $\frac{dI}{dG^\dagger}$ will only have entries in the positions for nodes V_1 and Y_1 , corresponding to λ_1 and
 161 λ_3 , and these will be of opposite polarity. Thus, the evaluation of $[\frac{dC}{dG^\dagger}]_{\text{Adjoint}}$ yields $(\lambda_1 - \lambda_3)V^\dagger$.
 162 This shows that the update to conductance G^\dagger is the product of the voltage drop at steady-state and
 163 the variation due to the injected current across the two nodes it is connected to.

164 Let's now compare this result to the one from EP. In EP, the update equation for G^\dagger , as applied
 165 to this circuit is $[\frac{dC}{dG^\dagger}]_{\text{EP}} \approx \frac{1}{2\beta} [(V_\beta^\dagger)^2 - (V_0^\dagger)^2]$, where V_0^\dagger is the voltage drop across G^\dagger before
 166 the current was injected, and V_β^\dagger is the voltage drop after the current was injected. If we express
 167 V_β^\dagger as $V_0^\dagger + \Delta V$, then $[\frac{dC}{dG^\dagger}]_{\text{EP}} \approx \left(\frac{\Delta V}{\beta}\right) (V_0^\dagger) + \frac{1}{2\beta} (\Delta V)^2$. If the current injected is very small,
 168 $(\Delta V)^2 \rightarrow 0$, then $[\frac{dC}{dG^\dagger}]_{\text{Adjoint}} \approx [\frac{dC}{dG^\dagger}]_{\text{EP}}$, where $(\lambda_1 - \lambda_3) \approx \frac{\Delta V}{\beta}$.

169 The circuit was optimized using using both the adjoint and the EP methods [4]. A plot of the log
 170 of the mean absolute values of all the conductances across iterations is shown in Fig. 5. Initially,
 171 the gradients from both methods exhibit close proximity. Over time, however, the approximations
 172 inherent in the EP method introduce biases that cause its gradients to diverge from the true gradient.
 173 In contrast, the gradients obtained via the adjoint method continue to get smaller, indicating that the
 174 circuit is getting closer to the desired state.

175 7 Conclusion

176 A central theme of this work is the impact of constraints in optimization problems on neural network
 177 design and training rules. We explored two types of constraints: nonlinear implicit equations and
 178 differential equations. The adjoint method emerged as a unifying framework, linking Neural ODEs
 179 and EP and offering efficient gradient computation. This has practical implications, especially in
 180 analog computation, where we have interpreted the adjoint variable as the variation in node voltages
 181 in an analog circuit.

182 Looking ahead, we aim to focus on the efficient implementation of Neural ODEs in analog hardware.
 183 Additional future work could explore other types of constraints that are more suitable for
 184 hardware design in systems other than analog circuits, opening up new avenues for efficient and
 185 adaptable neural network training.

186 References

- 187 [1] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary
 188 differential equations, 2019.
- 189 [2] Benjamin Scellier and Yoshua Bengio. Equilibrium propagation: Bridging the gap between
 190 energy-based models and backpropagation, 2017.
- 191 [3] Matthias Heinkenschloss. Implicit constraints, 2002. Available at <https://www.cmor-faculty.rice.edu/~cox/neuro/ImplConstraints.pdf>, accessed on 2023-09-29.
- 193 [4] Ebana framework. Available at <https://github.com/mawatfa/ebana>.