A Truly Sparse and General Implementation of Gradient-based Synaptic Plasticity

Jamie Lohoff^{1,2}, Anil Kaya^{1,2}, Florian Assmuth^{1,2}, and Emre Neftci^{1,2}

¹ Jülich Research Center, Peter-Grünberg-Institute 15, 52428 Jülich, Germany













Gradient-based Synaptic Plasticity

Long-term potentiation (LTP) is regulated by coincidence of action potentials (AP) and postsynaptic potentials (EPSP)



Markram et al., Science, 1997

Synaptic plasticity rules can be made compatible with a local approximation of gradient descent

Gradient-based synaptic weight upd

Three Factor Rules modulates learning with higher level information (reward, error, etc.)

late:
$$\Delta W \propto \frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial S^t} \Theta'(U^t) \frac{\partial U^t}{\partial W}$$

Gradient-based Synaptic Plasticity

Synaptic plasticity rules can be made compatible with a local approximation of gradient descent



E-prop is capable of online learning and could efficiently implemented, but in practice:

- Use BPTT + stop-gradient -> not online, memory-intensive
- Hand-implement eligibility traces -> does not scale well (in manpower etc.)

Gradient-based synaptic weight update: $\Delta W \propto \frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial S^t} \Theta'(U^t) \frac{\partial U^t}{\partial W}$

Recurrence Relation!



We have the best of both worlds...

What is Automatic Differentiation?

- How to compute derivatives algorithmically up to machine precision?
- Applications in ML, Robotics, Computational Fluid Dynamics, Finance, ...
- AutoDiff != gradient backpropagation

h(x) = f(g(x))

Uses chain rule just multiplication of Jacobians

single variable: $h'(x) = f'(g(x)) \cdot g'(x)$

multi variable: $Dh(x) = Df(g(x)) \cdot Dg(x)$

Elemental derivatives are hard-coded





Andreas Griewank • Andrea Walther



Evaluating Derivatives

Principles and Techniques of Algorithmic Differentiation



defelemental(lax.sin_p, lambda x: jnp.cos(x)) defelemental(lax.asin_p, lambda x: 1./jnp.sqrt(1.0 - x**2)) defelemental(lax.cos_p, lambda x: -jnp.sin(x)) defelemental(lax.acos_p, lambda x: -1./jnp.sqrt(1.0 - x**2)) defelemental2(lax.tan_p, lambda out, primal: 1.+out**2) defelemental(lax.atan_p, lambda x: 1./(1. + x**2))



Graph View on Automatic Differentiation

- Systematic construction of AutoDiff algorithms
- Algorithms tailored to functions _____ minimize computational resources

$f(x_1, x_2) = \begin{pmatrix} \log \sin x_1 x_2 \\ x_1 x_2 + \sin x_1 x_2 \end{pmatrix}$



Tape $v_{-1} = x_1$ $v_0 = x_2$ $v_1 = v_{-1}v_0$ $v_2 = \sin v_1$ $v_3 = \log v_2$ $v_4 = v_1 + v_2$

Constructing the computational graph

Tape $v_{-1} = x_1$ $v_0 = x_2$ $v_1 = v_{-1}v_0$ $v_2 = \sin v_1$ $v_3 = \log v_2$ $v_4 = v_1 + v_2$



Adding the Elemental Derivatives

Tape $v_{-1} = x_1$ $v_0 = x_2$ $v_1 = v_{-1}v_0$ $v_2 = \sin v_1$ $v_3 = \log v_2$ $v_4 = v_1 + v_2$

$$c_{10} = v_{-1}$$

 $c_{21} = \cos v_1$
 $c_{32} = 1/v_2$
 $c_{41}, c_{4,2} = 1$



Constructing the AutoDiff Algorithm

Tape $v_{-1} = x_1$ $v_0 = x_2$ $v_1 = v_{-1}v_0$ $v_2 = \sin v_1$ $v_3 = \log v_2$ $v_4 = v_1 + v_2$ $c_{1-1} = v_0$ $c_{10} = v_{-1}$ $c_{21} = \cos v_1$ $c_{32} = 1/v_2$ $c_{41}, c_{4,2} = 1$



Eliminating the 2nd Vertex
Tape

$$v_{-1} = x_1$$

 $v_0 = x_2$
 $v_1 = v_{-1}v_0$
 $v_2 = \sin v_1$
 $v_3 = \log v_2$
 $v_4 = v_1 + v_2$
 $c_{1-1} = v_0$
 $c_{10} = v_{-1}$
 $c_{21} = \cos v_1$
 $c_{32} = 1/v_2$
 $c_{41}, c_{4,2} = 1$
 $c_{31} = c_{32}c_{21}$
 $c_{41} = c_{41} + c_{42}c_{21}$



How to find good elimination orders?

Lohoff & Neftci. Optimizing Automatic Differentiation with Deep Reinforcement Learning. 38th Conference on Advances in Neural Information Processing (NeurIPS) 2024. Spotlight Paper.





And does it actually work?



MLP evaluation times for different modes and batch sizes RobotArm evaluation times for different modes and batch sizes



Task	Forward	Reverse	Markowitz	Alp
RoeFlux_1d	620	364	407	
RobotArm_6DOF	397	301	288	
HumanHeartDipole	240	172	194	
PropaneCombustion	151	90	111	
Random function g	632	566	451	
BlackScholes	545	572	350	
RoeFlux_3d	1556	979	938	
Random function f	17728	9333	12083	6
2-layer MLP [†]	10930	392	4796	398
Transformer [†]	135010	4688	51869	4831

Task	Forward	Reverse	Markowitz
RoeFlux_1d	$3.03\substack{+0.17 \\ -0.27}$	$3.08\substack{+0.17 \\ -0.23}$	$2.87\substack{+0.22\\-0.66}$
RobotArm_6DOF	$8.85\substack{+0.21\\-0.17}$	$8.48\substack{+0.32\\-0.20}$	$8.55\substack{+0.38\\-0.36}$
HumanHeartDipole	$16.97\substack{+1.23 \\ -2.39}$	$16.87^{+1.45}_{-3.90}$	$16.42^{+1.29}_{-2.73}$
PropaneCombustion	$36.91\substack{+2.87 \\ -1.50}$	$36.47\substack{+1.45\\-0.51}$	$36.99^{+1.55}_{-1.11}$
Random function g	$82.41^{+2.85}_{-2.97}$	$81.64\substack{+2.82\-3.56}$	$82.97^{+1.38}_{-1.17}$
BlackScholes	$5.04\substack{+0.23\\-0.33}$	$5.02\substack{+0.30\\-0.39}$	$5.03\substack{+0.23\\-0.29}$
RoeFlux_3d	$72.26\substack{+3.22\\-6.02}$	$83.98\substack{+5.69\\-6.68}$	$91.27^{+11.59}_{-16.49}$
Random function f	$12.42\substack{+0.66\\-0.25}$	$11.54\substack{+0.15\\-0.27}$	$20.20\substack{+0.26\\-0.31}$
2-layer MLP [†]	$760.63\substack{+80.01\-53.30}$	$29.67\substack{+1.33 \\ -4.05}$	$317.65\substack{+53.30\\-19.34}$
2-layer MLP [†] (GPU)	$11.04\substack{+0.31\\-0.13}$	$0.30\substack{+0.02\\-0.03}$	$1.01\substack{+0.02\\-0.05}$
Transformer [†]	$990.09\substack{+35.42\-31.11}$	$39.07\substack{+4.59\-7.67}$	$498.94\substack{+20.34\\-19.62}$
Transformer [†] (GPU)	$21.38\substack{+0.22\\-0.06}$	$0.29\substack{+0.02\\-0.03}$	$16.17^{+15.04}_{-0.09}$





Tensorized Vertex Elimination - Sparse AutoDiff

 $f(\mathbf{W}, \mathbf{x}) = \tanh(\mathbf{W} \cdot \mathbf{x})$



Simple LIF Layer:

 $\mathbf{u}_{t+1} = \alpha \mathbf{u}_t + \mathbf{W} \cdot \mathbf{z}_t$ $\mathbf{z}_{t+1} = \Theta(\mathbf{u}_t - \vartheta)$



Temporal Credit Assignment

Recurrent Neural Network: $h_t = f(h_{t-1}, x_t, \theta)$

Gradient:
$$\frac{\mathrm{d}L}{\mathrm{d}\theta} = \sum_{t} \frac{\mathrm{d}L}{\mathrm{d}h_t} \frac{\partial h_t}{\partial \theta}$$

 $\frac{\mathrm{d}L}{\mathrm{d}h_t} = \frac{\partial L}{\partial h_{t+2}} \frac{\partial h_{t+2}}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_t} + \frac{\partial L}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_t} + \frac{\partial L}{\partial h_t}$ $\frac{\mathrm{d}L}{\mathrm{d}h_t} = \left(\frac{\partial L}{\partial h_{t+2}}\frac{\partial h_{t+2}}{\partial h_{t+1}} + \frac{\partial L}{\partial h_{t+1}}\right)\frac{\partial h_{t+1}}{\partial h_t} + \frac{\partial L}{\partial h_t}$ $c_t = (c_{t+2}H_{t+1} + d_{t+1})H_t + d_t$

 $BPTT: c_t = c_{t+1}H_t + d_t$

Compute Scaling: $\mathcal{O}(n^2T)$

Memory Scaling: $\mathcal{O}(nT)$



 h_{t-1}





More Temporal Credit Assignment

Gradient:
$$\frac{\mathrm{d}L}{\mathrm{d}\theta} = \sum_{t} \frac{\partial L}{\partial h_t} \frac{\mathrm{d}h_t}{\mathrm{d}\theta}$$

 $\frac{\mathrm{d}h_{t+2}}{\mathrm{d}\theta} = \frac{\partial h_{t+2}}{\partial h_{t+1}} \frac{\partial h_{t}}{\partial h_{t}} \frac{\partial h_{t}}{\partial \theta} + \frac{\partial h_{t+2}}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial \theta} + \frac{\partial h_{t+2}}{\partial \theta}$ $\frac{\mathrm{d}h_{t+2}}{\mathrm{d}\theta} = \frac{\partial h_{t+2}}{\partial h_{t+1}} \left(\frac{\partial h_{t+1}}{\partial h_{t}} \frac{\partial h_{t}}{\partial \theta} + \frac{\partial h_{t+1}}{\partial \theta}\right) + \frac{\partial h_{t+2}}{\partial \theta}$

 $G_{t+2} = H_{t+2}(H_{t+1}G_t) + F_{t+1}) + F_{t+2}$

RTRL: $G_t = H_t G_{t-1} + F_t$

Compute Scaling: $\mathcal{O}(n^4T)$

Memory Scaling: $\mathcal{O}(n^3)$





Gradient-based Synaptic Plasticity (again)



 H_t becomes diagonal, so G_t stays diagonal:



		RTRL	BPTT
Scaling Laws:	Compute	$\mathcal{O}(n^4T)$	$\int \mathcal{O}(n^2 T)$
	Memory	$\mathcal{O}(n^3)$	O(nT)

Simple LIF Neuron: $\begin{aligned} \mathbf{u}_{t+1} &= \alpha \mathbf{u}_t + \mathbf{W} \cdot \mathbf{z}_t \\ \mathbf{z}_{t+1} &= \Theta(\mathbf{u}_t - \vartheta) \end{aligned}$ $\mapsto h_t = \mathbf{u}_t$

New Compute Scaling: $\mathcal{O}(n^2T)$ New Memory Scaling: $\mathcal{O}(n^2)$

Zenke and Neftci. Brain-inspired learning on neuromorphic substrates. Proceedings of the IEEE 109, Issue 5. 2021.

Bellec et al. A solution to the learning dilemma for recurrent networks of spiking neurons. Nature Communications 11. 2020.

James M Murray Local online learning in recurrent networks with random feedback. eLife 8:e43299. 2019.







Our Method in Action:





Outlook into the Future



• Generalize to multiple layers to enable large-scale training & meaningful applications

Kaiser et al. Synaptic Plasticity Dynamics for Deep Continuous Local Learning (DECOLLE). Frontiers of Neuroscience 14. 2020.

Bohnstingl et al. Online spatiotemporal learning in deep neural networks. IEEE TNNLS 34, Issue 11. 2023.

Chavlis & Poirazi. Dendrites endow artificial neural networks with accurate, robust and parameter-efficient learning. Arxiv:2404.03708. 2024.

- Test more intricate neuron types, train SNNs on very long time-horizons •
- Use Vertex Elimination formalism to find **novel approximation methods** (e.g. with RL)
- Include **HW constraints** in the search for good gradient-based Synaptic Plasticity

Really Fast and Memory-Graphax = efficient SNN Training

Application to **Dendritic Neural Networks** to make them scalable (with Yiota Poirazis Group)