

A Cortical Motif for Learning with Vectorized Errors

Qianli Liao^{* 1,5}, Liu Ziyin^{* 3,4}, Yulu Gan^{* 1,2}, Brian Cheung^{1,2}, Mark Harnett^{5,6}, and Tomaso Poggio^{1,2,5,6}

¹*Center for Brains, Minds, and Machines, MIT*

²*CSAIL, MIT*

³*Research Laboratory of Electronics, MIT*

⁴*Physics & Informatics Laboratories, NTT Research*

⁵*McGovern Institute, MIT*

⁶*Department of Brain and Cognitive Sciences, MIT*

February 25, 2026

Abstract

Recent experiments in mouse cortex suggest that instructive/error information is **vectorized**—distributed across neuron-specific signals in distinct neuronal populations—while converging evidence indicates that **both ascending and descending pathways** exhibit experience-dependent synaptic plasticity. These observations motivate learning architectures in which error is carried by a dedicated feedback stream and synapses across streams co-adapt. We introduce **Self-Assembling Motif (SAM)**, a minimal two-stream circuit motif of four classes of plastic connections trained by local heterosynaptic rules, without requiring matched forward–feedback representations. Starting from random connectivity, SAM self-organizes to implement **SGD-like learning** without weight transport or explicit derivatives, extending a line of biologically motivated alternatives to backpropagation. We prove that, under mild stationarity conditions, SAM approximates stochastic gradient descent with a learned positive semidefinite preconditioner, and we show competitive performance across vision benchmarks. SAM makes falsifiable predictions about the relationship between cell-type-specific error signals and bidirectional synaptic plasticity in cortex.

1 Introduction

How cortical circuits assign credit across many synapses remains a central open problem in neuroscience. Gradient-based learning in artificial neural networks provides an existence

30 proof that neuron-specific error signals can, in principle, support efficient multi-layer learn-
31 ing, but standard backpropagation relies on globally coordinated mechanisms across layers
32 or areas whose cortical implementation is unclear [17, 25].

33 Two recent experimental trends sharpen the biological constraints on any cortical credit-
34 assignment theory. First, instructive or error-related information in mouse cortex appears
35 to be *vectorized*: learning is guided by structured, neuron-resolved signals carried by specific
36 circuit elements rather than by a single broadcast scalar [1, 7, 8]. Such findings motivate
37 architectures in which error information is routed through dedicated pathways that can
38 deliver heterogeneous, unit-specific teaching signals.

39 Second, converging evidence indicates that synapses in both ascending and descending
40 pathways are plastic and experience-dependent, including synapses implicated in top-down
41 and cross-stream interactions [1, 4, 5, 7, 9, 19]. In particular, there is substantial evidence
42 that spines in layer 1 (where many “feedback” inputs make their synapses) on both L5
43 and L2/3 pyramidal neurons are reasonably plastic across the adult mouse cortex, and
44 that they show enhanced plasticity during learning paradigms [1, 7, 10, 13, 18, 24, 27, 28].
45 This suggests that effective learning circuits should not treat feedback connections as fixed
46 scaffolding and also not as matched one-to-one with the forward connections; instead, forward
47 representations and feedback/teaching routes are likely to co-adapt during learning.

48 Despite these constraints, many backprop-inspired circuit proposals assume correspon-
49 dence between forward and feedback representations—often requiring equal dimensionality,
50 or an implicit one-to-one pairing—to deliver neuron-specific error components to appropri-
51 ate targets [3, 11, 15, 16, 21, 22]. Whether such precise matching can be established and
52 maintained across development and learning seems still uncertain at best, particularly in
53 circuits where feedback pathways differ in size, cell types, and laminar organization [5, 19].

54 Here we introduce **Self-Assembling Motif (SAM)**, a minimal two-stream motif with
55 four classes of plastic connections trained by local heterosynaptic rules. SAM represents
56 error as a vector propagated through a feedback stream while removing the requirement for
57 matched forward-feedback representations: the forward and feedback pathways may have
58 different widths, and inter-stream mappings are learned through plasticity from initially ran-
59 dom connectivity. Starting from random connectivity, SAM self-organizes to support SGD-
60 like learning using only local heterosynaptic plasticity, without weight transport or explicit
61 derivatives, in the spirit of prior biologically motivated proposals [2, 3, 11, 15, 16, 21, 22]. We
62 prove that, under mild stationarity conditions, SAM’s updates approximate stochastic gra-
63 dient descent via a learned positive semidefinite preconditioner. Empirically, SAM performs
64 competitively on standard vision benchmarks [6, 12, 14, 20, 26, 29, 30] and yields falsifi-
65 able predictions linking cell-type-specific teaching signals to bidirectional synaptic plasticity
66 across cortical streams [5, 19].

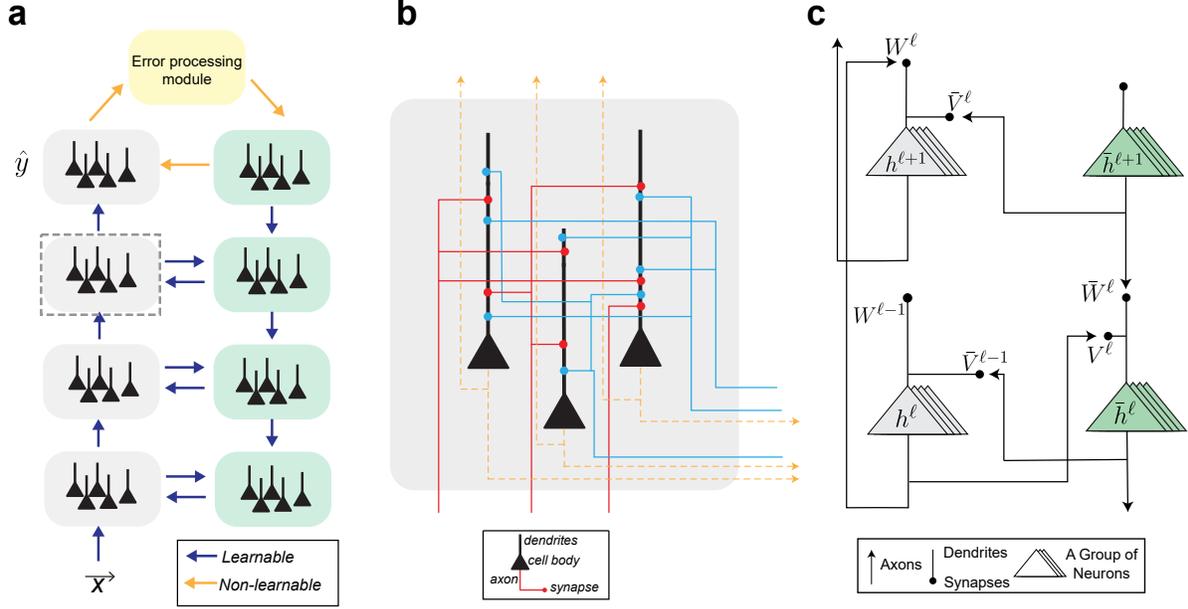


Figure 1: **A synaptic motif for interactions between ascending and descending information streams.** (a) Schematic of the upstream-downstream architecture. The upstream consists of a standard fully-connected neural network with multiple layers, roughly corresponding to a multi-region cortical processing pathway (e.g. V1-V2-V4-IT) [19]. The vector output of the upstream network goes to an error processing module (potentially corresponding to frontal cortices). This module computes a local error signal. This *vector* error signal is sent to the feedback (downstream) pathway, which processes information layer by layer downwards. The orange arrows represent non-learnable (identity) connections. Blue arrows represent learnable connections, each parameterized by a fully-connected weight matrix. (b) A detailed view of the smallest unit of the connection motif in panel a, indicated by grey dashed box. (c) A mathematical description of the unit in panel b. Each arrow represents full connections between two groups of neurons, parameterized by a weight matrix. Every h is a vector of activations. The inter-stream matrix V allows upstream and downstream layers to have different widths.

2 Results

2.1 A minimal two-stream motif with four plastic pathways

We consider an upstream (feedforward) stream and a downstream (feedback) stream that interact through four classes of synaptic pathways: upstream-to-upstream ($u \rightarrow u$), downstream-to-downstream ($d \rightarrow d$), upstream-to-downstream ($u \rightarrow d$), and downstream-to-upstream ($d \rightarrow u$) (Fig. 1), consistent with known feedforward/feedback organization in sensory cortex [5, 19]. Learning is driven by a local heterosynaptic rule of the form

$$\Delta W \propto S_{\text{post}} S_{\text{pre}}^{\top}, \quad (1)$$

where the postsynaptic factor is supplied by feedback/instructive inputs rather than by self-activation alone, consistent with the broader biological motivation for multi-factor plasticity

76 and feedback-gated learning [17, 25]. The full update rules for the four pathways are given in
77 Appendix A (and summarized in Methods). The motif of the figure is the simplest possible.
78 In reality, interneurons are also likely to be involved. Mechanisms such as burst-dependent
79 synaptic plasticity previously proposed ([23]) may also be involved increasing the effectiveness
80 of SAM.

81 **2.2 Vectorized error is represented by a self-learning feedback** 82 **stream**

83 In SAM, the instructive signal is a *vector* propagated through the downstream stream,
84 initialized at the output as $\bar{h}^L = \epsilon(\hat{y})$ and transmitted to lower levels to gate synaptic up-
85 dates. This architecture supports neuron-resolved vector signals without requiring matched
86 forward–feedback representations, aligning with the idea that biologically plausible learning
87 may rely on structured feedback pathways rather than exact weight transport [16, 17, 25].

88 To test whether vectorization is functionally necessary, we compared SAM to an ablated
89 variant in which the instructive signal is collapsed to a scalar modulatory factor shared
90 across neurons within a layer (Methods). Collapsing the vector signal substantially degrades
91 learning, consistent with the idea that neuron-specific error components provide essential
92 credit-assignment structure.

93 **2.3 Bidirectional plasticity enables self-assembly and robust learn-** 94 **ing**

95 Because SAM does not assume pre-established correspondences between streams, effective
96 learning requires that cross-stream mappings and feedback routes co-adapt with forward
97 representations. This is also compatible with increasing evidence that synapses implicated
98 in top-down or feedback circuitry can be experience-dependent [1, 4, 7, 9]. We therefore
99 tested the necessity of plasticity in each pathway by selectively freezing subsets of synapses
100 (Fig. 2).

101 We observe three robust trends. First, updating only the feedforward pathway W yields
102 limited improvement and tracks a linear baseline, indicating that cross-stream plasticity is
103 needed to learn nonlinear relations. Second, among inter-stream connections, making the
104 downstream-to-upstream pathway \bar{V} plastic provides the largest single gain, suggesting that
105 learning benefits strongly from adaptive routing of vectorized teaching signals into upstream
106 populations. Third, allowing all four classes of connections (W, \bar{W}, V, \bar{V}) to be plastic re-
107 covers the strongest performance and most reliably approaches the behavior of SGD-trained
108 networks [3, 11, 16] (Table 1; Fig. 2). Together, these results mirror emerging biological evi-
109 dence that synapses mediating ascending, descending, and cross-stream interactions exhibit
110 experience-dependent plasticity [4, 5, 9, 19].

111 3 Theory

112 3.1 SAM approximates preconditioned stochastic gradient descent

113 SAM updates the feedforward weights using a heterosynaptic, three-factor rule in which
 114 the presynaptic factor is the upstream activity and the postsynaptic factor is supplied by
 115 a feedback/instructive signal, following the general motivation for alternatives to backprop
 116 that avoid exact weight transport [3, 11, 15–17, 25]. Writing the upstream activations at
 117 layer ℓ as h^ℓ and the feedback-stream activities as \bar{h}^ℓ , the feedforward synapses obey an
 118 update of the generic form

$$\Delta W^\ell = \eta \bar{h}^\ell (h^{\ell-1})^\top - \eta \gamma W^\ell, \quad (2)$$

119 where η is a learning rate and γ is a weight-decay coefficient.

120 **Key idea.** If the feedback stream delivers an *approximately linear* transformation of the
 121 true backprop error at each layer, then the local rule in Eq. (2) becomes a *preconditioned* ver-
 122 sion of the SGD gradient [17, 25]. In SAM, this linear transformation is itself learned through
 123 plastic inter-stream connectivity, yielding an adaptive, data-dependent preconditioner.

124 **Theorem 1** (Informal: preconditioned-SGD equivalence). *Assume (i) the inter-stream map-*
 125 *plings vary slowly relative to the feedforward weights (a stationarity/separation-of-timescales*
 126 *condition), (ii) the learned feedback-to-upstream map has full row rank so that error compo-*
 127 *nents span the relevant subspace, and (iii) the feedback-stream activity \bar{h}^ℓ is a stable linear*
 128 *transform of the layerwise backprop error signal δ^ℓ up to higher-order terms. Then the SAM*
 129 *update on the feedforward weights can be written as*

$$\Delta_{\text{SAM}} W^\ell = H^\ell \Delta_{\text{SGD}} W^\ell - \eta \gamma W^\ell + o(1). \quad (3)$$

130 where H^ℓ is a learned positive semidefinite matrix (a preconditioner) determined by inter-
 131 stream connectivity.

132 **Form of the preconditioner.** In the simplest (and empirically relevant) regime, the
 133 preconditioner takes the form

$$H^\ell \approx \bar{V}^\ell (\bar{V}^\ell)^\top, \quad (4)$$

134 where \bar{V}^ℓ denotes the learned mapping that routes feedback-stream signals into the upstream
 135 population at layer ℓ .¹ Equation (4) makes two points explicit: (i) $H^\ell \succeq 0$ by construction,
 136 and (ii) the “effective learning rate structure” for W^ℓ is set by cross-stream plasticity.

137 **Sketch of derivation.** Under the above assumptions, the feedback factor can be expressed
 138 as $\bar{h}^\ell \approx \bar{V}^\ell \delta^\ell$ (up to higher-order terms). Substituting into Eq. (2) gives

$$\Delta W^\ell \approx \eta (\bar{V}^\ell \delta^\ell) (h^{\ell-1})^\top - \eta \gamma W^\ell.$$

¹The exact definition of \bar{V}^ℓ and the conditions under which Eq. (4) holds are given in Appendix A.

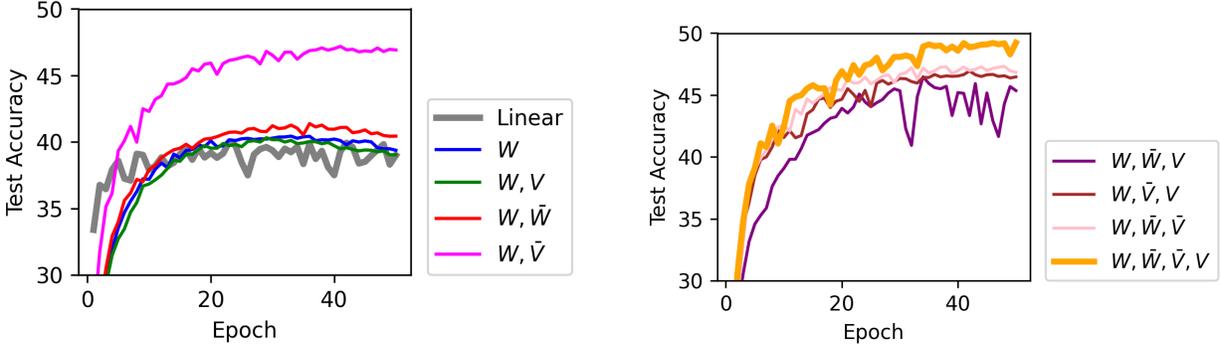


Figure 2: **Ablation studies.** Ablation study on roles of \bar{V} , \bar{W} and V on CIFAR-10 [12]. Here, we make a subset of all interconnections that are not trainable. The legend labels indicate trainable components. "Linear" refers to a linear network trained with SGD. We see that (1) making \bar{V} plastic is more important than making other connections plastic, and (2) making everything plastic significantly improves the performance further.

139 Since the SGD gradient has the standard form $\Delta_{\text{SGD}}W^\ell \propto \delta^\ell(h^{\ell-1})^\top$, we obtain Eq. (3)
 140 with a learned left-multiplication by H^ℓ induced by inter-stream routing. The $o(1)$ term
 141 collects the approximation error from imperfect stationarity and residual nonlinear coupling
 142 between streams. A limitation of the present analysis is that the convergence result relies on
 143 stationarity assumptions for inter-stream mappings; characterizing the full coupled dynamics
 144 in biophysical terms remains a goal for future work.

145 **Mechanistic interpretation.** The learned matrix H^ℓ acts as a data- and time-dependent
 146 preconditioner: SAM does not merely approximate SGD, but does so with an *adaptive*
 147 *geometry* set by experience-dependent inter-stream synapses. Biologically, this links im-
 148 provements in credit assignment to measurable learning-dependent changes in cross-stream
 149 connectivity, and predicts that manipulations that disrupt \bar{V}^ℓ plasticity should specifically
 150 reduce the "conditioning" benefits of learning even when feedforward plasticity remains in-
 151 tact [1, 4, 7, 9]. Consistent with this view, $\bar{V}^\ell(\bar{V}^\ell)^\top$ develops structured correlations and
 152 effectively low-rank spectra during learning (Fig. S1a).

153 A full statement, proof, and the precise assumptions are provided in Appendices A and D.

154 4 Methods

155 **Model architecture.** We use a two-stream network with an upstream (feedforward) path-
 156 way and a downstream (feedback) pathway connected by four classes of synapses ($u \rightarrow u$,
 157 $d \rightarrow d$, $u \rightarrow d$, $d \rightarrow u$), consistent with canonical feedforward/feedback organization in cortex
 158 [5, 19]. The downstream pathway may be wider than the upstream pathway, consistent with
 159 anatomical asymmetries observed in some systems [5].

Table 1: **Performance of the proposed method compared with SGD and biologically plausible algorithm baselines.** Bold denotes the best performing algorithm, and italics denote the best runner-up algorithm. SGD and the SAM are comparable, while significantly outperforming existing biologically plausible learning algorithms. Baselines include Feedback Alignment [16], Weight Mirroring [3], and the Kolen–Pollack algorithm [11]. For each experiment, we conduct five runs.

| | CIFAR10 | MNIST | Fashion MNIST | Chest MNIST | Path MNIST | SVHN | STL10 |
|--------------|-------------------------------|-------------------------------|-------------------------------|-------------------------------|-------------------------------|-------------------------------|-------------------------------|
| SGD | <i>50.36</i> ± <i>0.16</i> | <i>98.44</i> ± <i>0.09</i> | <i>96.80</i> ± <i>0.05</i> | 73.26 ± 0.14 | 70.11 ± 0.10 | <i>78.20</i> ± <i>0.13</i> | <i>41.00</i> ± <i>0.07</i> |
| FA | 47.85 ± 0.20 | 97.80 ± 0.19 | 95.75 ± 0.20 | 70.44 ± 0.27 | 69.90 ± 0.10 | 78.03 ± 0.15 | 39.29 ± 0.23 |
| WM | 49.35 ± 0.13 | 98.40 ± 0.23 | 96.00 ± 0.25 | 70.14 ± 0.08 | 70.00 ± 0.15 | 78.08 ± 0.15 | 39.48 ± 0.19 |
| KP Algorithm | 48.16 ± 0.09 | 98.45 ± 0.09 | 96.13 ± 0.07 | 71.03 ± 0.10 | 69.51 ± 0.05 | 78.20 ± 0.05 | 40.32 ± 0.11 |
| SAM (ours) | 51.91 ± 0.21 | 98.85 ± 0.10 | 97.23 ± 0.08 | 73.21 ± 0.11 | 70.23 ± 0.15 | 78.50 ± 0.07 | 41.50 ± 0.10 |

160 **Learning rule and ablations.** All synapses are updated by local heterosynaptic rules
 161 (Appendix A). For the *scalar-error ablation*, the layerwise feedback vector is replaced by a
 162 shared scalar modulatory factor (e.g., the mean or norm of the feedback vector within the
 163 layer), holding all other components fixed.

164 For *plasticity ablations*, we selectively freeze subsets of the four connection classes while
 165 keeping W plastic.

166 **Datasets and training protocol.** We evaluate on standard vision benchmarks and re-
 167 port accuracy averaged across multiple runs (Table 1); datasets include CIFAR-10 [12],
 168 MNIST [14], Fashion-MNIST [26], SVHN [20], STL-10 [6], and MedMNIST subsets such as
 169 ChestMNIST/PathMNIST [29, 30]. Full hyperparameters are in Supplementary Methods.

170 5 Discussion

171 SAM makes three testable predictions: (i) instructive/error signals should be vector-valued
 172 at the population level and causally gate plasticity at targeted synapses [1, 7]; (ii) synapses
 173 mediating cross-stream interactions should exhibit learning-linked plasticity in both direc-
 174 tions [4, 5, 9]; and (iii) efficient learning should not require layerwise one-to-one match-
 175 ing between forward and feedback representations, but should instead be compatible with
 176 mismatched pathway widths and learned inter-stream mappings [5, 19]. Across diverse vi-
 177 sion benchmarks, SAM matches or exceeds standard biologically motivated baselines and
 178 is competitive with SGD (Table 1), motivating its consideration as a cortical learning mo-
 179 tif [17, 25]. These predictions can be operationalized as a circuit-level “compatibility test”
 180 (Table 2). Many backprop-inspired models effectively assume a layerwise correspondence be-

Table 2: **Circuit-level falsifiable criteria for cortical credit assignment.** We summarize four pathway types between an upstream (u) and downstream (d) stream: $u \rightarrow u$, $d \rightarrow d$, $u \rightarrow d$, $d \rightarrow u$. A learning proposal specifies both (i) whether each pathway exists at appreciable strength (connectivity) and (ii) whether it is experience-dependent (plasticity). SAM predicts that all four pathways are present and plastic; several prior schemes correspond to strict subsets [3, 11, 16, 17].

| | Connectivity | | | | Plasticity | | | |
|-------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| | $u \rightarrow u$ | $d \rightarrow d$ | $u \rightarrow d$ | $d \rightarrow u$ | $u \rightarrow u$ | $d \rightarrow d$ | $u \rightarrow d$ | $d \rightarrow u$ |
| SGD | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Feedback Alignment [16] | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Weight Mirroring [3] | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ |
| SAM | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

181 tween forward and feedback representations, which typically implies restricted inter-stream
 182 connectivity (often $u \leftrightarrow d$ mappings fixed or absent) and/or restricted plasticity in feed-
 183 back pathways [16, 21, 22]. In contrast, SAM predicts that all four pathway types are both
 184 present and experience-dependent. This yields an eight-criterion experimental program—
 185 four connectivity tests and four plasticity tests—that can distinguish SAM-like learning
 186 from feedback-alignment-like alternatives in identified circuits [16, 17, 25].

187 **Experimental identification and distinction from feedback-alignment family.** A
 188 practical way to distinguish SAM-like circuits from feedback-alignment-like mechanisms is to
 189 test, separately, (a) the existence of each of the four pathway classes and (b) whether synapses
 190 in each class are plastic over learning [5, 19]. In feedback alignment, neuron-specific teaching
 191 depends on a fixed (or non-plastic) feedback scaffold and typically presumes a one-to-one
 192 correspondence between upstream and downstream representations [16, 21]. In SAM, by
 193 contrast, the inter-stream mappings are learnable: the circuit predicts measurable, learning-
 194 dependent changes in cross-stream synapses ($u \rightarrow d$ and $d \rightarrow u$) that accompany improvements
 195 in behavioral performance or task accuracy [1, 4, 7, 9].

196 References

- 197 [1] E. Abs, R. B. Poorthuis, D. Apelblat, K. Muhammad, M. B. Pardi, L. Enke, D. Kushin-
 198 sky, D. L. Pu, M. F. Eizinger, K.-K. Conzelmann, I. Spiegel, and J. J. Letzkus.
 199 Learning-related plasticity in dendrite-targeting layer 1 interneurons. *Science*, 361
 200 (6406):eaat4799, 2018. doi: 10.1126/science.aat4799.
- 201 [2] Nasir Ahmad, Marcel A. J. van Gerven, and Luca Ambrogioni. Gait-prop: A biologically
 202 plausible learning rule derived from backpropagation of error. In *Advances in Neural In-*
 203 *formation Processing Systems*, volume 33, 2020. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper/2020/hash/7ba0691b7777b6581397456412a41390-Abstract.html)
 204 [cc/paper/2020/hash/7ba0691b7777b6581397456412a41390-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/7ba0691b7777b6581397456412a41390-Abstract.html).
- 205 [3] Mohamed Akrouf, Collin Wilson, Peter Conway Humphreys, Timothy P. Lillicrap, and
 206 Douglas B. Tweed. Deep learning without weight transport. In *Advances in Neural*
 207 *Information Processing Systems*, volume 32, 2019.

- 208 [4] Victoria M. Bajo, Fernando R. Nodal, David R. Moore, and Andrew J. King. The de-
209 scending corticocollicular pathway mediates learning-induced auditory plasticity. *Nature*
210 *Neuroscience*, 13(2):253–260, 2010.
- 211 [5] Farran Briggs. Role of feedback connections in central visual processing. *Annual Review*
212 *of Vision Science*, 6(1):313–334, 2020.
- 213 [6] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in
214 unsupervised feature learning. In *Proceedings of the Fourteenth International Conference*
215 *on Artificial Intelligence and Statistics*, pages 215–223. JMLR Workshop and Conference
216 Proceedings, 2011.
- 217 [7] Guy Doron, Jiyun N. Shin, Naoya Takahashi, Moritz Drüke, Christina Bocklisch, Salina
218 Skenderi, Lisa de Mont, Maria Toumazou, Joshua Ledderose, Michael Brecht, Richard
219 Naud, and Matthew E. Larkum. Perirhinal input to neocortical layer 1 controls learning.
220 *Science*, 370(6523):eaaz3136, 2020. doi: 10.1126/science.aaz3136.
- 221 [8] Valerio Francioni, Vincent D. Tang, Enrique H. S. Toloza, Norma J. Brown, and
222 Mark T. Harnett. Vectorized instructive signals in cortical dendrites during a brain-
223 computer interface task. *bioRxiv*, January 2025. doi: 10.1101/2023.11.03.565534.
224 URL <https://www.biorxiv.org/content/10.1101/2023.11.03.565534v2>. Version
225 2 posted January 13, 2025.
- 226 [9] Min Fu and Yi Zuo. Experience-dependent structural plasticity in the cortex. *Trends*
227 *in Neurosciences*, 34(4):177–187, 2011. doi: 10.1016/j.tins.2011.02.001.
- 228 [10] Anthony J. G. D. Holtmaat, Joshua T. Trachtenberg, Linda Wilbrecht, Gordon M. G.
229 Shepherd, Xiao-Qing Zhang, Graham W. Knott, and Karel Svoboda. Transient and
230 persistent dendritic spines in the neocortex in vivo. *Nature*, 435(7045):1020–1025, 2005.
231 doi: 10.1038/nature03720.
- 232 [11] John Kolen and Jordan Pollack. Back-propagation without weight transport. In *Pro-*
233 *ceedings of the IEEE International Conference on Neural Networks (ICNN)*, volume 3,
234 pages 1375–1380. IEEE, 1994.
- 235 [12] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny
236 images. Technical report, University of Toronto, 2009.
- 237 [13] Matthew E. Larkum. A cellular mechanism for cortical associations: an organizing
238 principle for the cerebral cortex. *Trends in Neurosciences*, 36(3):141–151, 2013. doi:
239 10.1016/j.tins.2012.11.006.
- 240 [14] Yann LeCun. The MNIST database of handwritten digits, 1998. Available at [http:](http://yann.lecun.com/exdb/mnist/)
241 [//yann.lecun.com/exdb/mnist/](http://yann.lecun.com/exdb/mnist/).
- 242 [15] Qianli Liao, Joel Leibo, and Tomaso Poggio. How important is weight symmetry in
243 backpropagation?, 2015.

- 244 [16] Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, and Colin J. Akerman. Ran-
245 dom synaptic feedback weights support error backpropagation for deep learning. *Nature*
246 *Communications*, 7(1):13276, 2016.
- 247 [17] Timothy P. Lillicrap, Adam Santoro, Luke Marris, Colin J. Akerman, and Geoffrey
248 Hinton. Backpropagation and the brain. *Nature Reviews Neuroscience*, 21(6):335–346,
249 2020.
- 250 [18] Hiroshi Makino and Takaki Komiyama. Learning enhances the relative impact of top-
251 down processing in the visual cortex. *Nature Neuroscience*, 18(8):1116–1122, 2015. doi:
252 10.1038/nn.4061.
- 253 [19] Nikola T. Markov, Mária Ercsey-Ravasz, David C. Van Essen, Kenneth Knoblauch,
254 Zoltán Toroczkai, and Henry Kennedy. Anatomy of hierarchy: feedforward and feedback
255 pathways in macaque visual cortex. *Journal of Comparative Neurology*, 522(1):225–259,
256 2014.
- 257 [20] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y.
258 Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS*
259 *Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- 260 [21] Arild Nøkland. Direct feedback alignment provides learning in deep neural networks.
261 In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- 262 [22] Arild Nøkland and Lars H. Eidnes. Training neural networks with local error signals.
263 In *Proceedings of the International Conference on Machine Learning*, pages 4839–4850.
264 PMLR, 2019.
- 265 [23] Alexandre Payeur, Jordan Guerguiev, Friedemann Zenke, Blake A. Richards, and
266 Richard Naud. Burst-dependent synaptic plasticity can coordinate learning in
267 hierarchical circuits. *Nature Neuroscience*, 24:1010–1019, 2021. doi: 10.1038/
268 s41593-021-00870-x.
- 269 [24] S. Tattikota, K. K. A. Cho, Y. Liu, W. Hu, V. Barrera, J. L. Stein, and W.-B. Gan.
270 Pyramidal neurons in different cortical layers exhibit distinct dynamics and plasticity
271 of apical dendritic spines. *Frontiers in Neural Circuits*, 11:88, 2017. doi: 10.3389/fncir.
272 2017.00088.
- 273 [25] James C. R. Whittington and Rafal Bogacz. Theories of error back-propagation in the
274 brain. *Trends in Cognitive Sciences*, 23(3):235–250, 2019.
- 275 [26] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset
276 for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 277 [27] Tonghui Xu, Xinzhu Yu, Andrew J. Perlik, Willie F. Tobin, Jonathan A. Zweig, Kelly
278 Tennant, Theresa Jones, and Yi Zuo. Rapid formation and selective stabilization of
279 synapses for enduring motor memories. *Nature*, 462(7275):915–919, 2009. doi: 10.1038/
280 nature08389.

- 281 [28] Guang Yang, Feng Pan, and Wen-Biao Gan. Stably maintained dendritic spines are
282 associated with lifelong memories. *Nature*, 462(7275):920–924, 2009. doi: 10.1038/
283 nature08577.
- 284 [29] Jiancheng Yang, Rui Shi, and Bingbing Ni. MedMNIST classification decathlon: A
285 lightweight AutoML benchmark for medical image analysis. In *IEEE 18th International*
286 *Symposium on Biomedical Imaging (ISBI)*, pages 191–195, 2021.
- 287 [30] Jiancheng Yang, Rui Shi, Dong Wei, Ziyin Liu, Lina Zhao, Bin Ke, Hanspeter Pfister,
288 and Bingbing Ni. MedMNIST v2: A large-scale lightweight benchmark for 2d and 3d
289 biomedical image classification. *Scientific Data*, 10(1):41, 2023.

291 A Circuit and Self-assembly

292 We propose a circuit architecture that is grounded in a simple, biological synaptic motif
 293 (Fig. 1). It is based on two interacting streams of information flow: an ascending stream
 294 representing forward pathways and a descending stream representing feedback pathways,
 295 mirroring cortical connectivity between and within areas. These streams are defined by a set
 296 of synaptic weights, which obey local heterosynaptic plasticity rules. Unlike the standard
 297 computer implementation of gradient descent – backpropagation – our model does not require
 298 the explicit computation of derivatives, and does not require weight symmetry between
 299 forward and feedback pathways.

300 **Notation.** Let θ be the parameters of the model. We will use $\Delta f(\theta)$ to denote the difference
 301 of the quantity $f(\theta)$ after one step of the learning algorithm: $\Delta f(\theta) := f(\theta_{t+1}) - f(\theta_t)$. Any
 302 vector without a transpose superscript is treated as a column vector (e.g. $\nabla F(\theta)$ is a column
 303 vector for a scalar function F).

304 We assume that each pathway is organized in layers indexed with $\ell \in \{1, \dots, L\}$, where L
 305 is the depth of the pathway in terms of modules. Let $h^\ell \in \mathbb{R}^{d_\ell}$ denote the neuron’s cell body
 306 electric potential (i.e., before applying nonlinearity) of the ℓ -th layer upstream pathway and
 307 $\bar{h}^\ell \in \mathbb{R}^{\bar{d}_\ell}$ the downstream pathway. For reasons that will become clear, the upstream is also
 308 referred to as the forward pathway, and the downstream is the feedback pathway. Note that
 309 it is, in general, the case that the two pathways have different numbers of neurons, and so
 310 d_ℓ is not necessarily equal to \bar{d}_ℓ . The simplest such biological motif is shown in Figure ??,
 311 where four synaptic connection matrices connect two consecutive layers of the two pathways.
 312 Here, $W^\ell \in \mathbb{R}^{d_{\ell+1} \times d_\ell}$ is the connection from the ℓ -th layer upstream pathway to the $\ell + 1$ -th
 313 layer upstream pathway, $\bar{W}^\ell \in \mathbb{R}^{\bar{d}_\ell \times \bar{d}_{\ell+1}}$ the connection from $\ell + 1$ -th layer downstream to
 314 the ℓ -th layer downstream. $V^\ell \in \mathbb{R}^{\bar{d}_\ell \times d_\ell}$ and $\bar{V}^\ell \in \mathbb{R}^{d_{\ell+1} \times \bar{d}_{\ell+1}}$ are inter-stream connections:
 315 V goes from the downstream to the upstream, and \bar{V} from the upstream to the downstream.
 316 The algorithm’s Pseudo-code is in Appendix D.

317 The computation of the two pathways takes place sequentially, where the upstream com-
 318 putation is performed first and eventually used for inference, and the downstream computa-
 319 tion takes place after inference and is used for synaptic weight updates:²

$$h^{\ell+1} = W^\ell D_u^\ell h^\ell \quad (5)$$

$$\bar{h}^\ell = \bar{W}^\ell D_d^{\ell+1} \bar{h}^{\ell+1}, \quad (6)$$

320 where D matrices are diagonal matrix that functions as the neuron nonlinearity and may be
 321 different for upstream (D_u , corresponding to h) and downstream (D_d , corresponding to \bar{h})
 322 pathways and for different layers. The first layer of $h^1 := x$ is the input data, and $h^L := \hat{y}$ is

²This sequential ordering is consistent with, for example, the observation that inactivating V1 also leads to the inactivation of V2, while deactivating V2 has no effect on the activations of V1 (Anderson & Martin, 2009).

323 the output of the network. The input to the downstream pathway is an error signal obtained
 324 from an objective function F : $\bar{h}^L := \epsilon(\hat{y}) = -\nabla_{\hat{y}} F(\hat{y})$, which is computed from the output
 325 of the upstream pathway. The upstream propagates this signal in the reverse direction. One
 326 example of F is the MSE loss, where $F(\hat{y}) = \|y - \hat{y}\|^2$, where y is the correct label.

327 For the first layer, D is identity (i.e., no nonlinearity on the input). For later layers, we
 328 choose D_u to be a diagonal zero-one matrix:

$$D_u^\ell = \text{diag}(\mathbb{1}_{(h^\ell)_1}, \dots, \mathbb{1}_{(h^\ell)_d}), \quad (7)$$

329 where $\mathbb{1}_x = 1$ if $x > 0$ and is zero if $x \leq 0$, corresponding to a ReLU nonlinearity. For
 330 D_d , we focus on the simplest case where $D_d = I$ is the identity matrix.

331 Another choice of D_d is the Jacobian of a ReLU nonlinearity based on forward activa-
 332 tions, which requires $V = \bar{V} = I$ (namely a one-to-one correspondence between upstream
 333 and downstream neurons). With this choice of D_d , V and \bar{V} , our algorithm reduces to:
 334 1. standard backpropagation, if \bar{W} is hardwired to W throughout training. 2. Feedback
 335 Alignment (Lillicrap et al., 2016), if \bar{W} is random and fixed. 3. the KP-algorithm (Kolen &
 336 Pollack, 1994), if W is learned using our learning rule, which converges to the exact gradient
 337 descent provably.

338 **Heterosynaptic Learning:** our local learning rule takes the general form of $\Delta W \propto$
 339 $S_{post} * S_{pre}$, where the change ΔW of a synapse is proportional to the product of postsynaptic
 340 signal S_{post} and presynaptic signal S_{pre} . This is similar to Hebbian learning, but we have
 341 extra emphasis that S_{post} refers to the signals sent by feedback synapses to the postsynaptic
 342 neuron instead of those arising from self-activation due to forward synapses. Feedback signals
 343 are informative in learning while self-activation signals only lead to non-informative arbitrary
 344 self-strengthening.

345 In most situations, S_{pre} represents the presynaptic potential on an axon, mathematically
 346 taking the form of $D_{pre} * h_{pre}$, where D_{pre} is the nonlinearity and h_{pre} is the electric potential
 347 of the presynaptic cell body. The nature of S_{post} on the other hand is an open question.
 348 There are at least two possibilities: (1) S_{post} is simply proportional to the postsynaptic
 349 potential h_{post} and is not affected by postsynaptic neuron’s activity (2) S_{post} is proportional
 350 to the postsynaptic potential modulated by postsynaptic neuron’s activity, namely $S_{post} \propto$
 351 $D_{post} * h_{post}$. This makes the learning rule $\Delta W \propto (D_{post} h_{post}) * (D_{pre} h_{pre})$, which is even
 352 more symmetric.

353 This simple rule guides the learning of all four directions of synapses in our model:

$$\Delta W^\ell = \eta_W (D_u^{\ell+1} (\bar{V}^\ell D_d^{\ell+1} \bar{h}^{\ell+1})) (D_u^\ell h^\ell)^T - \gamma W^\ell \quad \propto (D_{W_{post}} S_{W_{post}}) (D_{W_{pre}} S_{W_{pre}})^T, \quad (8)$$

$$\Delta \bar{W}^\ell = \eta_{\bar{W}} (D_d^\ell (V^\ell D_u^\ell h^\ell)) (D_d^{\ell+1} \bar{h}^{\ell+1})^T - \gamma \bar{W}^\ell \quad \propto (D_{\bar{W}_{post}} S_{\bar{W}_{post}}) (D_{\bar{W}_{pre}} S_{\bar{W}_{pre}})^T, \quad (9)$$

$$\Delta V^\ell = \eta_V (D_d^\ell \bar{h}^\ell) (D_u^\ell h^\ell)^T - \gamma V^\ell \quad \propto (D_{V_{post}} S_{V_{post}}) (D_{V_{pre}} S_{V_{pre}})^T, \quad (10)$$

$$\Delta \bar{V}^\ell = \eta_{\bar{V}} (D_u^{\ell+1} h^{\ell+1}) (D_d^{\ell+1} \bar{h}^{\ell+1})^T - \gamma \bar{V}^\ell \quad \propto (D_{\bar{V}_{post}} S_{\bar{V}_{post}}) (D_{\bar{V}_{pre}} S_{\bar{V}_{pre}})^T, \quad (11)$$

354 where η is the time constant of learning (i.e., the learning rate), γ is the strength of local
 355 regularization terms that encourage the synapses to be weak when unused. The four sets of
 356 rules are almost completely symmetric along all directions – the only asymmetry between
 357 W/\bar{W} and V/\bar{V} stems from the fact that W and \bar{W} have the overwriting effect on the

neurons that they connect to, while V and \bar{V} only serve to provide learning signals but do not overwrite the target neurons’ activations. This design choice may be explored further in the future. We further discuss the biological evidences for our learning rules in Appendix B.

Both the learning rate and the regularization strength are labeled with a different subscript to emphasize that they could have different time constants. In this work, we always keep the regularization strengths uniform. We empirically search over 64 combinations of the learning rates and use the best one in the experiment.³ We found it particularly beneficial to fix the ratio between the four learning rates, while tuning the overall factor for different tasks.

A key feature of this architecture is that it allows both the inference network and the learning circuit, in addition to all their connections, to be simultaneously self-assembled. We thus name this algorithm ”Self-Assembling Motif” (**SAM**).

B Gradient Learning in Overparametrized Downstream Pathway

The circuits we describe can be viewed as a generalized version of SGD, with a matrix form and learnable learning rate – in this sense, the algorithm is learning the learning algorithm itself. The following theorem is an informal statement of this result. The full formal detail and proof are provided in Appendix A.

Theorem 2. *Consider a ReLU neural network with an arbitrary width and depth and with an overparametrized downstream pathway. If both V^ℓ and \bar{V}^ℓ are full-rank for all ℓ , then, for any x such that for all $\ell \in [L]$, $\Delta V^\ell = O(\epsilon)$ and $\Delta \bar{V}^\ell = O(\epsilon)$,*

$$\Delta_{\text{SAM}}W^\ell = H\Delta_{\text{sgd}}W^\ell - \gamma W^\ell + O(\epsilon), \quad (12)$$

for a positive definite matrix $H = \bar{V}^\ell(\bar{V}^\ell)^T$, and $\Delta_{\text{sgd}}W^\ell = -\nabla_W F$ is the SGD update.

Remark. This theorem essentially shows that when V and \bar{V} reach stationarity, the algorithm will run like a SGD algorithm with a matrix learning rate H and weight decay. The fact that matrix H is PSD is crucial from a mathematical perspective, as it guarantees that the loss objective of training will decrease in expectation and that the stationary points of the circuit will be identical to that of SGD. For example, besides SGD itself, Adam or Natural Gradient Descent can also be seen as having a learning rate corresponding to a PSD matrix.

It is worth noticing that this algorithm is different from SGD. H is not only a layer-wise matrix learning rate: it is explicitly dependent on time and is also being learned by the algorithm! This is consistent with observed metaplasticity in biological systems (Abraham & Bear, 1996; Abraham, 2008) – though it extends the idea in a fundamental way since the learning algorithm itself is being learned during the training of the forward network. In Figure 2a, we show examples of the H matrix after training.

³See Section E.

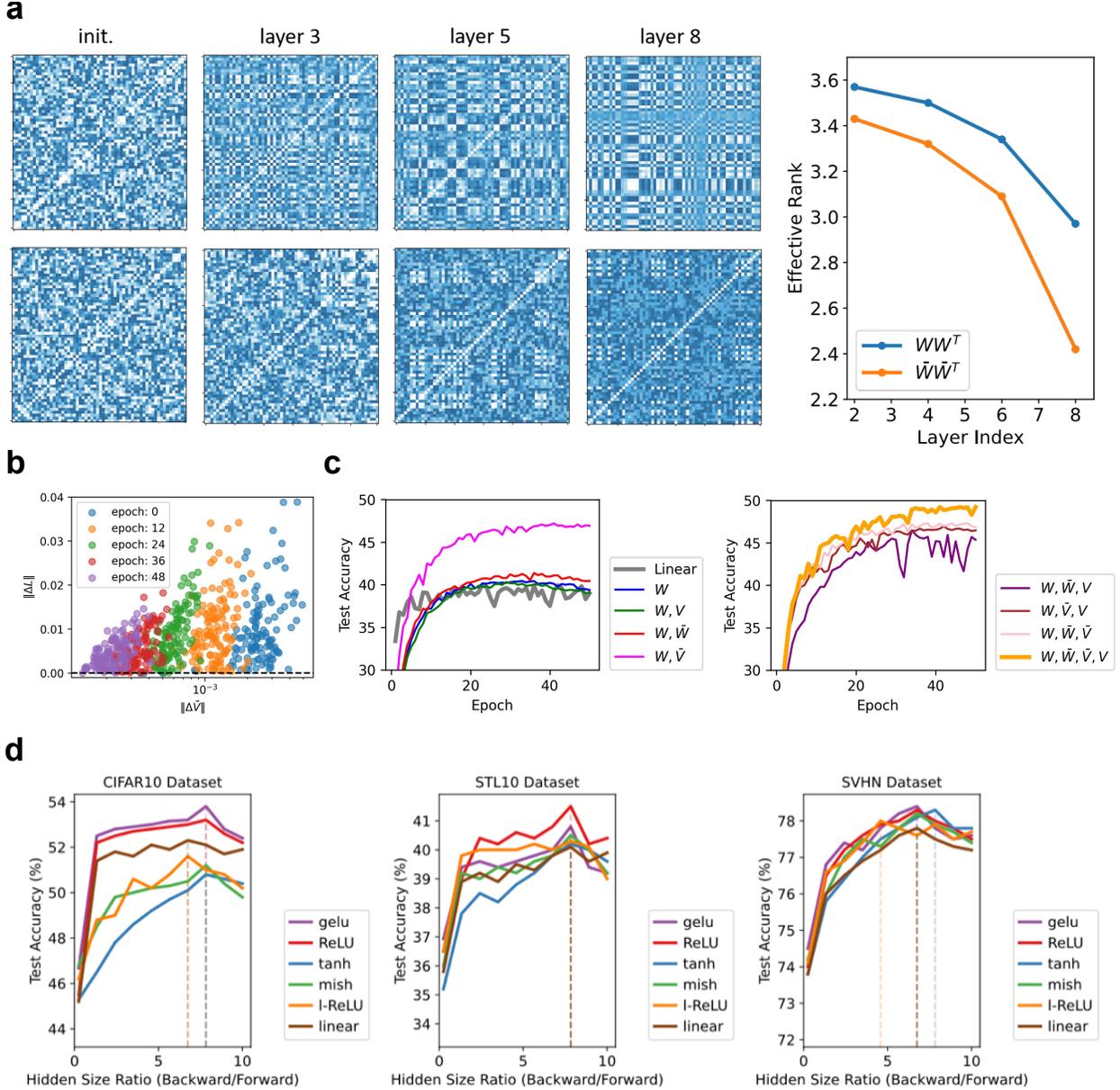


Figure S1: **Detailed analyses, variants and ablation studies of SAM.** (a) The matrices $\bar{V}\bar{V}^T$ (upper) and WW^T (lower) before and after training. Also, recall that $\bar{V}\bar{V}^T = H$ is the matrix learning rate for W (Theorem 1). The neurons within the same layer become correlated after training. The spectra of the weight matrices of the two pathway are interestingly found to be effectively low-rank (right). This visualization was done using the SVHN dataset with a network having 8 hidden layers. (b) Both \bar{V} and ΔL becomes smaller and smaller over time, a condition required for the algorithm to simulate gradient descent (Theorem 1). (c) Ablation study: see caption in main text. (d) The impact of the width of the downstream pathway on performance. We use different activation functions in the forward pathway. The width of the upstream pathway was kept constant at 200, while varying the width of the downstream pathway. We find that having a wider downstream network improves the performance of the feedforward network up to an overparametrization ratio of 7.5.

392 Another interesting aspect of the algorithm (see Theorem 1) is that it allows (but does
 393 not require) the backward pathway to be overparametrized. That is, the number of down-
 394 stream connections can be larger than the number of forward connections. This is consistent
 395 with observations in biological circuits (Briggs, 2020), but it is inconsistent with previous
 396 biologically inspired algorithms of learning, which require the same number of connections
 397 for ascending and descending pathways. Also, it turns out that the backward pathway is
 398 highly flexible in the choice of the activation; with any choice of D_d , the theorem holds.
 399 This prediction of the robustness of the backward pathway is numerically supported by the
 400 experiment in Section C.2.

401 When is it possible to satisfy the condition that V and \bar{V} are stationary? There are many
 402 possible ways. Two notable cases are when V and \bar{V} are permutation matrices, and D_d is
 403 the Jacobian matrix of D_u . In this case, H is identity, and so our algorithm is completely
 404 identical to SGD. In a different case, the feedback pathway is very expressive (e.g., by having
 405 a large number of feedback neurons), and so it is capable, from an approximation perspective,
 406 of approximating the forward net. Figure 2b shows how $\Delta\bar{V}$ approaches zero during training,
 407 and how this accompanies the decrease and $\|\Delta L\|$, consistent with the theory.

408 C Main Observations

409 C.1 SAM is an Efficient Learning Algorithm

410 We evaluate our algorithm on seven widely used image classification datasets: CIFAR10
 411 (Krizhevsky et al., 2009), MNIST (Yann, 1998), Fashion MNIST (Xiao et al., 2017), ChestM-
 412 NIST (Netzer et al., 2011), PathMNIST (Yang et al., 2021, 2023), SVHN (Netzer et al.,
 413 2011), and STL10 (Coates et al., 2011). For details about these datasets, see Appendix C.
 414 We compare our algorithm’s performance against the standard backpropagation and other
 415 biologically plausible algorithms. These include Feedback Alignment (FA) (Lillicrap et al.,
 416 2016), which uses random fixed weights for error propagation instead of transposed weights;
 417 Weight Mirroring (WM) (Akrouf et al., 2019a), a method enforcing symmetry between the
 418 forward and backward passes by mirroring weights; and the KP algorithm (Akrouf et al.,
 419 2019a), which adjusts synaptic weights without requiring weight transport, simplifying net-
 420 work architecture and enhancing learning in large networks. While the primary goal of this
 421 paper is not to pursue State-of-the-Art (SOTA) performance but to design a neural network
 422 algorithm that can be deployed in the brain, our approach achieves SOTA among biologically
 423 plausible algorithms on all datasets (see in Table 1). Notably, our algorithm outperforms
 424 models trained with backpropagation on six out of the seven datasets. The only exception
 425 is ChestMNIST, where it performs 0.02% lower than backpropagation.

426 C.2 Overparameterization in the Downstream Pathway is Benefi- 427 cial

428 In the nervous system, the number of backprojections can be much larger than that of the
 429 forward stream. For example, the feedback connections V1 to LGN is can be roughly 10

430 times the RGC synapses from the retina (Briggs, 2020). We have observed a potential com-
 431 putational explanation of this phenomenon in our simulations. Our models perform better as
 432 the feedback pathway becomes overparametrized and achieves the best performance across
 433 multiple datasets when the parameter count of the downstream pathway is approximately 5
 434 to 8 times that of the upstream pathway. Specifically, as shown in Figure 2d for CIFAR10,
 435 performance peaks when the parameter count of the downstream pathway is 8 times that
 436 of the upstream pathway. We perform experiments on the three most challenging datasets
 437 from the seven datasets used in Table 1: CIFAR10 (Krizhevsky et al., 2009), SVHN (Netzer
 438 et al., 2011), and STL10 (Coates et al., 2011).

439 C.3 Learned Representations are Structurally Low-Rank

440 Now, we show some examples of the learned weight matrices. We train a width-60 network
 441 on SVHN dataset for 50 epochs. See Figure 2a. As is common in deep learning, we plot
 442 the matrix WW^T and also $\bar{V}\bar{V}^T$. We see that before training, these matrices are primarily
 443 diagonal matrices due to i.i.d. Gaussian initialization. After training, correlations between
 444 neurons and weights appear. We characterize the rank of a matrix through a continuous
 445 metric referred to as the effective rank (Roy & Vetterli, 2007). We find that these weight
 446 matrices also become effectively low-rank.

447 Figure 2a can be compared with similar matrices trained with SGD in Ziyin et al. (2024)
 448 and Huh et al. (2021), for example, and they seem to have similar high-level structures,
 449 suggesting that our algorithm can lead to a meaningful representation learning. Detailed
 450 analysis of the learned representation using our algorithm may be an interesting future
 451 problem.

452 C.4 Interconnection Plasticity is Crucial

453 A remaining question is whether learning at some of the interconnections is more important
 454 than at others, and if so, which is more important? Additionally, are any interconnections
 455 not important at all? We perform this ablation study on CIFAR-10, where we initialize
 456 all the parameters randomly but only update a subset of all the connections. See Figure
 457 2c. Note that W is always plastic. To gauge how much the network has learned about
 458 the nonlinear relations in the data, we also compare with the learning trajectory of a linear
 459 model trained with SGD.

460 We discover the following phenomena:

- 461 • Only updating W works but only works as well as the linear model; making other parts
 462 plastic is thus important for learning nonlinear functions;
- 463 • If only one additional connection is made plastic, \bar{V} is the most important one, improving
 464 performance from 40% to roughly 46%;
- 465 • Making \bar{V} plastic is similar to making W , V plastic together, suggesting that some loss of
 466 learning capabilities due to connection lesions can be compensated by other connections
 467 (but not all);

468 • Making all four connections plastic improves the performance further, to the SGD level
 469 (cf. Table 1); this means that all connections need to be plastic to achieve the best
 470 performance.

471 Therefore, none of the pathways we introduced are redundant, and this partially explains the
 472 outperformance of our proposed algorithm over the existing two-pathway biologically plau-
 473 sible learning algorithms. Another reflection of the fact that there is an asymmetry between
 474 different pathways is that different connections require different learning time constants to
 475 work. In particular, the \bar{V} matrices require the largest learning rates, usually four times that
 476 of W (see Appendix).

477 D Theorems and Proofs

478 The following Lemma shows that the interstream connections evolve towards becoming
 479 aligned with each other after training.

480 **Lemma 1.** *If $\Delta V^\ell = 0$ and $\Delta \bar{V}^\ell = 0$ for all l , then*

$$\bar{V}^{\ell-1} = (V^\ell)^T. \quad (13)$$

481 *Proof.* This is a consequence of the dynamics of $\bar{V}^{\ell-1}$ and \bar{V}^ℓ being essentially identical. By
 482 definition, we have that

$$\Delta V^\ell = \bar{h}^\ell (h^\ell)^T - \gamma V^\ell, \quad (14)$$

$$\Delta (\bar{V}^{\ell-1})^T = \bar{h}^\ell (h^\ell)^T - \gamma (\bar{V}^{\ell-1})^T. \quad (15)$$

483 This implies that

$$\Delta (V^\ell - (V^{\ell-1})^T) = -\gamma (V^\ell - (V^{\ell-1})^T), \quad (16)$$

484 which is an exponential decay to zero. This finishes the proof. \square

485 The following definition states what it means for the downstream pathway to be “over-
 486 parametrized.”

487 **Definition 1.** *The downstream pathway is said to be overparametrized if for all l and $V^\ell \in$
 488 $\mathbb{R}^{\bar{d}_\ell \times d_\ell}$,*

$$d_\ell \leq \bar{d}_\ell. \quad (17)$$

489 *We also say that the matrix V^ℓ is overparametrized if this condition holds.*

490 The following Lemma is a technical step in the theorem proof.

491 **Lemma 2.** *If $V^\ell D_u^\ell (W^\ell)^T = \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T$, then,*

$$P \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T = \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T, \quad (18)$$

492 *where $P = ((V^\ell)^T)^+ (V^\ell)^T$ is a projection matrix and A^+ denotes the pseudoinverse of the
 493 matrix A .*

494 *Proof.* We have

$$P\bar{W}^\ell D_d^{\ell+1}(\bar{V}^\ell)^T = PV^\ell D_u^\ell(W^\ell)^T = V^\ell D_u^\ell(W^\ell)^T = \bar{W}^\ell D_d^{\ell+1}(\bar{V}^\ell)^T. \quad (19)$$

495 This finishes the proof. \square

496 Now, we explain briefly that the sample-wise gradient of a ReLU network can be written
 497 in a specific form. As implied in the main text, the output of a ReLU network can be written
 498 in the following form

$$f(x) = W^L D_u^{L-1} \dots D^1 W^1 x \quad (20)$$

499 where D_u^ℓ is the Jacobian of the ReLU activation at the i -th layer.

500 Let $F(f(x), y)$ be the loss function for the data point x and label y . The above discussion
 501 implies that the gradient of F with respect to W can be written as

$$\nabla_{W^\ell} F = \underbrace{(\nabla_f^T F W^L D_u^{L-1} \dots D^\ell)^T}_{\text{error signal}} \underbrace{((D_u^{\ell-1}) W^{\ell-1} \dots W_1 x)^T}_{\text{forward signal}}. \quad (21)$$

502 We will show that the downstream pathway will become self-assembled in a way that it
 503 computes the error signal.

504 Now, we are ready to prove the main theorem. For the ease of reference, we state
 505 the theorem again here.

506 **Theorem 3.** *Consider a ReLU neural network with an arbitrary width and depth and with*
 507 *an overparametrized downstream pathway. If both V^ℓ and \bar{V}^ℓ are full-rank for all ℓ , then, for*
 508 *any x such that for all $\ell \in [L]$, $\Delta V^\ell = O(\epsilon)$ and $\Delta \bar{V}^\ell = O(\epsilon)$,*

$$\Delta_{\text{SAM}} W^\ell = H \Delta_{\text{sgd}} W^\ell - \gamma W^\ell + O(\epsilon), \quad (22)$$

509 for a positive definite matrix $H = \bar{V}^\ell (\bar{V}^\ell)^T$, and $\Delta_{\text{sgd}} W^\ell = -\nabla_W F$ is the SGD update.

510 *Proof.* By definition of the algorithm,

$$\Delta V^\ell = \bar{h}^\ell (h^\ell)^T - \gamma V^\ell = \bar{W}^\ell D_d^{\ell+1} \bar{h}^{\ell+1} (h^\ell)^T - \gamma V^\ell, \quad (23)$$

$$\Delta(\bar{V}^\ell)^T = \bar{h}^{\ell+1} (h^{\ell+1})^T - \gamma (\bar{V}^\ell)^T = \bar{h}^{\ell+1} (h^\ell)^T D_u^\ell (W^\ell)^T - \gamma (\bar{V}^\ell)^T. \quad (24)$$

511 If $\Delta V^\ell = O(\epsilon)$ and $\Delta \bar{V}^\ell = O(\epsilon)$, we have

$$(\Delta V^\ell) D_u^\ell (W^\ell)^T = O(\epsilon), \quad (25)$$

512

$$\bar{W}^\ell D_d^{\ell+1} \Delta(\bar{V}^\ell)^T = O(\epsilon). \quad (26)$$

513 Thus, we have

$$V^\ell D_u^\ell (W^\ell)^T = \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T + O(\epsilon). \quad (27)$$

514 Because V^ℓ is full-rank and overparametrized, $G_\ell := (V^\ell)^T V^\ell$ is an invertible matrix. We
 515 can thus multiply $G_\ell^{-1} (V^\ell)^T$ on the left to obtain that

$$D_u^\ell (W^\ell)^T = G_\ell^{-1} (V^\ell)^T \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T + O(\epsilon). \quad (28)$$

516 Now, consider the matrix product between different layers⁴

$$(W^\ell D_u^\ell W^{\ell-1} D_u^{\ell-1})^T = (V^{\ell-1})^+ \bar{W}^{\ell-1} D_d^\ell (\bar{V}^{\ell-1})^T G_\ell^{-1} (V^\ell)^T \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T + O(\epsilon). \quad (29)$$

517 However, by Lemma 1, at stationarity, we must have

$$\bar{V}^{\ell-1} = (V^\ell)^T, \quad (30)$$

518 and so

$$(\bar{V}^{\ell-1})^T G_\ell^{-1} (V^\ell)^T = V^\ell G_\ell^{-1} (V^\ell)^T = P^2 = P, \quad (31)$$

519 where $P = ((V^\ell)^T)^+ (V^\ell)^T$ is a projection matrix, and the $+$ superscript denotes the pseudo-
520 doinverse. Thus,

$$(W^\ell D_u^\ell W^{\ell-1} D_u^{\ell-1})^T = (V^{\ell-1})^+ \bar{W}^{\ell-1} D_d^\ell P \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T + O(\epsilon). \quad (32)$$

521 Using Lemma 2, we have that

$$P D_d^\ell \bar{W}^\ell (\bar{V}^\ell)^T = D_d^\ell \bar{W}^\ell (\bar{V}^\ell)^T, \quad (33)$$

522 which implies that

$$(W^\ell D_u^\ell W^{\ell-1} D_u^{\ell-1})^T = (V^{\ell-1})^+ \bar{W}^{\ell-1} D_d^\ell \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T + O(2\epsilon). \quad (34)$$

523 Namely, the intermediate projections P have no effect. This argument can be applied repet-
524 itively to show that

$$(W^\ell D_u^\ell \dots W^{\ell-n} D_u^{\ell-n})^T = (V^{\ell-n})^+ \bar{W}^{\ell-n} D_d^{\ell-n+1} \dots \bar{W}^\ell D_d^{\ell+1} (\bar{V}^\ell)^T + O((\ell-n)\epsilon) \quad (35)$$

525 for any l and $n < l$. This means that the backward net has essentially become the transpose
526 of the forward net.

527 Now, by the definition of our algorithm, we have

$$\bar{V}^\ell \bar{h}^{\ell+1} = \bar{V}^\ell \bar{W}^{\ell+1} D_d^{\ell+2} \dots \bar{W}^{L-1} D_d^L \bar{h}^L \quad (36)$$

$$= \bar{V}^\ell V^{\ell+1} (W^{L-1} D_u^{L-1} \dots W^\ell D_u^\ell)^T ((\bar{V}^\ell)^T)^{-1} \bar{h}^L \quad (37)$$

$$= \bar{V}^\ell (\bar{V}^\ell)^T (W^{L-1} D_u^{L-1} \dots W^\ell D_u^\ell)^T ((\bar{V}^\ell)^T)^{-1} \bar{h}^L \quad (38)$$

$$= Z (W^{L-1} D_u^{L-1} \dots W^\ell D_u^\ell)^T \bar{h}^L + O(\ell\epsilon), \quad (39)$$

528 where $Z = \bar{V}^\ell (\bar{V}^\ell)^T$ is positive definite, we have used the definition $\bar{V}^L = I$. By definition,
529 \bar{h}^L is the label gradient.

530 Therefore,

$$D_u^\ell \bar{V}^\ell \bar{h}^{\ell+1} (h^\ell)^T D_u^\ell = \underbrace{D_u^\ell Z D_u^\ell}_{\text{effective learning rate: } H} D_u^\ell (W^\ell)^T D_u^{L-1} (W^{L-1})^T \bar{h}^L (h^\ell)^T D_u^\ell + O(\ell\epsilon), \quad (40)$$

$$= H \nabla_{h^{\ell+1}} F \quad (41)$$

531 The proof is complete. □

532

⁴Note that $(V^\ell)^+ = G_\ell^{-1} (V^\ell)^T$.

533 Note that this update rule in Eq. (40) is zero if and only if the SGD update is zero. It
 534 thus has the same stationary points as SGD. To see this, note that D_u^ℓ is a projection matrix,
 535 and Z is full-rank and symmetric. Thus, the null space of H is the same as the null space
 536 of D_u^ℓ . Therefore, if $\nabla_{h^{\ell+1}} F \neq 0$, it must be the case that $H\nabla_{h^{\ell+1}} F \neq 0$.

537 **Proposition 1.** *Let Φ denote the empirical loss. If H is positive semi-definite, and if we*
 538 *update the parameters by*

$$\dot{\theta} = -H\nabla_{\theta}\Phi, \tag{42}$$

539 *then, $\Phi(\theta(t))$ is a monotonic decreasing function of t for any initialization $\theta(0)$.*

540 *Proof.* For the first part, consider the time evolution of

$$\dot{\Phi} = (\nabla\Phi)^T \dot{\theta} = -(\nabla\Phi)^T H \nabla\Phi \leq 0. \tag{43}$$

541 □

542 E Experimental Setup

543 All experiments are conducted with PyTorch on one NVIDIA A100 80GB GPU. The batch
 544 size is set to 256.

545 **Datasets.** We use six datasets to evaluate our method. *CIFAR-10* (Krizhevsky et al.,
 546 2009) consists of 60,000 color images in 10 different classes, with each image having a reso-
 547 lution of 32x32 pixels. The dataset is divided into 50,000 training images and 10,000 testing
 548 images. The classes are mutually exclusive and include common objects such as airplanes,
 549 cars, cats, and dogs.

550 The *MNIST* (Yann, 1998) contains 70,000 grayscale images of handwritten digits (0-9),
 551 each with a resolution of 28x28 pixels. The dataset is divided into 60,000 training images and
 552 10,000 testing images. MNIST is known for its simplicity and has been a standard dataset
 553 for testing machine learning algorithms.

554 *Fashion MNIST* (Xiao et al., 2017) contains 70,000 grayscale images of clothing items
 555 like T-shirts, trousers, shoes, and bags, with each image having a resolution of 28x28 pixels.
 556 The dataset is also split into 60,000 training and 10,000 testing images. Fashion MNIST
 557 is considered more challenging than MNIST due to the variability in the visual features of
 558 clothing items.

559 The *SVHN* (Netzer et al., 2011) contains over 600,000 color images of digits (0-9), each
 560 with a resolution of 32x32 pixels. The dataset is divided into a training set of 73,257 images,
 561 a testing set of 26,032 images, and an additional 531,131 images for extra training. SVHN
 562 is challenging due to the varying digit sizes, orientations, and complex backgrounds.

563 The ChestMNIST (Yang et al., 2021, 2023) comprises chest X-ray images for multi-
 564 label classification tasks. Derived from the NIH ChestX-ray14 dataset, it includes images
 565 annotated with 14 different pathological labels, such as pneumonia and emphysema. It is
 566 widely used in medical imaging studies to develop and evaluate machine learning models
 567 capable of diagnosing multiple conditions from X-ray images.

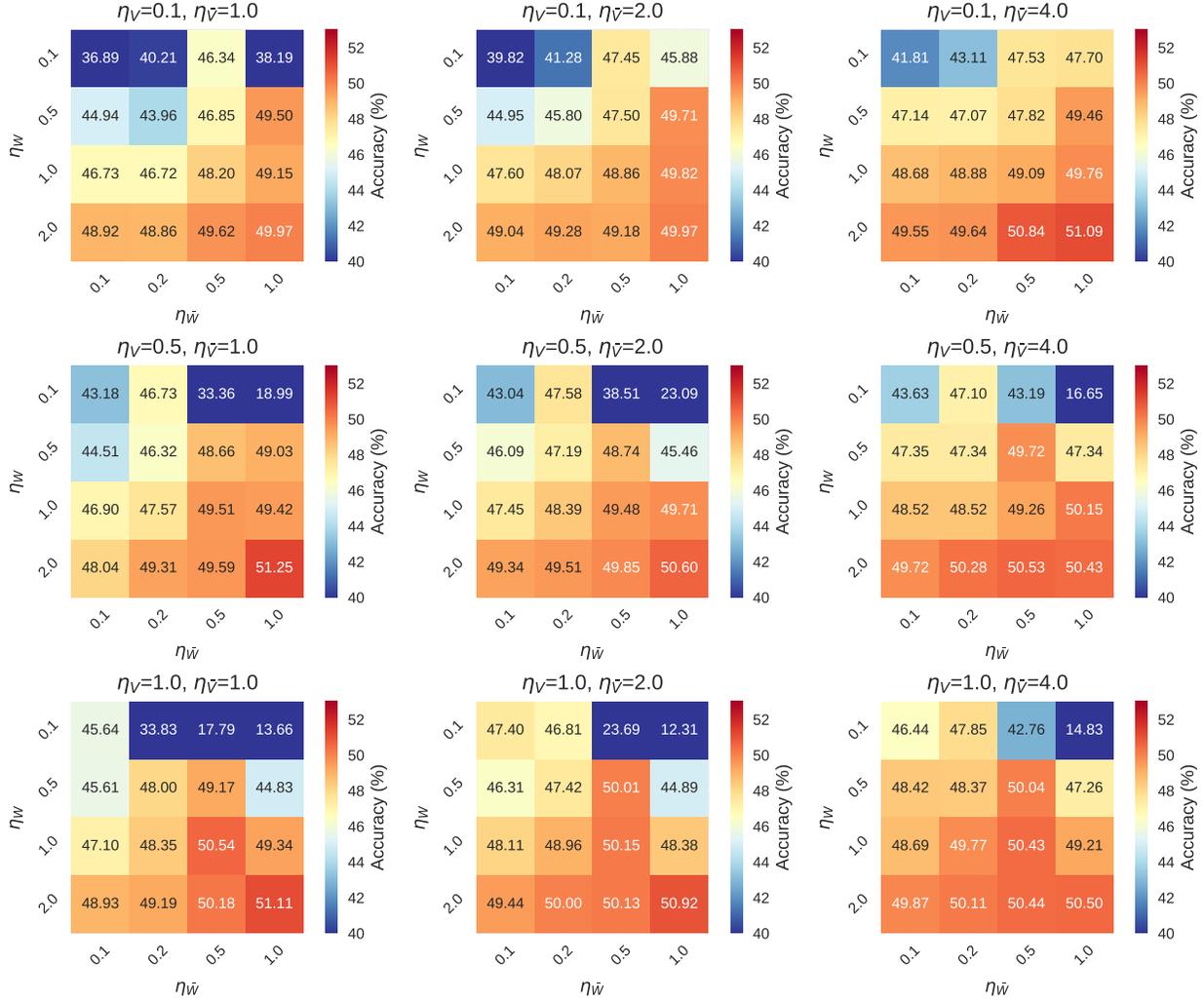


Figure S2: Performance of a four-layer MLP for different choices of learning rates. Interestingly, the best performance is achieved when η_V is the smallest.

568 PathMNIST (Yang et al., 2021, 2023) is a pathology image dataset designed for multi-
 569 class classification tasks. It originates from the NCT-CRC-HE-100K dataset and includes
 570 histological images of nine different types of tissue from colorectal cancer samples, such as
 571 adenocarcinoma and lymph node tissue. This dataset is used primarily for training models to
 572 differentiate between various cancerous tissue types based on their histopathological features.

573 *STL-10* (Coates et al., 2011) consists of 13,000 labeled images across 10 different classes
 574 and 100,000 unlabeled images. Each image has a resolution of 96x96 pixels, which is signifi-
 575 cantly larger than CIFAR-10 images, making the dataset more challenging. The dataset has
 576 a training set of 5,000 labeled images and a test set of 8,000 labeled images.

577 **Baselines.** We compare our method with four baselines, including SGD and three biologically-
 578 motivated algorithms. In the four baseline we assume that there is one feedback pathway for
 579 each feedforward pathway – a strong assumption on developmental mechanisms that SAM

580 does not require.

581 *Stochastic Gradient Descent (SGD)* is the standard optimization method that updates
582 parameters iteratively using a subset of data to minimize the objective function.

583 *Weight Mirroring* (Akrouf et al., 2019a) is a technique generally used in neural networks
584 to enforce symmetry between forward and backward passes. This can be implemented by
585 mirroring the weights used in the forward pass for use in the backward pass, which often
586 helps with learning efficiency and stability.

587 *Feedback Alignment* (Lillicrap et al., 2016) offers an alternative to traditional backprop-
588 agation by using random, fixed weights instead of transposed forward-pass weights to prop-
589 agate error signals.

590 *Kolen-Pollack (KP) algorithm* (Kolen & Pollack, 1994) enables synaptic weight adjust-
591 ments in neural networks without weight transport, simplifying the architecture and enhanc-
592 ing learning in large networks.