
Beyond Directed Acyclic Computation Graph with Cyclic Neural Network

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper investigates a fundamental yet overlooked design principle of artificial
2 neural networks (ANN): We do not need to build ANNs layer-by-layer sequentially
3 to guarantee the Directed Acyclic Graph (DAG) property. Inspired by biological in-
4 telligence, where neurons form a complex, graph-structured network, we introduce
5 the transformative Cyclic Neural Networks (Cyclic NN). It emulates biological
6 neural systems' flexible and dynamic graph nature, allowing neuron connections in
7 any graph-like structure, including cycles. This offers greater flexibility compared
8 to the DAG structure of current ANNs. We further develop the Graph Over Multi-
9 layer Perceptron, the first detailed model based on this new design paradigm. We
10 experimentally validate the advantages of Cyclic NN on widely tested datasets in
11 most generalized cases, demonstrating its superiority over current layer-by-layer
12 DAG neural networks. With the support of Cyclic NN, the Forward-Forward training
13 algorithm also firstly outperforms the current Back-Propagation algorithm. This
14 research illustrates a transformative ANN design paradigm, a significant departure
15 from current ANN designs, potentially leading to more biologically similar ANNs.

16 1 Introduction

17 Artificial intelligence (AI) has reshaped our daily lives and is expected to have a much greater impact
18 in the foreseeable future. Lying behind the most profound AI applications [19, 15, 16, 8], artificial
19 neural networks (ANN) such as multi-layer perceptron (MLP) [17], convolution neural network
20 (CNN) [12] and Transformer [21] are designed specifically for different domains to fit the training
21 data. Regardless of the network structure, neural networks are stacked layer-by-layer to form deep
22 ANNs for greater learning capacity. It has been a *de facto* practice until now that data is first fed into
23 the input layer and then propagated through all the stacked layers to obtain the final representations at
24 the output layer. This paper seeks to answer a fundamental question in ANNs: "Do we really need to
25 stack neural networks layer-by-layer sequentially?"

26 To answer this question, let's first examine the evidence from biological intelligence (BI). Neurosci-
27 entists have studied the biological neurons for decades. The connectome of *C. elegans* is the most
28 thoroughly studied biological neural system, and biologists depicted the most detailed connection
29 between 302 biological neurons [23, 4] as shown in Figure 1(a). Rather than being stacked layer-
30 by-layer, all the neurons form a complicated connection graph, where each can connect to several
31 other neurons within the system. We cannot even determine which neuron serves as the input/output
32 within the neural system to process information. The same findings have also been observed in
33 the latter more complicated neural systems, such as the biology neural connectome of drosophila
34 larva [24], zebrafish [3], mouse [20] and the human brain [18]. Observed biological intelligence
35 exhibits graph-structured, flexible, and dynamic neural systems, which are apparently different from
36 the current layer-by-layer ANNs we build nowadays, as depicted in Figure 1(b).

37 The learning rules actually cause the difference
 38 in the neural system structure between BI and
 39 AI. The Hebb’s Rule [6], depicted as “Neurons
 40 that fire together wire together”, is recognized as
 41 the fundamental learning way of biological neurons.
 42 The Spike-Timing-Dependent Plasticity
 43 (STDP) learning is then proposed to further consider
 44 the relative spiking time of pre-synapse and
 45 post-synapse neurons. Both learning rules of BI
 46 are localized, *i.e.*, the learning occurs on each
 47 neuron within its local influence scope. The
 48 localized learning rules grant the flexibility of
 49 each neuron on its connections to other neurons,
 50 which leads to the complicated graph-structured BI system.
 51 Conversely, for AI systems, the backward
 52 propagation (BP) algorithm [17] has dominated the
 53 training of ANNs. Data is fed into the ANNs
 54 from the input layer, forward propagates layer
 55 by layer to the last layer, calculates a global
 56 loss for the whole ANN based on the ground-truth
 57 labels, and then reversely backward propagates
 58 the error signals layer by layer to the input
 59 layer. In this procedure, ANNs are trained by a
 60 global loss function, and the ANNs must guarantee
 61 the error from global loss can be back-propagated
 62 layer by layer.

63 Cyclic NN distinguishes itself from the current
 64 layer-by-layer ANNs in several aspects. 1) More
 65 flexible neuron connections. Cyclic NN greatly
 66 increases the design space of ANNs beyond the
 67 DAG structure. In Cyclic NN, the information
 68 flow is not as unidirectional as in DAG. Former
 69 neurons can also adjust based on the information
 70 encoded by the latter neurons, which largely
 71 enhances information communication within the
 72 network. The flexible connection design also
 73 makes Cyclic NN more like the biological neural
 74 system. 2) Localized training. Instead of
 75 current dominating global loss-guided BP-based
 76 training, Cyclic NN is based on localized
 77 training, *i.e.*, each neuron is optimized with
 78 its own local loss function. There is no
 79 gradient propagating between neurons. Localized
 80 training has its unique advantages. It frees
 81 the need to build DAG dependency between
 82 neurons, which is the bedrock of supporting
 83 cycles within the network. Also, each neuron
 84 is optimized independently without waiting
 85 gradients from the latter layers. 3) Computational
 86 neuron. Different from current ANNs that a
 87 neuron is considered as a d dimension to 1
 88 dimension vector mapping; the neuron within
 89 Cyclic NN is considered the computational
 90 neuron with greater computation capacity
 91 because it is the optimization unit to fit the
 92 local task, which requires more parameters. This
 93 paper uses a linear layer to parameterize each
 94 computational neuron to fit the local classification
 95 task. It is also evident by the study of biological
 96 neuron [2], which empirically proves the learning
 97 capacity of a biological neuron is much larger
 98 than a d dimension to 1 dimension vector
 99 mapping function as the neuron defined within
 100 current ANNs. We take this observation and
 101 propose the computational neuron in Cyclic NN
 102 with more capable computation to fit the local
 103 optimization task. In summary, our contributions
 104 can be summarized as follows:

- 83 • Conceptually, we compare BI and AI to investigate a fundamental yet overlooked design principle:
 84 We do not need to satisfy the DAG constraint when designing ANNs.
- 85 • Methodologically, we propose the transformative Cyclic NN, a novel ANN design paradigm that
 86 supports a much more flexible connection between neurons, which discards current directed acyclic
 87 computation graph constraints.
- 88 • We test the novel design paradigm on the most generalized case and propose Graph Over Multi-layer
 89 Perceptron, the first detailed model based on Cyclic NN.
- 90 • Experimentally, we demonstrate the advantage of the proposed Cyclic NN on widely tested datasets.
 91 At the same time, we are the first to beat the current dominating BP training using the FF training
 92 algorithm by the supported flexible network design proposed in this paper.

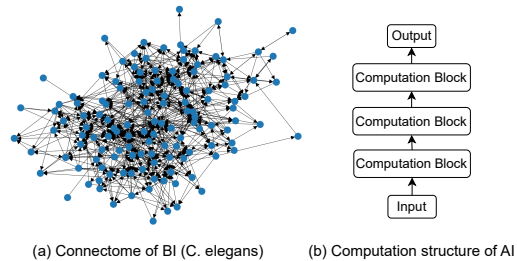


Figure 1: Neuron connection between Biology Neural Network and Artificial Neural Network

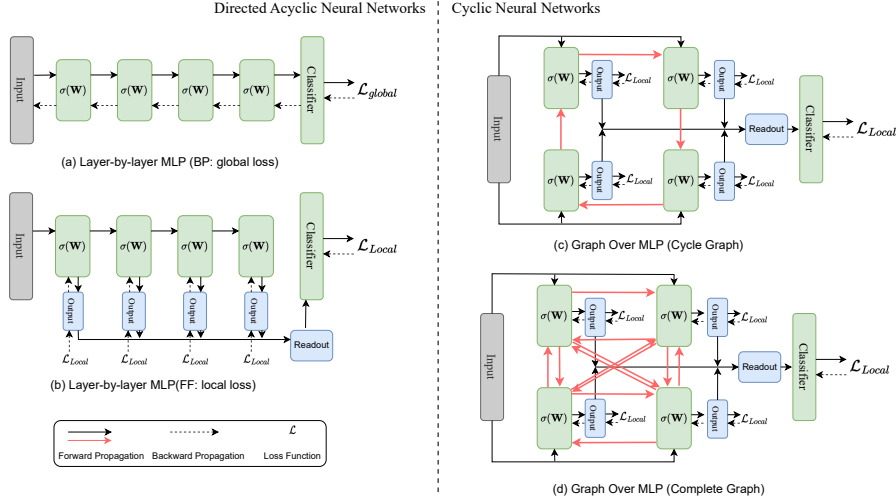


Figure 2: Comparison between different types of MLP structure.

93 2 Cyclic NN: Graph Over Multi-layer Perceptron

94 We propose the first Cyclic NN under the most generalized case, Graph Over Multi-Layer Perceptron
 95 (GOMLP), to show the design principle of Cyclic NN. As shown in Figure 2(c) and (d), GOMLP is
 96 designed by building a graph structure over the multi-layer perception to solve the classification task.

97 2.1 Input Construction

98 For the classification task, each sample is symbolized as the feature-label pair (\mathbf{h}_i, y_i) , where \mathbf{h}_i
 99 is the representation of sample i and y_i is the corresponding label. To enable the local optimization
 100 illustrated in Section 2.3, a fusion function is used to construct the input as:

$$\begin{aligned} \mathbf{h}_{\text{pos}} &= f_{\text{fusion}}(\mathbf{h}, \mathbf{y}_{\text{true}}) = \mathbf{h} \parallel \mathbf{y}_{\text{true}}, \\ \mathbf{h}_{\text{neg}} &= f_{\text{fusion}}(\mathbf{h}, \mathbf{y}_{\text{false}}) = \mathbf{h} \parallel \mathbf{y}_{\text{false}}, \\ \mathbf{h}_{\text{neu}} &= f_{\text{fusion}}(\mathbf{h}, \mathbf{y}_{\text{neutral}}) = \mathbf{h} \parallel \mathbf{y}_{\text{neu}}, \end{aligned} \quad (1)$$

101 \mathbf{h}_{pos} , \mathbf{h}_{neg} , and \mathbf{h}_{neu} are the constructed input for local optimization of different parts. f_{fusion} is a
 102 function to fuse information between feature and label, which is defined as a concat function (\parallel). \mathbf{y}_{true}
 103 is the one-hot vector of ground-true label, $\mathbf{y}_{\text{false}}$ is the one-hot vector of a randomly sampled false
 104 label. For \mathbf{y}_{neu} , we place an $\frac{1}{\text{Class Number}}$ on all the dimensions of one-hot vector to indicate $\mathbf{h}_{\text{neutral}}$ is
 105 neutral to all classes. f_{fusion} can be designed as any proper function to fuse information of the input
 106 feature and the label. In our study, we design it as a simple concat function same as [7].

107 2.2 Computation Graph

108 The computation graph \mathcal{G} contains the computational neurons \mathcal{V} and the synapses \mathcal{E} . Each computa-
 109 tional neuron $N \in \mathcal{V}$ is a local module for calculation and optimization, while synapse S defines how
 110 the information propagates between computational neurons. \mathcal{G} can be defined as a graph generator:

$$\mathcal{G} = \text{Generator}(|\mathcal{V}|, |\mathcal{E}|). \quad (2)$$

111 The above-generated graph \mathcal{G} denotes a general graph structure. Meanwhile, to justify the effective-
 112 ness of the proposed Cyclic NN, we test multiple graph generators in this paper, including the Chain
 113 graph (Figure 2(b)), Cycle graph (Figure 2(c)), Complete graph (Figure 2(d)), Watts-Strogatz (WS)
 114 graph [22] and Barabási–Albert (BA) graph [1].

115 **Neuron Update.** In GOMLP, each neuron is parameterized by a linear layer. At each propagation,
 116 neuron N is updated by $\mathbf{h}_{\text{out}} = \sigma(\mathbf{W}\tilde{\mathbf{h}}_{\text{in}})$ (we omit the notation of N in equation for simplicity)
 117 where σ is the Relu activation function [14], $\mathbf{W} \in \mathbb{R}_{d_{\text{out}}^N \times d_{\text{in}}^N}$ is N 's parameter. d_{out}^N is N 's output
 118 dimension, which is a pre-defined dimension size, and d_{in}^N is N 's input dimension, which is defined
 119 by the output of N 's pre-synapse neurons. $\tilde{\mathbf{h}}_{\text{in}}$ is the normalized input as $\tilde{\mathbf{h}}_{\text{in}} = \frac{\mathbf{h}_{\text{in}}}{\|\mathbf{h}_{\text{in}}\|_2}$, where \mathbf{h}_{in} is
 120 computational neuron N 's input.

121 **Synapse Propagation.** Each synapse $S = (N_i \rightarrow N_j)$ is a directional edge from computational
122 neuron N_i to N_j , which indicates N_i is the pre-synapse neuron of N_j and $\mathbf{h}_{out}^{N_i}$ (the output of
123 N_i) will be propagated to N_j . Assume for neuron N , we obtain a set of pre-synapse neurons
124 (N_1, N_2, \dots, N_n) based on the topology of \mathcal{G} . Then, in each propagation, N receives the output
125 of all its pre-synapse neurons along the synapses and fuse the information to form its input by
126 a concatenation function as $\mathbf{h}_{in} = \mathbf{h} \parallel \mathbf{h}_{out}^{N_1} \parallel \mathbf{h}_{out}^{N_2} \parallel \dots \parallel \mathbf{h}_{out}^{N_n}$, where \parallel is the concat function, \mathbf{h} is
127 the input representation constructed in Section 2.1. Then we can obtain $\mathbf{h}_{in, pos}$, $\mathbf{h}_{in, neg}$, $\mathbf{h}_{in, neu}$ by
128 providing \mathbf{h}_{pos} , \mathbf{h}_{neg} , \mathbf{h}_{neu} separately. As we relax the layer-by-layer restriction, the differentiation
129 between the input/hidden/output layers is also relaxed. We directly put the input \mathbf{h} to all computational
130 neurons. Thus, the input dimension size of N , $d_{in}^N = d_{\mathbf{h}} + d_{out}^{N_1} + d_{out}^{N_2} + \dots + d_{out}^{N_n}$.
131 **Readout Layer.** Readout layer collects information from all computational neurons and decides on
132 the classification. The input of the readout layer is the concat function of all computational neurons
133 as $\mathbf{h}_{in}^{readout} = f_{readout}(\mathbf{h}_{out}^{N_*}) = \parallel_{i=1}^{|\mathcal{V}|} (\mathbf{h}_{out}^{N_i})$, where \parallel is the concat function. Then, the readout layer
134 casts the representation to output dimension as $\hat{\mathbf{y}} = \text{Softmax}(\mathbf{W}_{readout} \mathbf{h}_{in}^{readout})$. where $\mathbf{W}_{readout} \in$
135 $\mathbb{R}^{\text{Class Number} \times d(\mathbf{h}_{in}^{readout})}$ is the parameter of the readout layer and $\hat{\mathbf{y}}$ is the prediction vector on classes.

136 2.3 Local Optimization

137 **Computational Neuron Optimization.** Computational neurons are optimized to differentiate the
138 positive examples from negative ones. For computational neuron N , its optimization involves $\mathbf{h}_{in, pos}$
139 and $\mathbf{h}_{in, neg}$. After the computational neuron update, we can get $\mathbf{h}_{out, pos}$ and $\mathbf{h}_{out, neg}$, respectively.
140 Then, following [7], a goodness score is calculated as $p(\mathbf{h}) = \sigma(\sum_i h_i^2 - \theta * d(\mathbf{h}))$, where $p(\mathbf{h})$ is
141 the goodness score of \mathbf{h} , $d(\mathbf{h})$ is the dimension size of \mathbf{h} , σ is the Relu activation function and θ
142 is the threshold. The binary cross-entropy loss is used to optimize each computational neuron as
143 $\mathcal{L}_N = -\frac{1}{|\mathcal{D}|} \sum_{\mathcal{D}} (\log(p(\mathbf{h}_{out, pos})) - \log(p(\mathbf{h}_{out, neg})))$, where \mathcal{D} is the dataset. The optimization of
144 computational neurons aims to increase the neuron’s output for positive samples while decreasing
145 the neurons’ output for negative samples. It enables each computational neuron its own ability to
146 differentiate positive examples from negative ones.

147 **Readout Layer Optimization.** To relieve the label leakage issue, the readout layer is only opti-
148 mized with $\mathbf{h}_{neutral}$, and we use a multi-class cross-entropy loss to optimize the readout layer as
149 $\mathcal{L}_{Readout}(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{|\mathcal{D}|} \sum_{|\mathcal{D}|} \sum_{c=1}^C y_c \log(\hat{y}_c)$, where C is the number of classes, \mathbf{y} is the one-hot
150 vector of ground-truth label and $\hat{\mathbf{y}}$ is the prediction.

151 During the inference time, we pair each test sample with the neutral label to construct \mathbf{h}_{neu} . It then
152 propagates through the GOMLP to obtain its representation on each computational neuron. Finally,
153 we predict its class with the largest logit from the output of the readout layer.

154 3 Experiments

155 We conduct experiments on
156 MNIST [11], NewsGroup [10]
157 and IMDB [13] dataset, and the
158 setup is illustrated in the supple-
159 mentary material. FF-Cycle, FF-
160 WSGraph, FF-BAGraph, and FF-
161 Complete are different versions
162 of GOMLP, where the training is
163 FF and only the graph generator
164 defined in Eq. 2 differs.

Table 1: Error rate (%) \downarrow on different datasets.

Train	Graph	MNIST	NewsGroup	IMDB
BP	Chain*	1.77 \pm 0.16	42.11 \pm 0.92	17.16 \pm 0.19
FF	Chain	1.83 \pm 0.2	43.88 \pm 0.28	18.75 \pm 0.92
BP	Chain	1.74 \pm 0.11	38.85 \pm 0.42	17.27 \pm 0.13
FF	Cycle	1.80 \pm 0.14	43.54 \pm 0.41	18.97 \pm 0.49
FF	WSGraph	1.70 \pm 0.17	38.28 \pm 0.13	17.93 \pm 0.28
FF	BAGraph	1.64 \pm 0.08	38.41 \pm 0.14	18.20 \pm 0.67
FF	Complete	1.54 \pm 0.05	38.266 \pm 0.06	17.58 \pm 0.20

165 3.1 Overall Comparison

166 The overall experiment result is shown in Table 1. We show the error rate of different methods on
167 different datasets (the lower, the better). Best performance is marked bold. From the table, we can
168 have several interesting and exciting findings:

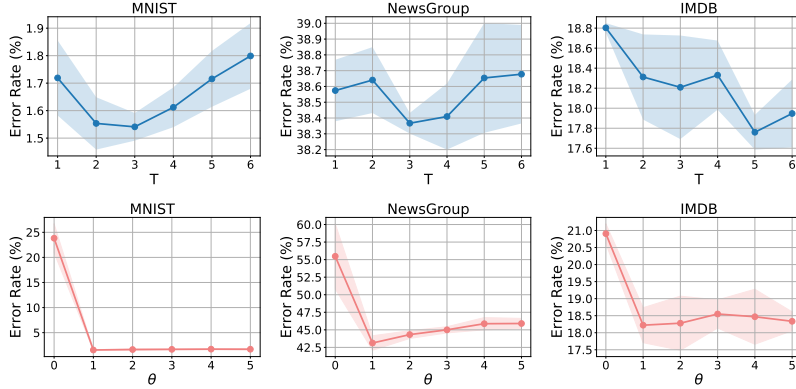


Figure 3: Parameter sensitivity of T and θ

- 169 • FF-Complete achieves the best performance on MNIST and NewsGroup datasets and comparable
170 results to the best one on the IMDB dataset. It is the first FF-trained model that outcompetes the
171 BP-trained model. It is an exciting observation of the effectiveness of the FF algorithm compared
172 with the BP algorithm.
- 173 • FF-Chain performs worse than BP-Chain* on all datasets. This observation is on par with [7], where
174 the FF lags behind the BP training algorithm when they both follow layer-by-layer organization as
175 a chain graph. However, we can surpass BP-Chain* when organizing the computational neurons
176 as a graph structure. This finding inevitably reveals the advantages of GOMLP by organizing
177 multi-layer perceptron as a flexible graph structure.
- 178 • FF-Cycle achieves similar performance with FF-Chain on three datasets. It is reasonable because
179 these two methods have only one edge difference. When we build more complex graphs (WSGraph,
180 BAGraph, Complete Graph), we can observe much better performance immediately. It shows the
181 benefits of enriching the communication between computational neurons by the GOMLP.
- 182 • BP-Chain is better than BP-Chain* in most cases. Compared with BP-Chain*, BP-Chain further
183 adds layer-wise optimization directly from the final loss. It indicates the advantageous layer-wise
184 optimization, which provides new guidelines when designing layer-by-layer neural networks.

185 In summary, the experiment results answer that we do not need to stack neural networks layer-by-layer
186 sequentially, and we can organize the neural networks as a flexible, complex graph structure like the
187 brain. More excitingly, we can outperform the current *de facto* layer-by-layer neural network design
188 paradigm with the Cyclic NN, and provide a totally new way of building ANNs.

189 3.2 Hyper-parameter Sensitivity

190 Experiment results of hyper-parameter (T and θ) sensitivity is shown in Figure 3. T controls
191 the number of propagation between computational neurons. Larger T indicates more times the
192 information is propagated. We can observe an error rate trend that first decreases and then increases
193 on all three datasets. When T is small, computational neurons can not draw sufficient lessons
194 from each other. When T is large, computational neurons are over-propagated, which leads to the
195 over-smoothing problem. θ controls the goodness threshold of each computational neuron. We can
196 observe a sharp error rate decrease when θ increases from 0 to 1, and then it gets stable with larger
197 θ . It indicates the existence of the goodness threshold matters more than the threshold value. When
198 $\theta = 0$, there is little room to optimize the computational neuron towards the negative sample, which
199 can lead to the training collapse as the computational neuron can not differentiate the negative sample.

200 4 Conclusion

201 In summary, this research introduces Cyclic NN, a novel ANN architecture inspired by the complex,
202 graph-like neural networks in biological intelligence. This transformative design diverges from
203 traditional directed acyclic ANN structures. Our findings, demonstrated through the Graph Over
204 Multi-layer Perceptron model and validated on various datasets, showed enhanced performance over
205 conventional DAG networks. This significant development paves the way for more flexible and
206 biologically realistic AI systems, representing a major shift in ANN design.

207 **References**

- 208 [1] Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Reviews*
209 *of modern physics*, 74(1):47, 2002.
- 210 [2] David Beniaguev, Idan Segev, and Michael London. Single cortical neurons as deep artificial
211 neural networks. *Neuron*, 109(17):2727–2739, 2021.
- 212 [3] Paul Brooks, Andrew Champion, and Marta Costa. Mapping of the zebrafish brain takes shape.
213 *Nature Methods*, 19(11):1345–1346, 2022.
- 214 [4] Steven J Cook, Travis A Jarrell, Christopher A Brittin, Yi Wang, Adam E Bloniarz, Maksim A
215 Yakovlev, Ken CQ Nguyen, Leo T-H Tang, Emily A Bayer, Janet S Duerr, et al. Whole-animal
216 connectomes of both caenorhabditis elegans sexes. *Nature*, 571(7763):63–71, 2019.
- 217 [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
218 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
219 2018.
- 220 [6] Donald Olding Hebb. *The organization of behavior: A neuropsychological theory*. Psychology
221 press, 2005.
- 222 [7] Geoffrey Hinton. The forward-forward algorithm: Some preliminary investigations. *arXiv*
223 *preprint arXiv:2212.13345*, 2022.
- 224 [8] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ron-
225 neberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al.
226 Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- 227 [9] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
228 *arXiv:1412.6980*, 2014.
- 229 [10] Ken Lang. Newsweeder: Learning to filter netnews. In *Machine learning proceedings 1995*,
230 pages 331–339. Elsevier, 1995.
- 231 [11] Yann LeCun, Bernhard Boser, John Denker, Donnie Henderson, Richard Howard, Wayne
232 Hubbard, and Lawrence Jackel. Handwritten digit recognition with a back-propagation network.
233 *Advances in neural information processing systems*, 2, 1989.
- 234 [12] Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series.
235 *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- 236 [13] Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher
237 Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting*
238 *of the association for computational linguistics: Human language technologies*, pages 142–150,
239 2011.
- 240 [14] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines.
241 In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages
242 807–814, 2010.
- 243 [15] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.
244 08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- 245 [16] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
246 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on*
247 *Machine Learning*, pages 8821–8831. PMLR, 2021.
- 248 [17] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by
249 back-propagating errors. *nature*, 323(6088):533–536, 1986.
- 250 [18] Alexander Shapson-Coe, Michał Januszewski, Daniel R Berger, Art Pope, Yuelong Wu, Tim
251 Blakely, Richard L Schalek, Peter H Li, Shuohong Wang, Jeremy Maitin-Shepard, et al. A
252 petavoxel fragment of human cerebral cortex reconstructed at nanoscale resolution. *Science*,
253 384(6696):eadk4858, 2024.

- 254 [19] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur
255 Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of
256 go without human knowledge. *nature*, 550(7676):354–359, 2017.
- 257 [20] Olaf Sporns and Edward T Bullmore. From connections to function: the mouse brain connec-
258 tome atlas. *Cell*, 157(4):773–775, 2014.
- 259 [21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
260 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
261 *processing systems*, 30, 2017.
- 262 [22] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *nature*,
263 393(6684):440–442, 1998.
- 264 [23] John G White, Eileen Southgate, J Nichol Thomson, Sydney Brenner, et al. The structure of the
265 nervous system of the nematode *caenorhabditis elegans*. *Philos Trans R Soc Lond B Biol Sci*,
266 314(1165):1–340, 1986.
- 267 [24] Michael Winding, Benjamin D Pedigo, Christopher L Barnes, Heather G Patsolic, Youngser
268 Park, Tom Kazimiers, Akira Fushiki, Ingrid V Andrade, Avinash Khandelwal, Javier Valdes-
269 Aleman, et al. The connectome of an insect brain. *Science*, 379(6636):eadd9330, 2023.

270 **Supplementary Materials**

271 **Datasets.** We conduct experiments
 272 on three widely studied datasets from
 273 computer vision and natural language
 274 processing domains. Data statistics
 275 are shown in Table 2. For each
 276 dataset, the training and test split fol-
 277 lows the original setting. We further
 278 extract 20% samples from the training
 279 data as validation sets to tune hyper-
 280 parameters.

Table 2: Dataset Statistics

Dataset	MNIST	NewsGroup	IMDB
# Training Samples	50,000	9,314	20,000
# Validation Samples	10,000	2,000	5,000
# Test Samples	10,000	7,532	25,000
# Dimensions	784	788	770
# Classes	10	20	2

- 281 • MNIST¹ [11]. It contains handwritten digits from 0-9, which is the most accessible and used
 282 datasets in the field of machine learning.
- 283 • NewsGroup² [10]. It is a collection of approximately 20,000 newsgroup documents, partitioned
 284 across 20 different newsgroups. This dataset is widely used for experiments in text applications of
 285 machine learning techniques, such as text classification and text clustering.
- 286 • IMDB³ [13]. It is a movie review dataset crawled from IMDB. It is the most widely studied dataset
 287 for binary sentiment classification.

288 For MNIST, we directly use its flattened pixel values as the input of all methods and replace the first
 289 10 pixels with labels as the fusion function, which is the same as [7] and leads to an input dimension
 290 of $28 * 28 = 784$. For NLP datasets (NewsGroup, IMDB), we use BERT [5] to encode the sentences
 291 into a fixed-length tensor (768) as the input. The fusion function is the concat function, which leads
 292 to an input dimension of $768 + 20 = 788$ for NewsGroup and $768 + 2 = 770$ for IMDB dataset.

293 **Baselines.** In this paper, we aim to reveal the advantages of graph-structured multi-layer perceptron.
 294 We compared GOMLP with a variant of different methods, which can be differentiated by two
 295 attributes (Training and Graph). Training indicates the training method, where BP indicates Backward
 296 Propagation [17] and FF indicates the Forward-forward algorithm [7]. The graph indicates the graph
 297 structure of computational neurons. We keep 4 computation neurons for all methods during the
 298 experiments. The special cases are further illustrated as:

- 299 • BP-Chain*: Layer-by-layer neural networks trained with BP as depicted in Figure 2(a). It is the
 300 current default way of building and training ANNs.
- 301 • FF-Chain: Layer-by-layer neural networks trained with FF as depicted in Figure 2(b) same as [7].
- 302 • BP-Chain: A modified version of BP-Chain*, where we use the structure of Figure 2(b) and trained
 303 with BP. It adds direct local supervision on each layer.

304 **Experimental Setting.** We use Adam [9] optimizer to train the model until it converges. Learning rate
 305 and weight decay are tuned within (0.1,0.01,0.001) and (0.0, 1e-2, 1e-4, 1e-6, 1e-8), respectively. The
 306 early stop technique is applied to avoid overfitting, where we stop training if there is no improvement
 307 on the validation set for continuous 10 epochs. We report the mean and variance on 20 experiments
 308 with different random seeds. All experiments are conducted on GeForce 4090 GPU.

309 **Ablation Study.** This section studies the
 310 impact of different optimization modules
 311 within GOMLP, including the computa-
 312 tional neuron optimization \mathcal{L}_N and readout
 313 layer optimization $\mathcal{L}_{Readout}$. We conduct
 314 experiments on the FF-Complete structure,
 315 and the results are summarized in Table 3.
 316 We can have the following observations: 1)
 317 The error rate increases when removing any optimization module, indicating the usefulness of each
 318 component. 2) GOMLP falls to a very large error rate (nearly random guess) when removing $\mathcal{L}_{Readout}$.

Table 3: Error rate (%) ↓ of Ablation study.

Model	MNIST	NewsGroup	IMDB
FF-Complete	1.54	38.26	18.20
- \mathcal{L}_N	2.24	47.61	22.94
- $\mathcal{L}_{Readout}$	95.58	95.55	44.26

¹<http://yann.lecun.com/exdb/mnist/>

²<http://qwone.com/~jason/20Newsgroups/>

³<https://ai.stanford.edu/~amaas/data/sentiment/>

319 It is reasonable as we depend on the readout layer to complete the final classification task. Without op-
320 timization on the readout layer, GOMLP falls into random guess even with optimized computational
321 neuron's input. 3) The error rate increases by removing \mathcal{L}_N . It shows the computational neuron's
322 optimization can provide a more informative goodness score for the readout layer to complete the
323 classification task. \mathcal{L}_N and $\mathcal{L}_{\text{Readout}}$ complement each other within GOMLP, and they collectively
324 make the best performance.