

City University London
MSc Artificial Intelligence
Project Report
2024

Extending the Biological Plausibility and Efficacy of Hebbian Convolutional Neural Networks

Author: Julian Jimenez Nimmo
Student ID: 230066319
Supervisor: Dr Esther Mondragon
Submitted: 1st October, 2024

By submitting this work, I declare that this work is entirely my own except those parts duly identified and referenced in my submission. It complies with any specified word limits and the requirements and regulations detailed in the assessment instructions and any other relevant programme and module documentation. In submitting this work I acknowledge that I have read and understood the regulations and code regarding academic misconduct, including that relating to plagiarism, as specified in the Programme Handbook. I also acknowledge that this work will be subject to a variety of checks for academic misconduct.

Signed: Julian Jimenez Nimmo

Abstract

This research explores the integration of Hebbian learning, a biologically plausible learning rule, into Convolutional Neural Networks (CNNs) for image processing tasks. Hebbian learning operates on local neural information, offering an alternative to the biologically implausible and computationally intensive backpropagation algorithm widely used in deep learning. Recent studies have combined Hebbian learning with current deep learning architectures, with distributed competition showing particular promise. Our project significantly extends the capabilities of Hebbian learning in CNNs, demonstrating performance that surpasses traditional backpropagation methods. A key finding is that Hebbian learning, when implemented with a hard competition mechanism, achieves state-of-the-art performance while exhibiting clear indications of hierarchical learning. This represents a crucial step towards more biologically realistic artificial neural networks. Furthermore, we introduce novel, biologically inspired competition mechanisms that enhance both the performance and generalisability of learnt representations. The research also includes architectural modifications inspired by both biological systems and current deep learning architectures. These changes result in networks with fewer parameters while maintaining high accuracy. Additional visualisation tool to understand Hebbian learning and its representations were also developed for interpretability and explainability.

Contents

1	Introduction and Objectives	1
1.1	Purpose, Product and Beneficiaries	1
1.2	Research Questions	2
1.3	Objectives	2
1.4	Methods	2
1.5	Structure of Report	3
2	Context	5
2.1	Neuron Model	5
2.2	Discrete-Time Weight Update and Backpropagation	6
2.2.1	Theoretical Foundations of Backpropagation Learning	7
2.2.2	Loss Functions and Error Attribution	7
2.2.3	Gradients and Learning Dynamics	7
2.2.4	Backpropagation Algorithm	8
2.2.5	Advanced Optimisation Techniques	9
2.2.6	Biological Implausibilities	9
2.3	Convolutional Neural Networks	10
2.3.1	Biologically Inspired Architecture	10
2.3.2	Theoretical and Mathematical Definitions	12
2.4	Depthwise Separable Convolutions	13
2.5	Hebbian Learning Rules	14
2.5.1	Conceptual Overview of Hebbian Learning	14
2.5.2	Variants of Hebbian Learning	15
2.6	Competition between Neurons	16
2.6.1	Winner-Takes-All (WTA) Competition	16
2.6.2	Soft Winner-Takes-All Competition	17
2.7	Anti-Hebbian and Lateral Inhibition	17
2.7.1	Anti-Hebbian Learning	18
2.7.2	Lateral Inhibition	18
2.8	Unsupervised and Hybrid Hebbian-CNN Research	19
2.8.1	Unsupervised and Hybrid Hard-WTA	19
2.8.2	Unsupervised Gradients Hard-WTA	19
2.8.3	Unsupervised Hebbian-PCA	20
2.8.4	Unsupervised SoftHebb	20
2.8.5	Hebbian Resilience	20
2.9	Literature Suggested Improvements	20
3	Methods	22
3.1	Dataset Selection and Pre-Processing	22
3.1.1	Dataset Selection	22
3.1.2	Data Loading and Processing	22
3.2	Neural Network Model Overview	24
3.2.1	Code Structure and Organisation	24
3.3	Hebbian CNN Implementation	25
3.3.1	Hebbian File Structure	25
3.3.2	Baseline: Hebbian Learning rule	25
3.3.3	Hard-WTA Research Issues	26
3.3.4	Optimised Hard-WTA	26
3.3.5	Optimised SoftHebb	27

3.3.6	BCM Learning Rule	27
3.3.7	Temporal and Statistical Thresholds	27
3.3.8	Pre-synaptic Competition	28
3.3.9	Lateral Inhibition	29
3.3.10	Depthwise Separable Convolution	30
3.3.11	Residual Block	31
3.3.12	Dale's Principle Network	31
3.3.13	Visualisation Learning	32
3.4	Experimental Setup	33
3.4.1	Configurations	33
3.4.2	Training and Testing Details	35
3.4.3	Evaluation Metrics	35
4	Results	37
4.1	Initialisation of Weights	37
4.2	Configurations Quantitative Results	37
4.3	Configurations Qualitative Results	40
5	Discussion	48
5.1	Analysis of results	48
5.2	Objectives Evaluation	50
5.3	Key Findings	50
5.4	Answer to Research Question	51
6	Evaluation, Reflection and Conclusion	53
6.1	Evaluation and Reflection of Objectives and Findings	53
6.2	Evaluation of Limitations	53
6.3	Future Work	54
6.4	Conclusion	54
A	Appendix	59
A.1	Further Results	59
A.2	Code Implementation	61
A.2.1	Data.py	61
A.2.2	Hebb.py	67
A.2.3	Hebb_depthwise.py	77
A.2.4	Hebb_abs.py	85
A.2.5	Hebb_abs_depthwise.py	92
A.2.6	Model_hebb.py	100
A.2.7	Model_depthwise	105
A.2.8	Model_residual.py	111
A.2.9	Model_BackProp.py	115
A.2.10	Experiment_hebbian.py	121
A.2.11	Experiment_hebbian_depthwise.py	127
A.2.12	Experiment_hebbian_residual.py	133
A.2.13	Experiment_bp.py	140
A.2.14	Visualizer.py	148
A.2.15	receptive_field.py	157
A.2.16	receptive_field_residual.py	161
A.2.17	params.py	166
A.3	Architectures	166

1 Introduction and Objectives

Deep neural networks have achieved remarkable, often human-level, performance across diverse tasks in Artificial Intelligence (AI), including Computer Vision (CV), Natural Language Processing (NLP), and Reinforcement Learning (RL). These advancements, both in performance and architectural complexity, have been predominantly driven by the error backpropagation algorithm developed by Rumelhart et al. (1986), despite its biological implausibility.

Backpropagation, while effective, diverges significantly from biological learning mechanisms in several key aspects. Firstly, it assumes weights are symmetrical between forward and backward connections (synapses) between neurons, whereas biological synapses are unidirectional. Secondly, backpropagation updates weights globally after applying computations across all neurons through a backward pass, contrasting with the local, feedforward learning observed in biological neurons, which depends solely on the activity of locally connected neurons.

Current deep learning architectures trained with backpropagation face several practical challenges:

- **Resource Intensity:** They require substantial computational resources, including large volumes of labelled data (big data), high energy consumption, and extended training times, contributing to significant ecological impacts. Current Large-Language-Model (LLM) training has led to a rise in emissions from the leading AI companies, as stated by NPR Kerr (2024)
- **Limited Adaptability:** Despite proficiency in specific tasks, these models often struggle to adapt to new data or generalise effectively to unknown samples.
- **Vulnerability:** They exhibit fragility against adversarial attacks, introducing critical security vulnerabilities in deployed systems.

These limitations stand in stark contrast to biological intelligence, which demonstrates robustness to novel tasks and adversarial inputs, rapid learning from limited experiences, and remarkable energy efficiency.

Hebbian learning, a biologically inspired paradigm, attempts to bridge this gap by modelling neural plasticity through local learning rules. These rules update synaptic strengths based on the correlation between pre- and post-synaptic neural activities, more closely approximating biological learning processes.

1.1 Purpose, Product and Beneficiaries

The primary purpose of this project is to analyse and implement biologically-plausible learning rules from the Hebbian learning family and biologically-inspired components within a deep learning architecture. We aim to evaluate features learnt by these models and compare their performance against traditional backpropagation methods.

This research will involve implementing state-of-the-art (SOTA) Hebbian learning algorithms to identify key factors enabling effective learning, analyse differences between various Hebbian methods and propose and test improvements to narrow the performance gap with backpropagation.

The project will yield the following products:

- Creation of versatile neural network models capable of training via backpropagation, biologically inspired learning rules, or hybrid approaches.
- Updated, computationally optimised implementation of state-of-the-art biologically plausible learning algorithms.

- Novel learning algorithms and inter-neuronal competition mechanisms for training neural network architectures.
- Comprehensive comparative analysis of Hebbian and backpropagation learning in terms of performance, efficiency, and biological plausibility.

This research will benefit several stakeholder groups, with a focus on **AI Researchers** exploring alternative learning paradigms for continuous learning and computational efficiency, **Neuroscientists and Biologists** employing artificial neural networks to model neuronal behaviour in the brain, **Industry Practitioners** requiring models capable of continuous learning from unlabelled data or models with low energy consumption and rapid training times, and **Theoretical Neuroscientists** investigating the computational principles underlying biological learning to bridge the gap between artificial and biological neural networks.

1.2 Research Questions

The primary question driving this research project is:

Can Hebbian learning algorithms achieve representational learning comparable to backpropagation while maintaining biological plausibility?

Additionally, we will explore the following sub-questions:

- How do different Hebbian learning rules affect the network's ability to form hierarchical representations?
- What role does neuronal competition play in Hebbian learning?

1.3 Objectives

To address the research questions, we have established the following objectives:

- Develop a baseline model implementing a simple Hebbian learning rule:
 - Utilise datasets and parameters consistent with existing research.
- Implement and validate state-of-the-art Hebbian approaches, ensuring replication of reported performance, learning rules, and competition modes.
- Identify limitations in current state-of-the-art methods and propose extensions, focusing on both optimisation and performance enhancements.
- Conduct comparative analyses between Hebbian and backpropagation learning across multiple metrics through visualisations of filters for interpretability and other metrics:
 - Classification accuracy.
 - Learning speed and computational efficiency.
 - Learnt representations and differences in synaptic connections

1.4 Methods

We will employ Python 3. as the primary programming language, with PyTorch serving as the deep learning framework and GitHub as our version control. Additional libraries include:

- Matplotlib and Seaborn for data visualisation
- Scikit-learn and Umap for data manipulation and evaluation metrics
- NumPy for numerical computations

- Weights and Biases (WanDB) for real-time training visualisation and analysis

To accomplish our objectives, the following methodologies were implemented:

1. CIFAR-10 Dataset with ZCA Whitening
2. Hard-WTA Hebbian CNN implementation similar to Miconi (2021), Lagani, Falchi, Gennero & Amato (2022) implementations.
3. SoftHebb CNN implementation by Journé et al. (2022)
4. Visualisation tools to understand and interpret neural learning
5. Pre-Synaptic competition
6. Temporal Competition and Statistical Competition modes
7. BCM learning rule
8. Lateral Inhibition through Surround Modulation Kernel
9. Depthwise Separable Convolution Hebbian Architecture
10. Residual Hebbian Block
11. Hebbian Network following Dale's Principle

The implementation for this project can be found on our GitHub repository:

https://github.com/J Julian-JN/Master_Thesis_Hebbian_CNN

To ensure comparable experimental setups and learning rules, we will reference and adapt code from the following publicly available projects:

- HebbDemo: <https://github.com/GabrieleLagani/hebbdemo>
- Lagani Thesis Research: <https://github.com/GabrieleLagani/HebbianLearningThesis>
- SoftHebb: <https://github.com/NeuromorphicComputing/SoftHebb>

1.5 Structure of Report

The report is structured as follows:

- Section 2: Context
 - Introduction of a neuron in biological and artificial terms
 - Explanation of backpropagation and its biological implausibility
 - Overview of a Convolutional Neural Network architecture, its biological counterpart and its mathematical foundation
 - Biological evidence for Depthwise Separable Convolutional architecture
 - Overview and description of biologically-plausible Hebbian learning rule and its variants
 - Analysis of neuronal competition methods enabling specialisation
 - Survey of recent techniques, implementations and results in the field
- Section 3: Methods
 - Neural network architecture design considerations

- Dataset processing techniques
- Detailed description of implementation decisions for baseline, state-of-the-art, and novel extensions
- Design of a comprehensive visualisation board to analyse learnt features
- Elaboration on different configurations
- Detailed explanation of the experimental setup
- Section 4: Results
 - Comprehensive presentation of results from various configurations
 - Quantitative and qualitative analysis of results
 - Visualisation of learned features, activation patterns, and performance metrics
- Section 5: Discussion
 - Evaluation of results with respect to project objectives.
 - Discussion of new findings and answer to research questions
- Section 6: Conclusion
 - Reflection on key findings and their implications for the field.
 - Evaluation of limitations and potential avenues for future work.
 - Concluding remarks on the viability of Hebbian learning as an alternative to back-propagation and overview of the project.

This structure ensures a comprehensive treatment of the research topic, from theoretical foundations to practical implementation and analysis of results.

2 Context

2.1 Neuron Model

A neuron, in both biological and artificial contexts, functions as a fundamental unit of information processing and transmission (Lagani et al. 2023). In neuroscience, it refers to a specialised cell that processes and transmits information through electrical and chemical signals. In the field of Artificial Intelligence (AI) and Artificial Neural Networks (ANNs), a neuron is an abstracted mathematical model that mimics the basic function of its biological counterpart.

The learning process in both biological systems and ANNs involves modifying the strength of connections between neurons. In biology, this is achieved through synaptic plasticity, while in ANNs it's implemented through algorithms like backpropagation.

Neurons can form deep networks by organising into layers, with each consecutive layer processing increasingly abstract representations of the input data. In deep neural networks, multiple hidden layers between the input and output layers allow for the modelling of complex, hierarchical relationships in data, mimicking the layered structure observed in biological neural systems.

A simplified biological neuron model can be modelled artificially as presented in Figure 1, where synaptic connections or weights are described by a vector w . The neuron takes as input a vector x , and produces an output $y = f(\sum x^T w)$ where f is an activation function. In biological and neuroscience terms, input values x are also called pre-synaptic activations, while the output y is termed post-synaptic activation.

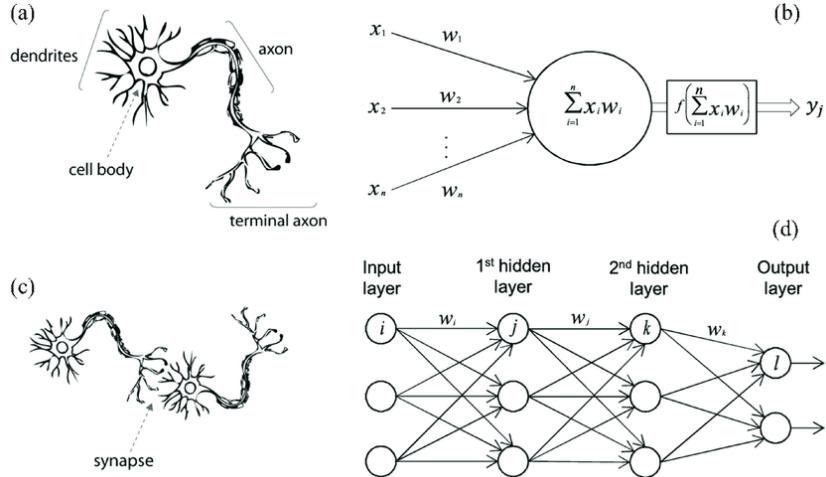


Figure 1: Comparison of biological and artificial neural models and networks (Churhe 2024)

Biological neurons operate under a multitude of physiological constraints (Pulvermüller et al. 2021), among which several principles are of importance in understanding their functionality. The principle of **locality** states that neurons are capable of direct interaction solely with their immediate neighbours via synaptic connections, constraining the scope of connectivity. A property of these synaptic connections is their capacity for **synaptic plasticity**, where the strength of synaptic connections can be modulated over time in response to activity.

Dale's Principle states that presynaptic neurons undergo either an exclusively excitatory or inhibitory influence on their postsynaptic counterparts,d. Furthermore, the propagation of information to neighbouring neurons in the network is restricted to excitatory synapses, establishing a unidirectional flow of excitatory signals in neural circuits (Shepherd & Grillner 2018).

ANNs trained with backpropagation fully respect plasticity, as it can increment or decrease the strength of a connection. The locality constraint in backpropagation is only maintained during a feedforward pass, as synaptic updates utilise global information from non-local neurons. Dale's Principle and the propagation constraint are typically not respected, as a neuron in backpropagation networks can have a mixture of excitatory and inhibitory connections to different neighbours, which has been proven to aid with representation learning.

However, new research by Cornford et al. (2020) using a backpropagation network which respects Dale's Principle and propagates information using excitatory synapses, with explicit excitatory neurons and inhibitory inter-neuron cells in a feedforward-inhibition architecture, has achieved comparable state-of-the-art performance.

2.2 Discrete-Time Weight Update and Backpropagation

In the context of ANNs, a learning rule is an algorithm that governs how the network adapts and modifies its synaptic connections in response to input data. The primary goal of a learning rule is to enable the network to improve its performance on a given task over time by learning from experience.

The discrete time learning rule is the theoretical algorithm behind ANNs. This rule operates on the principle that the network's state and weight updates occur at discrete time steps, rather than continuously across time. Current Von Neumann hardware, such as CPU and GPU, use discrete electrical signals while biological neurons operate with continuous electrical signals.

A general form for a discrete-time synaptic plasticity rule can be expressed as the change in a connection between neuron i and j , commonly referred to as synaptic strength or weight w , at time step t , based on a function f over inputs $x(t)$ and synaptic strengths $w(t)_{ij}$:

$$\Delta w(t)_{ij} = f(x(t), w(t)) \quad (1)$$

and the update to the weight as:

$$w(t+1)_{ij} = w(t)_{ij} + \Delta w(t)_{ij} \quad (2)$$

While the input $x(t)$ can be considered as the output of neuron i serving as input for neuron j , and could be expressed as $x(t)_i$, we adopt the simpler notation $x(t)$ for clarity. The time-dependent notation for weights and other terms in the equation can be simplified as we assume a discrete-time weight modification. For instance, $w(t)$ is represented as w , omitting the explicit time dependence. These notations simplification are adopted for all terms, unless not implicit.

A batch formulation of this general learning rule can be developed as follows, where multiple inputs are presented at the corresponding time step:

$$\Delta w = f(x^{(1)}, \dots, x^{(N)}, w) \quad (3)$$

All learning rules in this project work in discrete time steps given its focus on current hardware implementations, although extensions of these rules can be applied to continuous time scenarios, such as Hebbian extensions with Spike-Timing-Dependent-Plasticity (STDP), as performed in a CNN by Kheradpisheh et al. (2018). However, the backpropagation training paradigm behind the success of deep learning operate strictly in discrete time, given its theoretical and mathematical formulation.

2.2.1 Theoretical Foundations of Backpropagation Learning

The learning process in deep neural networks is driven by the efficient backpropagation algorithm (Rumelhart et al. 1986) and the gradient descent optimisation rule, which can be expressed in a simplified form as:

$$\Delta w = -\eta \nabla_w E(x, w) \quad (4)$$

where η represents the learning rate, E denotes the calculated loss function, and ∇_w symbolises the gradients associated with parameters w . This equation encapsulates how neural networks learn from data, but to understand how it modifies its weights we must define the concepts of loss functions, gradients and gradient descent, and backpropagation.

2.2.2 Loss Functions and Error Attribution

The loss function serves as a crucial metric in quantifying the discrepancy between a neural network's predictions and the actual target values or true labels, guiding the learning process (Rojas & Rojas 1996). The choice of loss function is context-dependent, varying based on the nature of the problem at hand as it decides how the task is optimised.

For regression tasks, the Mean Squared Error (MSE) is a common choice:

$$E_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (5)$$

Here, y_i represents the true value and \hat{y}_i the predicted value. The quadratic nature of MSE ensures that larger errors are penalised more heavily, driving the model to minimise significant deviations.

In classification scenarios, the Cross-Entropy loss, also known as Softmax loss, is typically employed:

$$E_{CE} = - \sum_{i=1}^C y_i \log(\hat{y}_i) \quad (6)$$

where C denotes the number of classes, y_i is the true probability of class i , and \hat{y}_i is the predicted probability. This loss function is particularly effective in multi-class problems, as it encourages the model to output probability distributions that match the true class distributions.

The loss function plays an essential role in assigning contributions to the error for each parameter in the network. By computing partial derivatives of the loss with respect to each parameter through gradients, we can determine how much each weight or parameter contributes to the overall error. This attribution of error is the key to backpropagation's success, also known as **Credit Assignment** (Bengio & Frasconi 1993), as it allows the network to understand which parameters need adjustment and to what degree.

2.2.3 Gradients and Learning Dynamics

Gradients represent the rate of change of the loss function with respect to each parameter in the network (Amari 1993). Mathematically, for each weight w_i , the gradient is defined as:

$$\frac{\partial E}{\partial w_i} = \lim_{\epsilon \rightarrow 0} \frac{E(w_i + \epsilon) - E(w_i)}{\epsilon} \quad (7)$$

This partial derivative indicates the inclination and magnitude of the steepest increase in the loss landscape for a given parameter. The negative of this gradient points in the direction of steepest descent, which is the direction we move to minimise the loss, as visualised in Figure 2.

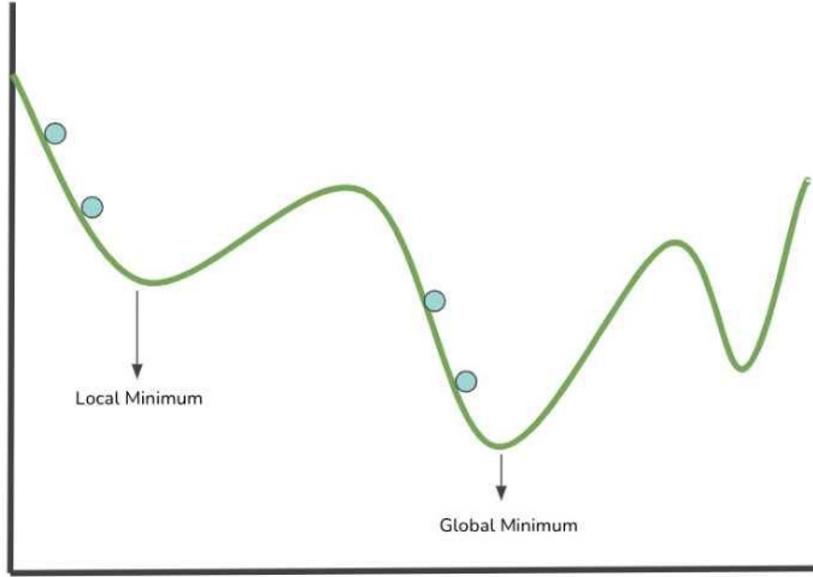


Figure 2: Simplified diagram of a loss landscape, where the aim of the backpropagation is to minimise this loss using gradients (Kucharlapati 2024)

Moving in the opposite direction of the gradient refers to the process of adjusting the network's parameters to reduce the loss. This is the foundation behind learning in current neural networks. By iteratively making small adjustments in the direction opposite to the gradient, the network gradually improves its performance, minimising the discrepancy between its predictions and the target values. This iterative process is known as **Gradient Descent**.

The computed gradients are directly used to update the weights of the network. A basic Gradient Descent weight update rule can be expressed as:

$$w_{new} = w_{old} - \eta \frac{\partial E}{\partial w} \quad (8)$$

This equation illustrates how each weight is adjusted in proportion to its contribution to the error, as determined by the gradient. The learning rate η in the gradient descent equation controls the size of these steps. A larger learning rate may lead to faster convergence but risks finding a local minimum, while a smaller learning rate ensures more stable learning but requires more iterations to reach the optimal point.

2.2.4 Backpropagation Algorithm

Backpropagation is an efficient algorithm for computing gradients in deep neural networks (Amari 1993), (Rumelhart et al. 1986). The term *propagating the error backwards* refers to the process of calculating how much each neuron in the network contributed to the overall error, starting from the output layer and moving towards the input layer.

This backward flow of information allows the network to adjust its weights and reduce the error. The process applies the chain rule of calculus to efficiently compute gradients for all parameters in the network.

For a neuron j in layer l , the gradient is computed as:

$$\frac{\partial E}{\partial w_{ji}^l} = \frac{\partial E}{\partial a_j^l} \frac{\partial a_j^l}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{ji}^l} \quad (9)$$

where a_j^l is the activation of neuron j in layer l , z_j^l is its weighted input, and w_{ji}^l is the weight of the connection from neuron i in layer $l - 1$ to neuron j in layer l .

This equation demonstrates how the error signal is decomposed and attributed to individual weights in the network. By applying this computation recursively through the layers, backpropagation efficiently calculates the gradients for all parameters, enabling the network to learn from its mistakes with gradient descent and incrementally improve its performance.

2.2.5 Advanced Optimisation Techniques

While the interplay between loss functions, gradients, and backpropagation forms the core of learning in neural networks, various algorithms have been developed to enhance convergence (Ruder 2016) and overcome challenges such as local minima. Optimizers improve upon basic gradient descent by addressing issues such as varying gradient magnitudes in different dimensions and the need for adaptive learning rates:

Stochastic Gradient Descent (SGD) updates weights using a single or small batch of samples, introducing beneficial noise that can help escape local minima. Momentum accelerates SGD by accumulating a velocity vector in directions of persistent reduction in the objective across iterations, helping to overcome areas of shallow gradients.

More sophisticated algorithms like Adam (Adaptive Moment Estimation) adapt the learning rate for each parameter, combining ideas from momentum-based methods and utilising moving averages of past gradients. These adaptive methods can lead to faster convergence and improved performance. A general form encompassing many of these optimizers can be expressed as:

$$w_{t+1} = w_t - \eta_t m_t \quad (10)$$

where η_t is the learning rate at time t , and m_t is a function of past gradients, specific to each optimizer.

2.2.6 Biological Implausibilities

Despite its effectiveness in ANNs, backpropagation faces several significant biological implausibilities that limit its direct application to neuromorphic hardware (Campbell et al. 2022). These issues highlight the divergence between current machine learning techniques and our understanding of biological learning processes (Lillicrap et al. 2020).

Gradient descent, as implemented in backpropagation, requires non-local error signals calculated from neurons that are not directly connected in the network (Song et al. 2020). This process, known as **Non-Local Plasticity**, involves propagating information backwards through the network, utilising information from neurons at different time steps and network locations. Such non-local information exchange is not observed in biological neural systems, where synaptic plasticity is primarily influenced by local interactions between pre- and post-synaptic neurons.

The reliance on non-local plasticity introduces computational inefficiencies in ANNs, as both forward and backward pass variables must be stored in memory. This requirement significantly increases the energy and resource demands of the learning process.

Another biological implausibility from backpropagation is the **Weight Transport** issue. The problem arises from backpropagation's requirement to use the transpose of the weight matrix

during the backward pass to update gradients. This requires the calculation of the transposed matrix for every neuron and its connections, further increasing computational demands. Biological neurons utilise unidirectional synapses, and the assumption of bidirectional connections with symmetric weights implied by weight transport lacks empirical support in neuroscience.

In backpropagation, the error for a neuron can only be computed after a complete forward and backward pass through the entire network. This phenomenon, known as **Update Locking**, impedes rapid learning as the network cannot process new samples until the bidirectional propagation is completed for the current sample. Update locking is biologically implausible because neurons in biological systems do not exhibit such precise coordination in modifying their synapses.

Backpropagation relies on a differentiable **Global Loss** function that provides a top-down supervision signal. This approach frames learning as a task of specialisation through loss minimisation. It can lead to issues such as overfitting, increased sensitivity to adversarial attacks, and limited generalisability. Biological learning processes are believed to involve more localised and distributed forms of error correction.

These biological implausibilities alongside its discrete timing constraint prevent the direct implementation of backpropagation-based learning rules in energy-efficient, brain-inspired neuromorphic hardware (Schuman et al. 2022). Such hardware is typically designed for local rule-based training paradigms with continuous signals that more closely align with our current understanding of biological neural plasticity. The discrepancies between backpropagation and biological learning mechanisms have spurred research into alternative learning algorithms that maintain the computational power of backpropagation while adhering more closely to biologically plausible principles.

2.3 Convolutional Neural Networks

Convolutional Neural Networks (CNN) are a deep neural network architecture designed for image or multi-dimensional data processing (LeCun et al. 1998), typically trained using the backpropagation training algorithm described in Section 2.2. A simple form of this architecture is presented in Figure 3.

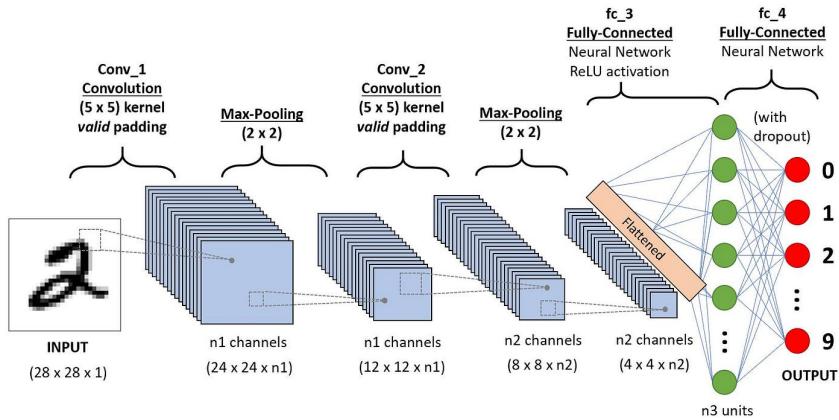


Figure 3: Overview of the CNN architecture (Mayur 2024)

2.3.1 Biologically Inspired Architecture

This convolutional architecture (Fukushima 1980) is inspired by the visual cortex in the brain, mimicking the processing of visual senses in biological systems. The primary visual cortex (V1)

found in mammals is organised in a manner which allows different neurons to specialise in detecting distinct features in small regions of the visual field known as receptive fields.

Neurons in the visual cortex are organised in a hierarchical structure where simple cells in the first regions respond to specific orientations of edges and colours, while more complex cells respond to more complex features by integrating information from multiple simple cells.

CNNs are designed to mimic this hierarchical processing, with convolutional layers simulating neurons in the visual cortex in different regions. Pooling layers aggregate the information from multiple neurons, preserving the most informative features while reducing the spatial resolution. For the hierarchical structure, CNNs use multiple layers to build up from simple patterns to more complex representations. This mimicry is illustrated in Figure 4.

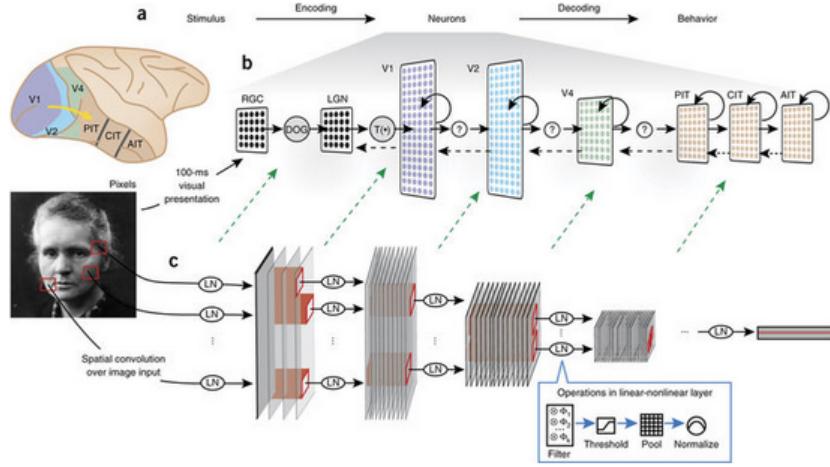


Figure 4: Comparisons between visual cortex and CNN (Lindsay 2018).

In CNNs, neurons are computational units that perform part of the convolution operation and have two key properties: each neuron is connected only to a small region of the input volume, known as its receptive field; and neurons in the same depth/channel share the same weights, implementing a key CNN concept called **weight sharing**.

Weight sharing is a fundamental concept in CNNs that ties together the ideas of neuron weights and convolutional filters. It significantly reduces the number of parameters in the network and provides translation invariance. In a given convolutional layer, multiple neurons (typically arranged in a 2D grid) share the same set of weights. This shared set of weights is equivalent to a single convolutional filter.

Each convolutional filter can be thought of as a single neuron that is replicated across the entire input space. The weights of the replicated neuron are aggregated together during an update for shared weights. The output of applying a single filter across the entire input creates a feature map. All neurons in this feature map share the same weights (filter), but operate on different spatial locations of the input.

A convolutional layer typically consists of multiple filters, each producing its own feature map. Neurons across different feature maps have different weights, allowing the network to detect various features. Mathematically, for any neuron in a particular feature map, its output can be represented as:

$$y_{i,j} = f(W * x_{i,j} + b) \quad (11)$$

where W is the shared weight matrix, $*$ denotes the convolution operation, $x_{i,j}$ is the input patch, b is a shared bias term, and f is the activation function.

2.3.2 Theoretical and Mathematical Definitions

To understand CNNs more intuitively (Wu 2017), the components of this architecture are broken down below:

Convolutional Layers: These can be thought of as feature detectors. Each convolutional filter slides over the input, looking for specific patterns (like edges, textures, or more complex shapes in deeper layers). In a mathematical notation, convolutional layers can be defined as below, where a filter/kernel K is applied to the input X :

$$Y_{i,j,c} = \sum_{m=1}^k \sum_{n=1}^k \sum_{c'=1}^{C_{\text{in}}} X_{i+m-1,j+n-1,c'} \cdot K_{m,n,c',c} \quad (12)$$

This 3D input has shape $H \times W \times C_{\text{in}}$ where H is height, W is width and C_{in} is the number of input channels. The kernel is a 4D tensor of shape $k \times k \times C_{\text{in}} \times C_{\text{out}}$ where k is the spatial dimensions of the filter, C_{in} is the number of input channels and C_{out} the number of output channels. The output Y will have size $H' \times W' \times C_{\text{out}}$ where H' and W' are determined by stride s and padding p . The result of the convolutional operation will have output shape defined by:

$$\text{Output size} = \left(\frac{H - k + 2p}{s} + 1 \right) \times \left(\frac{W - k + 2p}{s} + 1 \right) \times C_{\text{out}} \quad (13)$$

After each convolutional layer, an activation function is typically applied element-wise to introduce non-linearity. The Rectified Linear Unit (ReLU) is a common choice as it helps the network learn complex patterns by allowing for non-linear decision boundaries.:

$$f(x) = \max(0, x) \quad (14)$$

Pooling Layers: These layers perform down-sampling. They help the network focus on the most important information by reducing the spatial dimensions. The pooling layer consists of the following pooling operation:

$$Z_{i,j,c} = \max_{m,n \in \{1, \dots, k\}} Y_{i+m-1,j+n-1,c} \quad (15)$$

where $Z_{i,j,c}$ is the output of the pooling layer at position (i, j) at channel c . The output size of a pooling layer is derived from:

$$\text{Output size} = \left(\frac{H' - k}{s} + 1 \right) \times \left(\frac{W' - k}{s} + 1 \right) \times C_{\text{out}} \quad (16)$$

While max pooling is common, other pooling methods exist. **Average Pooling** takes the average value in each pooling window, while **Global Pooling** applies pooling across the entire spatial dimensions and is often used before fully connected layers.

Fully Connected Layers: These layers take the high-level features learnt by the convolutional and pooling layers and use them to make the final classification or prediction. The fully connected

layer uses flattened convolutional outputs (1D vector) f as the input:

$$o = W \cdot f + b \quad (17)$$

where w is the weight matrix and b is the bias vector.

Modern CNN architectures often incorporate additional elements. **Residual Connections** allow training for very deep networks by providing skip connections between layers (He et al. 2016), **Inception Modules** use multiple filter sizes in parallel to capture features at different scales (Szegedy et al. 2015) and **Depthwise Separable Convolutions** reduce computational cost by separating spatial and cross-channel correlations (Chollet 2017).

These advancements have led to state-of-the-art performance in various computer vision tasks, including image classification, object detection, and semantic segmentation.

2.4 Depthwise Separable Convolutions

Depthwise Separable Convolutions, developed by (Chollet 2017), is an architectural modification of a standard CNN which divides convolutional operations into two steps, seen in Figure 5: a depthwise convolution followed by a pointwise convolution. This architectural change reduces the number of parameters and computations required while maintaining similar performance and reducing overfitting.

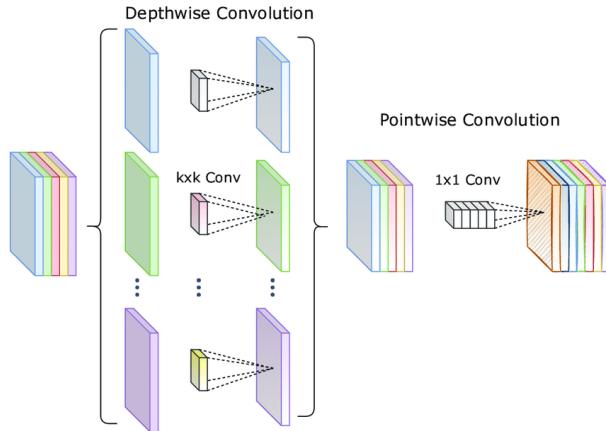


Figure 5: Depthwise Separable Convolution architecture (Sultonov et al. 2022).

Depthwise convolution step applies a single convolutional filter per input channel. If the input has M channels, it uses M filters. This step captures spatial relationships within each input channel separately and independently.

Pointwise convolution applies a convolution of dimensions 1×1 , which means it operates on a single pixel (or feature) at a time across all input channels. They take an input with dimensions (H, W, M) and produce an output with dimensions (H, W, N) . This step creates new features by computing linear combinations across channels.

Depthwise separable convolutions are considered more biologically plausible than regular convolutions. The separation of spatial (depthwise) and feature combination (pointwise) operations aligns with the specialisation observed in biological neurons. They create sparser connections between layers, which is more similar to biological neural networks than the dense connections in regular convolutions. The two-step process mimics the hierarchical nature of feature extraction in the visual cortex, where lower levels process basic spatial features, and higher levels combine these features.

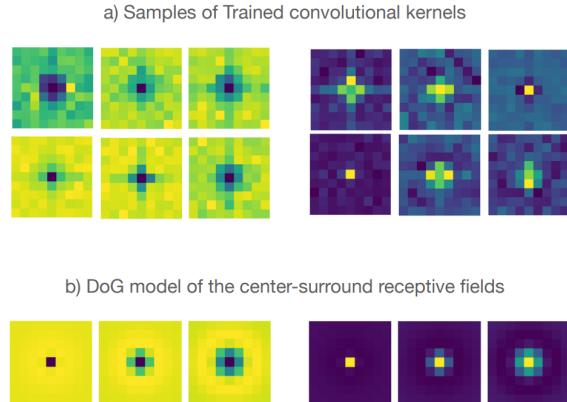


Figure 6: Centre-surround kernels which emerge in trained depthwise networks (Babaiee et al. 2024)

The depthwise step processes each channel independently, which is analogous to how different types of visual information (e.g., colour, orientation) are processed in parallel pathways in the early visual system. The pointwise step allows for the integration of information across different channels, similar to how neurons in higher cortical areas combine various types of visual information.

Recent research by Babaiee et al. (2024) has found the kernels of backpropagation trained depthwise separable convolutional networks exhibit centre surround on- and off-centre kernels, as seen in Figure 6. These kernels can be seen in biological receptive fields, marking a biological link between this convolution architecture and the visual cortex. They extended this discovery with an initialisation scheme to facilitate learning of these receptive fields.

2.5 Hebbian Learning Rules

Hebbian learning is a biologically-inspired alternative to backpropagation. Weights in synapses increase with respect to the activities in pre- and postsynaptic stimuli at a given time step (Lagani et al. 2023). This ability to modify its strength is called *synaptic plasticity*, which should be able to increase through *Long-Term Potentiation* or LTP, or decrease through *Long-Term Depression* or LTD.

2.5.1 Conceptual Overview of Hebbian Learning

At its core, Hebbian learning embodies the idea that "neurons that fire together, wire together." This principle suggests that the connection between two neurons strengthens when they activate simultaneously. In the context of artificial neural networks, this translates to adjusting the weights between neurons based on their correlated activity.

Key aspects of Hebbian learning include updates depend only on the activity of locally connected neurons, not on global network performance. No external "teacher" or error signal is required, allowing for **Unsupervised Learning**. The rule captures statistical correlations in the input data, and can provide a mechanism for storing and retrieving patterns producing **Associative Learning** (Magee & Johnston 1997).

These learning rules can be applied to discrete time Artificial Neural Networks (ANN), although biological neurons work in continuous time. The basic **Hebbian** Learning rule as described by Hebb can be expressed as:

$$\Delta w(t) = \eta y(x, w)x \quad (18)$$

where $y(x, w)$ or y is the post-synaptic activation of the neuron, which is a function of the input and the weights. It strictly follows the principle of *cells that fire together wire together*.

This learning rule applied to each individual neuron, and requires only local connections between neurons. This learning rule assumes weights are non-negative. While Equation 18 captures the essence of associative learning, it suffers from unbounded weights which can only grow in strength, leading to instability in neural network dynamics. This limitation has motivated the development of modified Hebbian rules.

2.5.2 Variants of Hebbian Learning

A modification to regulate unbounded growth is the **Grossberg Instar** learning rule (Grossberg 1976), which introduces a decay term proportional to the value of the weight and regulated by the value of the activation. Small weight modifications occur when the output is low, and vice versa.

$$\Delta w(t) = \eta y(x - w) = \eta(yx - yw) \quad (19)$$

When $\eta y(x, w) \leq 1$, this learning rule obtains the geometric interpretation of weights converging towards clusters of inputs, with each step acting proportionally to the similarity between input and weight. If a similar input is presented, the neuron will evoke a stronger response to it, and repeated presentation of inputs will allow convergence to the centroid of the input distribution.

This learning rule belongs to the family of *post-synaptic gating* rules, as the postsynaptic activity controls the extent of the update. In this form of Hebbian learning rules, the generic form can be expressed as:

$$\Delta w(t) = \eta y(x - \theta) \quad (20)$$

where the threshold θ determines how updates occur. This formulation allows for bidirectional synaptic plasticity, a key feature observed in biological neural systems. LTP is induced when the pre-synaptic activity is weaker than θ , otherwise LTD will occur.

Conversely *pre-synaptic gating* rules inverts the contribution of x and y . The threshold in these rules can be fixed on the weights, on current activations or on history of activations.

$$\Delta w(t) = \eta x(y - \theta) \quad (21)$$

The **Covariance rule** (Sejnowski & Tesauro 1989) combines both pre- and postsynaptic influences by using running averages of both pre- and post-synaptic activities, based on the biological concept of synaptic traces.

$$\Delta w(t) = \eta(y - \theta_y)(x - \theta_x) \quad (22)$$

It more accurately reflects the temporal dynamics of synaptic plasticity and incorporates the concept of metaplasticity, where the history of synaptic activity influences future plasticity.

The Bienenstock-Cooper-Munro **BCM** (Bienenstock et al. 1982) is a variation of the pre-synaptic learning rules, and introduces nonlinearity function ψ in the learning:

$$\Delta w(t) = \eta x \psi(y - \theta) \quad (23)$$

where θ_y and θ_x represent thresholds or running averages of post- and pre-synaptic activities. This rule also introduces metaplasticity and temporal dynamics to the synaptic modifications.

This threshold is not fixed, as this would lead to unstable updates as in the basic Hebbian rule. By allowing a dynamic threshold, learning can be stabilised. This threshold must grow faster than the postsynaptic activity as this activity can grow to large values. As such, non-linearity $\psi(y - \theta)$ is defined as $y(y - \theta)$, where θ is the moving average of the postsynaptic activities

A different suggestion to stabilise updates is the **Oja's** learning rule (Oja 1982), which normalises the weight vector at each update, keeping vector lengths constant but allowing directions to change. This allows the weights to perform online Principal Component Analysis (PCA) (Becker & Plumley 1996).

$$\Delta w(t) = \eta y_i (x - \sum_{j=1}^i y_j w_j) \quad (24)$$

This learning rule does add some biological implausibilities, as it requires knowledge of weights and postsynaptic activities of other neurons, which may not be directly available in biological systems.

2.6 Competition between Neurons

As neural networks are comprised of multiple interconnected neurons, it is important to have different neurons encode unique information. Competition mechanisms allow neurons to specialise and avoid convergence to similar weights as other neurons. These mechanisms play a crucial role in self-organising systems and contributes to the formation of efficient neural representations.

2.6.1 Winner-Takes-All (WTA) Competition

Winner-Takes-All or WTA is a hard competition (Rumelhart & Zipser 1985), where a singular winner is allowed to update their weights. It allows neurons to perform clustering on a set of data points and promotes a sparse representation of the data. The competition can be based on the maximum activation between neurons, or the similarity between input and weight vectors, using angular or euclidean distances.

Mathematically, the WTA mechanism can be formalised as follows:

$$i^* = \arg \max_i y_i \quad (25)$$

where i^* is the index of the winning neuron and y_i is the activation of neuron i .

This competition can be included for different spatial locations between within-layer neurons, locally connected neurons, or across the whole network. It can also be used across other dimensions, such as between neurons/filters of different channels in CNN, which is the predominant choice in Hebbian-CNN implementations Amato et al. (2019). This hard competition can also be modified to allow Top-K winner selection.

WTA competition induces a form of sparse coding, where only a small subset of neurons are active for any given input. This sparsity has been shown to have several advantages, improving

energy efficiency in neural computation, enhancing pattern separability to facilitate subsequent learning in classification tasks and Increasing storage capacity in associative memory models.

This competition can be included in any Hebbian rule, by substituting the postsynaptic activity y for a competition coefficient r , which can be expressed in different ways:

$$\Delta w(t) = \eta r(x - w) \quad (26)$$

where r can be expressed as either a binary mask denoting which neurons won the competition, or as a mask over the post-synaptic activities. This coefficient r will express the reward given to the neurons in different competition modes:

$$\text{binaryMask}_i = \begin{cases} 1 & \text{if } i = \arg \max_j y_j \\ 0 & \text{otherwise} \end{cases} \quad (27)$$

where y_j is the activation of neuron j . Now, we can express r_i in two forms:

$$r_i = \text{binaryMask}_i \quad (28)$$

$$r_i = y_i \cdot \text{binaryMask}_i \quad (29)$$

Equation 28 represents a hard Winner-Takes-All mechanism, where only the winning neuron (with highest activation) receives a non-zero r_i . Equation 29 represents an activity-masked Winner-Takes-All, where the winning neuron's r_i is set to its activation y_i , while all other neurons use $r_i = 0$.

2.6.2 Soft Winner-Takes-All Competition

A softer form of this competition was suggested (Nowlan 1989), as Hard-WTA enforces a quantised representation of the data as singular neurons respond when a specific pattern occurs in the data. Neural networks, both biological and trained on backpropagation, exhibit a distributed representation, where multiple neurons encode different patterns of the data conjointly.

Soft-WTA applies a Softmax function over the activations, allowing all neurons to learn in a distributed manner. Neurons with higher activation functions receive a higher score. A temperature parameter T can be included to control the variance of the activations.

$$\text{softmax}(y) = r = \frac{e^{y_i/T}}{\sum_{j=1}^n e^{y_j/T}} \quad (30)$$

As $T \rightarrow 0$, the Soft-WTA approaches the Hard-WTA, while as $T \rightarrow \infty$, it approaches uniform competition. The temperature parameter T provides an interpolation between these extremes. Similarly, this can be used directly instead of the postsynaptic activity or used to mask the postsynaptic activities.

2.7 Anti-Hebbian and Lateral Inhibition

Other forms of competition can be included to promote diversity between neurons, such as Anti-Hebbian Learning and Lateral Inhibition:

2.7.1 Anti-Hebbian Learning

Anti-Hebbian learning consists of the negation of the Hebbian update rules (Choe 2022), where synaptic strengths between neurons is weakened if they are activated together. The purpose of learning the opposite of Hebbian learning is to decorrelate activities between neurons. It has biological evidence in the CA1 region of the hippocampus, where certain interneurons exhibit Anti-Hebbian plasticity and contribute to the temporal coding of information

Mathematically, a basic Anti-Hebbian rule can be expressed as:

$$\Delta w_{ij} = -\eta x_i y_j \quad (31)$$

where Δw_{ij} is the change in synaptic weight from neuron j to neuron i , η is the learning rate, and x_i and y_j are the activations of neurons i and j respectively.

It induces competition by discouraging co-activation of similar neurons, while increasing diversity in the learnt representations by decorrelating neurons. This ensures each neuron learns a distinct pattern, as well as reducing the noise from data to help separate sensory data processing.

Hebbian and Anti-Hebbian learning was incorporated as a loss function for backpropagation networks in HaH by Cekic et al. (2022). Anti-Hebbian learning can be applied to all non-winners of a competition mode. SoftHebb applies this idea to all neurons except the maximally activated neuron which updates using Hebbian updates, helping decorrelate activities while using a distributed learning rule.

2.7.2 Lateral Inhibition

Lateral Inhibition is a related concept, which focuses on inhibiting the activity of neurons based on the activity of neighbouring neurons (Gabbott & Somogyi 1986). Similarly, this approach enhances the contrast and competition between neurons. It is achieved by preserving the most informative neurons active, while inhibition suppresses the activity of redundant neurons. This enhances sparsity in the representations.

Lateral inhibition is present in biological neural systems, particularly in sensory processing circuits. One of the most well-studied examples of lateral inhibition occurs in the retina. Horizontal cells provide inhibitory feedback to photoreceptors and bipolar cells.

In computational neuroscience and machine learning, lateral inhibition serves to suppress weakly activated neurons, and enhances the signal-to-noise ratio in neural representations. In visual processing, lateral inhibition can enhance edge detection and other low-level feature extraction processes. By promoting competition between neurons, lateral inhibition naturally leads to sparse neural codes, which have been shown to have computational and representational advantages.

This can be applied to neural networks through inhibitory connections which use Anti-Hebbian learning rules. WTA can be considered an extreme form of lateral inhibition, where neighbouring neurons are completely suppressed. Lateral anti-Hebbian connections between neurons in a CNN architecture were used by Pogodin et al. (2021) to showcase weight sharing can be accomplished in a locally-connected CNN. Using a sleep-like phase which uses Hebbian plasticity, it provided a local and biological solution to the weight sharing principle in neural networks.

Lateral inhibition was implemented by Hasani et al. (2019) as a convolutional operation, with a kernel replicating the surround modulation in biological visual systems. During modulation, each neuron excites nearby neighbours and inhibits far away neighbours within a range based on the level of its own activity. It was found to increase feature saliency and decrease redundancies while reducing training steps necessary for convergence.

2.8 Unsupervised and Hybrid Hebbian-CNN Research

Hebbian learning, an intrinsically unsupervised process, operates without the need for external feedback or supervisory signals to guide the learning trajectory. This learning rule facilitates online learning through continuous and incremental updates to synaptic strengths as each sample is presented.

To quantify the efficacy of this learning process, features extracted via unsupervised layers are typically passed through a single classifier layer. The MNIST (LeCun et al. 1998) and CIFAR-10 (Krizhevsky et al. 2009) datasets serve as standard benchmarks in this field, with backpropagation methodologies achieving over 99% on both datasets, thereby establishing a rigorous baseline for comparative analysis.

2.8.1 Unsupervised and Hybrid Hard-WTA

Recent investigations by Amato et al. (2019) have explored the application of Hebbian learning to Convolutional Neural Networks (CNNs) in an unsupervised context , yielding promising outcomes through diverse architectural configurations and competitive mechanisms. Lagani et al. achieved notable success, with 98.55% accuracy on MNIST and 65% on CIFAR-10 utilising a hard-WTA framework with 3 convolutional layers.

Their methodology employed cosine similarity as both the competitive criterion and activation function. Crucially, data whitening emerged as an essential pre-processing step for the enhanced performance with this hard competition. The researchers systematically evaluated architectures comprising up to six layers, observing a consistent degradation in performance with each additional layer.

Hybrid approaches integrating layers trained via backpropagation with those employing Hebbian learning have yielded intriguing insights (Amato et al. 2019). Empirical evidence suggests that the use of Hebbian learning in both initial and final layers can be executed without compromising performance. However, the integration of Hebbian learning within intermediate layers was associated with a decline in efficacy.

These results indicate that hard-WTA competition is suitable for shallow network architectures and for the fine-tuning of lower and higher layers in pre-trained networks. Notably, this approach achieves convergence in fewer epochs relative to backpropagation while reducing computational resource requirements. Nevertheless, its efficacy in training intermediate layers is limited, and performance decreases with depth, contrasting with the performance gained in backpropagation as depth increases.

2.8.2 Unsupervised Gradients Hard-WTA

Miconi (2021) integrated Hebbian learning and gradient-based networks, incorporating layer-wise loss functions equivalent to various Hebbian learning rules. This approach utilised hard-WTA competition based on post-synaptic activations alongside data whitening. Additional implementations of synaptic pruning, triangular activation functions, and adaptive thresholds to maintain homogeneous firing rates among neurons facilitated the propagation of sparse representations across layers, culminating in a 65% accuracy on CIFAR-10.

These findings set in motion the search for distributed Hebbian representations, as the localised representations in individual neurons might cause the efficiency decrements in deep networks and impede the formation of hierarchical representations.

2.8.3 Unsupervised Hebbian-PCA

Lagani, Falchi, Gennaro & Amato (2022) explored the Sanger/Oja learning rule without competition, termed HPCA (Hebbian Principal Component Analysis), to promote distributed coding capable of extracting principal components from input data. Employing the same experimental framework to their hard-WTA research but removing the requirement for data whitening, HPCA yielded comparable results and advantages. Significantly, it exhibited reduced performance degradation in deeper network configurations, maintaining accuracies around 60% in a six-layer network architecture.

2.8.4 Unsupervised SoftHebb

The SoftHebb model, introduced by Journé et al. (2022), incorporates soft-WTA competition as its foundational framework. This competitive mechanism, in conjunction with a novel plasticity learning rule, has demonstrated the capacity to minimise cross-entropy loss without explicit access to labels under specific conditions. The model employs softmax over post-synaptic activations for competitive selection and implements anti-Hebbian learning on all neurons except the maximally activated unit. The rule can be described below:

$$\Delta w_{ik}^{(SoftHebb)} = \eta \cdot y_k \cdot (x_i - u_k \cdot w_{ik}). \quad (32)$$

where w_{ik} is a synaptic weight from a pre-synaptic neuron i with activation x_i , y_k is the postsynaptic output of a neuron k , and u_k is the result of the post-synaptic softmax competition.

A key innovation is the introduction of per-neuron adaptive learning rates derived from weight vector norms, substantially accelerating learning such that it requires a single epoch. To achieve high performance, convolutional width-scaling where each successive layer increases in width by a factor of four, facilitates the formation of hierarchical representations across multiple layers. These advancements enabled the network to achieve 80% accuracy on CIFAR-10, representing the current state-of-the-art performance among unsupervised Hebbian approaches.

2.8.5 Hebbian Resilience

Comparative analysis of these Hebbian methodologies extend beyond performance metrics to encompass robustness evaluations (Gupta et al. 2022). When compared with backpropagation, Hebbian learning exhibits several notable advantages: convergence rates up to 20 times faster, performance improvements of up to 20% in limited data and computational resources scenarios, and enhanced resilience against both noise and adversarial attacks.

2.9 Literature Suggested Improvements

To improve the efficacy of these Hebbian methodologies, literature has suggested the extension of these local rules to new biologically-inspired convolutional architectures, with inter-layer backward feedback, top-down feedback connections from higher to lower layers, and recurrent connections.

Recurrent and top-down feedback has improved performance for models trained with backpropagation in object recognition tasks (Spoerer et al. 2017), reducing the error in scenarios with occlusion of objects and providing robustness against noise. A recurrent neural network trained on local reward-modulated Hebbian learning rule (Miconi 2017) showcases similar neural activity present in higher visual cortex, offering a plausible architectural model of the animal cortex.

Further improvements suggested by Lagani (2024) for bio-inspired deep learning include pre-synaptic competition, which introduces competition to the inputs x of a layer, BCM learning

rule implementation in a Hebbian-CNN, new competition modes based on statistical information of the data or temporal history of activations, and Information-theory based approaches.

In WTA competition especially in the convolutional setting, adjacent patches of an image can lead to neurons learning translated and correlated filters. Mechanisms which allow spatial decorrelation allow for richer and distributed spatial representations.

3 Methods

This chapter details the key implementation considerations that guided the development of our biologically plausible neural network model. The design philosophy encompasses several crucial aspects aimed at creating a robust, efficient, and user-friendly research tool.

- **Computational Efficiency:** A primary focus of the implementation was to optimise hardware utilisation, thereby reducing computational resource requirements and enhancing training and evaluation speeds. To achieve this, we prioritised a GPU-accelerated implementation leveraging the PyTorch framework, known for its efficient parallel processing capabilities in deep learning applications.
- **Code Readability:** To facilitate rapid development and integration within the machine learning and neuroscience communities, we placed significant emphasis on code readability. This was achieved through the use of descriptive variable naming conventions and comprehensive code documentation. Such practices ensured that researchers can easily understand, modify, and extend the codebase to suit their specific research needs.
- **Modular Architecture:** The implementation adopted a modular architecture, drawing inspiration from industry-standard practices. This approach involved dividing the system into distinct modules, each responsible for specific functionalities. The modular structure closely mirrors that of the PyTorch framework, providing an intuitive organisational system familiar to researchers experienced with PyTorch.
- **Parametric Flexibility:** A key feature of our implementation was its high degree of parametric flexibility. The system allowed for control over parameters at various levels of the model architecture. This design choice enables researchers to fine-tune the model with precision, accommodating a wide range of experimental setups.
- **Extensibility:** Recognising the rapid pace of advancements in the field, we designed the codebase with extensibility in mind. The implementation included clear documentation and adheres to clean coding practices, facilitating the seamless integration of future improvements and additional features.

3.1 Dataset Selection and Pre-Processing

3.1.1 Dataset Selection

CIFAR-10 was chosen as the competitive benchmark for Hebbian learning in our neural network implementations to analyse the differences that emerge in representations. MNIST and CIFAR-10 are two datasets which are used as a common benchmark in Hebbian research. While noticeable differences in the learnt representations and performance were encountered with the CIFAR-10 dataset, similar performances were achieved across previous methods when evaluated on MNIST, without any significant differences in the learnt representations.

The dataset is formed by 60,000 RGB images of size 32x32 pixels. There are 10 classes, which consist of animals and vehicles: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck. There are 600 images per class, and the training set is formed by 50,000 images, while the test set is formed by the remaining 10,000 images.

3.1.2 Data Loading and Processing

Data processing consisted of implementing data augmentations to the dataset. Given its highlighted importance in research, ZCA whitening was implemented as an additional data

augmentation, altering the whole dataset when the data is initialised and loaded. ZCA in our implementation was calculated in the following manner:

Given a dataset X with n samples and d features, where $X \in \mathbb{R}^{n \times d}$, we computed the mean of the data:

$$\mu = \frac{1}{n} \sum_{i=1}^n X_i \quad (33)$$

The mean from each data point was subtracted to obtain the centred data matrix \tilde{X} :

$$\tilde{X} = X - \mu \quad (34)$$

The next step was to compute the covariance matrix Σ of the centred data:

$$\Sigma = \frac{1}{n-1} \tilde{X}^T \tilde{X} \quad (35)$$

We then performed Singular Value Decomposition (SVD) on the covariance matrix Σ :

$$U, S, V^T = \text{SVD}(\Sigma) \quad (36)$$

U is a matrix of eigenvectors (or singular vectors). S is a diagonal matrix containing the eigenvalues (or singular values).

To whiten the data, we needed to compute the whitening matrix. The whitening matrix W_{ZCA} was given by:

$$W_{\text{ZCA}} = U \cdot (S + \epsilon I)^{-\frac{1}{2}} \cdot U^T \quad (37)$$

where $(S + \epsilon I)^{-\frac{1}{2}}$ is the inverse square root of the eigenvalue matrix with a small regularisation term ϵ to ensure numerical stability. ϵ was a small constant added to the eigenvalues to prevent division by zero and affected the degree of whitening applied to the dataset.

Finally, we applied the ZCA whitening transformation to the centred data:

$$X_{\text{ZCA}} = \tilde{X} \cdot W_{\text{ZCA}} \quad (38)$$

The ZCA matrix for a dataset was stored to avoid calculation of the dataset statistics as it was a costly computational process due to the iteration across the whole dataset and the large matrix multiplication involved. This implementation leveraged PyTorch's capability to perform these matrix operations efficiently on GPU hardware, creating an efficient implementation which can also be used for other datasets.

The pre-computed ZCA matrix was applied to the data samples during data loading as a data augmentation. If ZCA-whitening was not used, the data was instead normalised using the dataset's statistics, which were also pre-computed for efficiency. A random horizontal flip was included for improved generalisability, and the image was converted into a tensor data structure.

For the configurations involving Hard-WTA competition, ZCA whitening was crucial for achieving high performance. Research on Hebbian rules and Hard-WTA competition noted a significant performance drop when ZCA whitening was not implemented. To provide equal variance

while preserving the spatial structure of the data, ZCA (Zero-phase Component Analysis) is a preprocessing technique which decorrelates features from the input and can improve convergence and learning in ANNs. The implementation is described in Appendix A.2.1

3.2 Neural Network Model Overview

The selected ANN architecture for this project was the Convolutional Neural Network. This choice was motivated by several factors including its structural similarity to the visual cortex, its proven efficacy in object recognition and detection tasks, and the existing body of research applying Hebbian learning principles to CNNs.

In addition to standard CNN architectures, this study incorporated variations inspired by both biological neural networks and advanced machine learning techniques. Specifically, we evaluated the application of Hebbian learning to architectures incorporating residual blocks, as popularised by ResNet and MobileNet, and Depthwise Separable Convolutions.

3.2.1 Code Structure and Organisation

The codebase was organised into several key components to facilitate efficient training, evaluation, and adherence to the project's design goals.

Experiment Management: Files responsible for handling the training and evaluation processes were prefixed with `experiment_X.py`, where X is the configuration name. These modules handled the instantiation and modification of models and datasets based on experiment-specific parameters. They also managed the visualisation method calls to understand the learning processes across different network layers. Accuracy evaluation was conducted at each epoch during the training of the backpropagation classifier layer.

Model Architecture Definition: Files prefixed with `model_X.py` defined the overall network architecture. These modules specified individual layers and their sequential interactions based on experiment parameters. They also included methods for filter visualisation. Various architectural configurations were defined here, including implementations that mirror state-of-the-art research models, as well as those incorporating residual blocks or depthwise separable convolutions.

Layer Implementation: To achieve modularity, individual CNN layers were implemented in files prefixed with `hebb_X.py`. This design mirrors the abstraction level found in PyTorch, encapsulating the computational processes within each layer. These modules contained all calculations related to forward passes and weight updates, including different learning rules and competition mechanisms. Separate implementations were provided for standard and depthwise convolutions.

Data Management: The `data.py` module was responsible for creating training and test datasets. It also managed the application of data augmentation techniques and implemented ZCA whitening calculations.

Visualisation Utilities: The `visualiser.py` module provided a suite of functions for visualising statistical metrics, visual representations, and neural activity within the CNNs. These visualisation tools can be called at various stages of the learning process (before, during, or after training) to provide a comprehensive view of learned features and neuronal changes. The `receptive_field.py` module handled the representation of the receptive field of neurons in different layers.

3.3 Hebbian CNN Implementation

3.3.1 Hebbian File Structure

The Hebbian implementation files, prefixed with `hebb_X.py`, shared a common framework that encapsulates the core functionalities of Hebbian learning within the neural network while following our design principles. The primary components of this framework were as follows:

- The `forward()` method served as the main point for processing input through a Hebbian layer. It called the appropriate method to calculate postsynaptic activity y , called the update function to compute weight changes based on the current input and activation and returned the postsynaptic activity to be used as input for upcoming layers in the network.
- The `compute_activation()` method was responsible for calculating the postsynaptic activity y through a convolutional operation. It performed a convolution operation to generate postsynaptic activity, applying any necessary modifications to the weights based on the experimental configuration. It returned both the postsynaptic activity and the potentially modified weights.
- The `compute_update()` method was the core of the Hebbian learning process, calculating weight changes based on the chosen learning rule and competition mode. Its structure was designed for flexibility and modularity, where each combination of learning rule and competition mode was treated as a unique mode within this method. Individual functional components of these rules and competition mechanisms were implemented as separate methods. The method computed and stored the weight changes to be applied to the layer's synaptic connections.
- The `local_update()` method applied the stored weight changes to the synapses, updating the synaptic connections at a local level.

The learning rules and competition modes described in the following sections were implemented in the `compute_update()` method.

3.3.2 Baseline: Hebbian Learning rule

As a baseline, a simple Hebbian learning rule without any competition mechanisms was implemented. While the original Hebbian rule defined in Equation 18 was the simplest to implement directly, Grossberg Instar's variation defined in Equation 19 was selected instead as it lets weights grow and decrease in strength and serves as the standard rule in Hebbian-CNN research.

In our implementation, the post-synaptic activities y were calculated through a convolutional operation between the weights w and pre-synaptic activities or layer-inputs x . An additional activation function can be applied to this result, such as a ReLu or Cosine Similarity function, as performed by Lagani, Falchi, Gennaro & Amato (2022). With y calculated, the weight changes were computed using weights w , presynaptic input x and postsynaptic activity y in a method called `update_basic_hebbian()`. The change in weights Δ_w was stored in a temporary variable buffer for computational efficiency and posterior data analysis.

Weights were updated in our model using the stored weight changes. The `local_update()` method incorporated the weight changes as a gradient, which can then be combined with backpropagation-calculated gradients for a hybrid approach with an optimiser or used in a fully Hebbian approach. This was controlled by parameter α , where $\alpha = 1$, used in all our Hebbian experiments, dictated a fully Hebbian approach.

Our implementation for the calculated Δ_w for this rule is detailed in the following section, as it serves as the fundamental learning rule in all configurations except BCM learning mode. The

next step in the project was the state-of-the-art (SOTA) SoftHebb and Hard-WTA Hebbian learning rules as they presented the best performance through approaches involving different competition dynamics between neurons.

3.3.3 Hard-WTA Research Issues

The code implemented by Amato et al. (2019) was selected as our main inspiration, and we started with Hard-WTA competition. While both Amato et al. (2019) and Miconi (2021) used local Hebbian learning, Miconi (2021) always required the usage of a backpropagation optimiser as it used gradients, while Amato et al. (2019) allowed for training without an optimiser or a combination with backpropagation gradients.

The implementation by Amato et al. (2019) explicitly unfolded the image into patches. This operation is resource intensive, requiring increased memory requirements due to unfolding of tensor variables, reshaping and permuting dimensions of variables, and creation of variables based on the memory-intensive unfolded tensors. It also used non-matrix multiplication which is not as optimised for GPU hardware as matrix-based operations. It applied competition to neurons in the same spatial location but different channels, encouraging each filter to learn a different representation. These issues can be seen in the code repository for Amato et al. (2019).

3.3.4 Optimised Hard-WTA

We used the Conv2D operation provided through PyTorch as it allowed for optimised matrix operations to take place with implicit patches, working directly with the input without needing to unfold it into patches. Additional creation of variables was limited to avoid intense memory allocations. This represented our main step to optimise Hebbian learning with GPU operations, inspired by the implementation by Journé et al. (2022) and theoretical advancements behind FastHebb (Lagani, Gennaro, Fassold & Amato 2022), which used implicit patches and matrix operations.

These optimisations allowed for scalability with input size, improved memory allocation, reduced computations performed and faster speeds due to processing smaller variables. Memory was now constrained closer to the original input size, rather than significantly increasing the input memory requirements. The use of matrix-based operations allowed for efficient usage of GPU devices, increasing training speeds.

The Hard-WTA competition for selecting winner neurons in the `compute_wta_mask()` method first flattened the spatial dimensions into a single dimension and uses `argmax()` PyTorch function to locate the maximum value activation along the channel dimension for each spatial location. A binary mask was created, where winners are only preserved, and applied to the original activities.

With this masked activation acting as the postsynaptic activity y , the Grossberg Instar rule in Equation 19 was used in the `update_hardwt()` method. We applied A convolution between presynaptic input x and weights w for the yx term. For the yw term, a single value for each output channel representing the total activation of that channel across all examples in the batch and all spatial locations was calculated and applied to the weights. The result was normalised to avoid large changes in the weights.

We note yw term lost some locality in its updates as it sums over the spatial dimensions. It still served a similar purpose to the original rule, applying a modulation of the weight decay term, but applied a more global normalisation across batch and spatial dimensions at each channels in a layer. This change was implemented for computational efficiency and functionality in this convolutional context, and was also applied by Journé et al. (2022).

Cosine similarity between presynaptic activity x and synaptic connections w was also implemented

as an additional activation function during forward passes and during weight updates, as performed by Amato et al. (2019).

3.3.5 Optimised SoftHebb

We implemented SoftHebb by applying the soft competition to the postsynaptic activity using a Softmax function and an inverse temperature parameter to control the competition. We did not require additional optimisation as Journé et al. (2022) already used our optimisation techniques. The implementation can be seen in Appendix A.2.2.

The following step was to apply antihebbian learning to all neurons except the maximally activated unit. All neurons were initialised as antihebbian, and the index of the highest activation was used to convert it back to a Hebbian update. These competition mechanisms were implemented in the `compute_softwta_activations` method.

The same Grossberg Instar implementation detailed previously was applied to calculate weight changes. The only difference was the use of the result of the competition in the yx term instead of the masked activations, as detailed in Equation 32 and in the `update_softwta()` method.

3.3.6 BCM Learning Rule

Our implementation for the BCM learning rule from Equation 23 embraced the matrix-based and resource optimisation approach. An optional competition mode was applied to the postsynaptic activities for increased competition. Hard-WTA competition selection used in Section 3.3.4 was used as a competitive baseline.

The threshold θ was updated using an exponential moving average of the WTA activations. The parameter `theta_decay` determined the change in θ based on recent activity. The non-linearity $\psi(y - \theta)$ was computed as the product of the postsynaptic activity and the difference between this activity and the updated θ threshold. Each channel had its own θ threshold.

The final step in our implementation consisted of using a convolutional operation to calculate the correlation of presynaptic input x with non-linearity $\psi(y - \theta)$. The updates were normalised to maintain consistency with the other learning rules. The BCM implementation was computed in the `update_bcm()` method and can be seen in Appendix A.2.2.

3.3.7 Temporal and Statistical Thresholds

We created two distinct temporal and statistic competition modes inspired by the biological mechanisms of metaplasticity and homeostasis to expand on the concept of competition mechanisms.

Temporal Selection: Temporal competition was inspired by the concept of synaptic traces and metaplasticity, where the threshold for plasticity changes based on a history of activity called synaptic traces which stored in neuron. This phenomenon is known as metaplasticity. It promotes stability in representations learnt over time, selecting which neurons to activate when features appear consistently across time steps.

We updated synaptic traces based on recent activations in the `update_activation_history()` method, storing the post-synaptic activities as a moving average where a window parameter controlled the amount of activities stored. With the updated synaptic trace, we calculated the winners which surpass an activity threshold based on the synaptic trace. The threshold, calculated in the `compute_temporal_winners()` method, was a mean of the median activations in the synaptic trace favouring neurons which had been active consistently across time.

We then applied this competition mask back to the postsynaptic activities. An additional competition mode, such as Hard-WTA described in Section 3.3.4, could be applied for further specialisation. Finally, the Grossberg Instar implementation was used to calculate synaptic changes in the `update_temporal_competition()` method.

Statistical Selection: Adaptive threshold competition was inspired by homeostatic plasticity, which adjusts neural plasticity to maintain stable levels of activity. This avoided dominance of neurons and the vanishing activities of neurons. Typically homeostasis is implemented directly based on the firing rate and activations of neurons, as seen in Appendix A.2.2.

Our direct implementation of homeostasis maintained a running average of each neuron's activity and compares this to a target activity level, where parameters control the influence of recent activity and activities stored in the synaptic trace. Synaptic weights were gradually scaled up or down to bring each neuron's average activity closer to the target and to modify their sensitivity to inputs, as seen in the `synaptic_scaling()` method.

An indirect approach was to use the input statistics to govern the competition while maintaining stable activity across neurons. First similarities between weights and input were calculated, using a simple form of cosine similarity.

The threshold for competition calculated in the `compute_adaptive_threshold()` method was based on the similarity scores. If input patterns produced high similarity scores, the threshold increased, while low similarity scores decreased the threshold. This mechanism indirectly produced homeostasis by preventing over-activation when many high similarity scores occurred through a high threshold which restricts which neurons were updated. It prevented under-activation when low similarity scores occurred by lowering the threshold to facilitate the update of neurons.

This mechanism allowed for faster changes to unseen distributions, allowing more neurons to adapt and mimic sensory adaptation in biological systems. The threshold was computed based on the mean and standard deviation of the current input batch, allowing for rapid adaptation to changes. If the overall strength of inputs increased, the mean similarity score would increase, raising the threshold and promoting specialisation. If the variability of inputs changed, the standard deviation similarity score reflected these changes and affected the threshold to allow more neurons to adapt.

Winner neurons were selected with this threshold and optionally passed to a further competition mechanism. This additional competition mechanisms in both modes promoted efficient specialisation and sparse connectivity. The Grossberg Instar rule was again used to update the synaptic connections in the `update_adaptive_threshold()` method.

3.3.8 Pre-synaptic Competition

We designed a form of pre-synaptic competition, following the suggestions by Lagani (2024). The pre-synaptic input x was combined with parallel synaptic coupling $m_{i=1} = 1/w_1, m_{i=2}, \dots, m_{i=N}$ where i is the coupling to each neuron and N the total number of neurons. The concept of synaptic coupling is similar to weights in neural networks whereas in biological neural networks they refer to the strength of connections between neurons. However, unlike synapses they directly influenced each other, leading to competitive dynamics in neural signalling and learning.

In the context of CNN, synaptic couplings were represented as the kernel weights connecting input channels to output channels. Based on the recommendations from Lagani (2024), the weights were converted into synaptic couplings m by first inverting their value, to provide normalisation, competition and simulate electrical admittances:

$$m = \frac{1}{w} \quad (39)$$

A competition mechanism was applied to these couplings in the `compute_activations()` method and the result was used for convolutions and update calculations. The different competition modes which can be selected to introduce presynaptic competition were implemented in the `compute_presynaptic_competition()` method and are listed as:

1. Linear Competition, defined as:

$$w_{eff} = \frac{m}{\sum_i m_i + \epsilon} \quad (40)$$

where $m = \frac{1}{|w|+\epsilon}$, w is the original weight, and ϵ is a small constant (1e-6) to prevent division by zero. It normalises the weights to values between 0-1.

2. Softmax Competition, defined as:

$$w_{eff} = \text{softmax}(m) = \frac{e^{m_i}}{\sum_j e^{m_j}} \quad (41)$$

where $m = \frac{1}{|w|+\epsilon}$ as before. This creates a more pronounced competition where stronger connections are emphasised, and the weights summed together add up to 1.

3. L2 Norm Competition, defined as:

$$w_{eff} = \frac{m}{\sqrt{\sum_i m_i^2}} \quad (42)$$

where $m = \frac{1}{|w|+\epsilon}$. This ensures that the sum of squared effective weights equals 1.

In all cases, the competition was applied along the output channel dimension. This means that for each input channel and same spatial location in the kernel, the weights connecting to all output channels competed with each other. This promotes diversity among output channels, encouraging them to specialise in detecting different features from the same input locations.

Additional competition across different spatial and channel dimensions of the architecture were also implemented. Applying competition across spatial dimensions within each input-output channel pair encouraged each input-output channel pair to focus on specific spatial patterns. This was implemented in `compute_presynaptic_competition_spatial()` method.

Implementing competition across input channels for each output channel and spatial location encouraged each output neuron to be selective about which input features it responds to. This was implemented in `compute_presynaptic_competition_input()` method.

A more global competition across all dimensions simultaneously created a more intense competition where every weight competed with all others leading to very sparse but highly specialised connections. This was implemented in `compute_presynaptic_competition_global()` method.

3.3.9 Lateral Inhibition

We created a fixed kernel created through a difference of Gaussians (DoG) function to simulate the effect of lateral inhibition and surround modulation found in early layers of the visual cortex and the excitatory/inhibitory impact on neighbouring neurons. It was implemented in the `create_sm_kernel()` method.

This kernel models how neurons in the visual cortex respond to stimuli in their receptive field and surrounding areas. The centre (excitatory) region represents the classical receptive field,

while the surround (inhibitory) region models lateral inhibition from neighbouring neurons as seen in Figure 7. It is mathematically defined as:

$$K_{SM}(x, y) = \frac{1}{K_{center}} \left(\frac{G_e(x, y)}{2\pi\sigma_e^2} - \frac{G_i(x, y)}{2\pi\sigma_i^2} \right)$$

where:

$$\begin{aligned} G_e(x, y) &= \exp \left(-\frac{x^2 + y^2}{2\sigma_e^2} \right) \\ G_i(x, y) &= \exp \left(-\frac{x^2 + y^2}{2\sigma_i^2} \right) \\ K_{center} &= \left(\frac{G_e(x, y)}{2\pi\sigma_e^2} - \frac{G_i(x, y)}{2\pi\sigma_i^2} \right) \Big|_{x=0, y=0} \end{aligned}$$

where $K_{SM}(x, y)$ represents the surround modulation kernel at position (x, y) . $G_e(x, y)$ and $G_i(x, y)$ are the excitatory and inhibitory Gaussian functions, respectively. σ_e and σ_i are the standard deviations for the excitatory and inhibitory Gaussians. K_{center} is the value of the unnormalised kernel at the centre $(0, 0)$. The kernel was computed as the difference between the normalised excitatory and inhibitory Gaussians. This kernel strengthened synapses in the immediate neighbourhood of a neuron, and weakened synapses further from this neighbourhood.

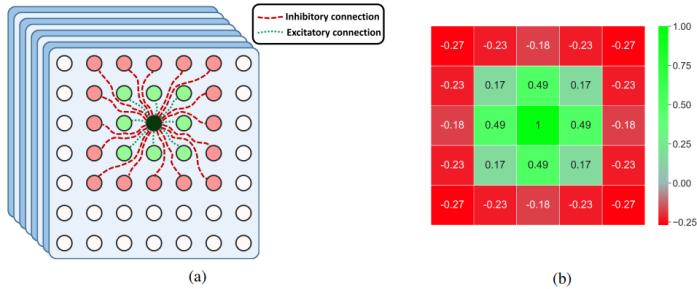


Figure 7: Surround modulation kernel, as presented in Hasani et al. (2019)

This fixed kernel was applied to the postsynaptic activities through a convolutional operation, enhancing contrast and edge detection while suppressing uniform backgrounds.

3.3.10 Depthwise Separable Convolution

We developed a novel Hebbian implementation of the biologically plausible Depthwise Separable Convolution to investigate the potential of biological rules which can generate centre-surround kernels observed in neurons within the visual cortex.

The depthwise convolution shared similarities with the previously described convolutional approach, with modifications in the application of updates and competition mechanisms to ensure filter independence. Given the one-to-one correspondence between input and output channels, we extracted this relationship from the convolutional operation output yx .

By extracting the diagonal elements along the input and output channel dimensions, we selected only the elements where the input channel index matches the output channel index, preserving the independence criteria. The yw term was also adjusted to preserve independence. In the

depthwise CNN, yw was computed separately for each channel, ensuring the update for each filter depended only on its own activations.

The competition modes were modified to shift the dimension of competition from neurons at the same spatial location across channels to neurons within the same channel across spatial locations. Additionally, to maintain independence, the number of output channels was required to match the number of input channels. These implementation modifications can be seen in Appendix A.2.3.

The pointwise convolution, implemented as a 1×1 convolution, was fully compatible with the previously established learning rules and competition mechanisms. This compatibility stems from the interaction between channels at identical spatial locations and lack of filter independence. The implementation required only a kernel size of 1.

This variation of the Hebbian-CNN was located in `hebb_depthwise.py`, and the model using this architecture was defined in `model_depthwise.py`, as shown in Appendix A.2.

3.3.11 Residual Block

We implemented additional feedforward connections linking layers at different sequential orders to enhance the biological plausibility of the model. This approach was inspired by the concept of skip connections and residual blocks found in advanced deep CNN architectures such as ResNet and MobileNet.

Convolutional blocks were constructed with a bottleneck structure utilising Depthwise Separable Convolutions for feature extraction. A skip connection was incorporated to link the input of the block to the output of the bottleneck. This block design is related to the approach employed in MobileNet, using depthwise separable convolutions for parameter reduction while integrating the advantages of residual blocks from ResNet to facilitate training of deeper networks.

The block employed an inverted bottleneck structure, built by three convolutional layers. The first layer was a pointwise convolution that expanded the number of channels, facilitating the learning of rich channel features. The second layer was a depthwise convolution which learnt the spatial information from the expanded feature space. The third was another pointwise convolution that projected the features back to a desired output channel count.

The skip connection was a pointwise convolution designed to match the input size to the output size of the third layer of the bottleneck. The output of this skip connection was added to the output of the bottleneck structure, enabling more efficient learning in deeper architectures. The block module was defined in the `HebbianResidualBlock()` class in (`model_residual.py`), as seen in Appendix A.2.

3.3.12 Dale's Principle Network

We modified the weight initialisation and update processes to investigate the learning dynamics of a purely excitatory-driven network that follows Dale's Principle, in comparison to mixed biologically implausible excitatory-inhibitory weights.

All synaptic weights were initialised to positive values with an absolute value function to ensure compliance with this principle. In any operation involving weights as a term, the absolute value function was applied to maintain this constraint. The change in weights can still produce both LTP and LTD, as both directions are essential for plasticity.

If a weight update resulted in a potentially inhibitory weight, two approaches were considered, the application of an absolute value function to ensure exclusive excitatory impact or a ReLU function to allow only excitatory connections to exist. This approach may lead to neurons

with activity levels close or equal to zero remaining near this level of activity during learning, potentially creating sparse connections at the expense of facilitating the formation of distributed representations.

To address this phenomenon, two alternative mechanisms were proposed, a homeostasis mechanism as detailed in Section 3.3.7 to maintain all neurons active or a probabilistic metaplasticity mechanism that pruned synapses near a local minimum and generated new synapses based on local activity. These alternatives could potentially mitigate the sparsity issue while maintaining the biological plausibility of the network.

This Hebbian variation and its mechanisms which respects Dale’s Principle were created in the `hebb_abs.py` and `hebb_abs_depthwise.py` files, further showcased in Appendix A.2.5.

3.3.13 Visualisation Learning

We developed a comprehensive suite of visualisation tools to enhance our understanding of Hebbian learning and neuronal feature detection in CNNs. These can be employed at any time during learning or for posterior analysis. These tools, detailed in Appendix A.2, provided innovative insights into the learning dynamics of Hebbian-CNN across various scales, from individual synapses to entire layers. All plots are stored using WanDB.

Descriptive statistical information was compiled for all neurons within each layer, encompassing both the pre-update weights and the weight changes. This analysis included measures such as mean, median, standard deviation, and the ratio of positive to negative weights or weight updates. The data was presented for a selected subset of neurons as well as for the entire layer, offering a multi-scale perspective on neuronal behaviour.

Violin plots were employed to depict the distribution of weight changes for a selected subset of neurons, providing a nuanced view of synaptic plasticity. Additionally, weight value distributions across entire layers were visualised to understand neuronal connectivity patterns and synaptic sparsity in different model configurations. These distributions also offered insights into the stability of weight configurations over time.

An innovative approach to visualising CNN filters was developed, focusing on weight changes per spatial location within each filter. This technique involved plotting the mean weight change across channels, resulting in a 2D representation that reveals emergent patterns and localised regions of LTP and LTD. For 1x1 convolutions, which lack spatial structure, bar plots were utilised to represent weight changes, with colour coding to indicate magnitude and direction of change.

Complementing this, a global LTP/LTD visualisation was created to illustrate the dominant learning direction for a selected number of neurons, providing a global view of synaptic plasticity across the network. Additionally, the features from the final convolutional layer were reduced to a 2D space through dimensionality reduction techniques such as UMAP to allow a visual representation of the different classes and their clusters. This enabled us to determine whether the model effectively groups inputs of the same class together while distinguishing them from inputs of different classes. These described tools were defined in the `visualiser.py` file.

The filters themselves were also directly visualised to provide insight into the features they represent. For the initial convolutional layer, full RGB visualisations offered an intuitive representation of low-level features. In deeper layers, where direct visualisation is less interpretable, the mean across channels was used to create 2D representations that only capture dominant patterns. It was defined in the `visualize_filters()` method in all `model_X.py` files.

To further comprehend neuronal selectivity, receptive fields were visualised using a Projected Gradient Ascent (PGA) method. In biological terms, the receptive field corresponds to the

region of sensory space to which a neuron responds. Receptive fields in CNNs are hierarchical and grow larger in deeper layers due to the cumulative effect of convolutional operations. This hierarchical structure should allow the network to capture increasingly complex and abstract features as information propagates through the layers.

In the first layer, receptive fields may correspond to simple edge detectors or spots of colour which match with the direct filter visualisations at the first layer, while in deeper layers they may represent complex object shapes or even entire objects which were not captured by the direct filter visualisations of deeper layers which only focus on patterns.

The PGA technique optimised an input image initialised with random noise and sized to match the neuron’s receptive field, to maximise the neuron’s activation. We forward propagated this image through the network to the neuron of interest and computed the gradient of the neuron’s activation with respect to the input image. Next, we updated the input image to maximise the neuron’s activation, subject to regularisation constraints to keep the image interpretable.

This process revealed the optimal stimuli for neurons at various levels of the architecture, demonstrating the hierarchical representations formed across layers. We implemented a custom implementation of PGA in the `receptive_fields.py` file, which incorporated L2 regularisation and gradient optimisation to aid with finding the optimal and explainable stimuli which activates a neuron.

3.4 Experimental Setup

3.4.1 Configurations

Tables 1 and 2 present a detailed overview of various configurations designed to evaluate competition dynamics and Hebbian learning rules in neural networks. We designed and ran each configuration to systematically explore different combinations of architectural designs, competition mechanisms, and learning rules, each designed to investigate specific aspects of neural network behaviour and performance.

The configurations primarily utilised two main architectural approaches: the SoftHebb architecture, based on the work of Journé et al. (2022), which incorporated padding and has larger number of filters per layer, and the Lagani, Falchi, Gennaro & Amato (2022) padding-free and smaller architecture with 3 or 4 convolutional layers. These architectural choices provided a foundation for comparing different network structures and their impact on learning and competition dynamics. More details for these architectures can be found in Appendix A.3.

A key focus of this experimental setup was the examination of various competition mechanisms. The configurations incorporated a range of approaches, including Soft Winner-Take-All (SoftWTA), Hard Winner-Take-All (HardWTA), presynaptic competition, temporal selection, statistical selection, and surround lateral inhibition. These mechanisms were combined, as seen in Table 1 where configurations like `SoftHebb-Presynaptic/HardWTA/Cos-Instar` integrated both HardWTA and presynaptic competition while using a Cosine Similarity Activation. This approach allowed for an understanding of how different competition strategies interact and influence network behaviour.

The learning rules employed in these configurations are predominantly the Grossberg Instar rule and the Bienenstock-Cooper-Munro (BCM) rule. The inclusion of a configuration using backpropagation with `SoftHebb-Backpropagation` provided a benchmark for comparison with traditional gradient-based learning methods. This diversity in learning rules enabled a comparative analysis of their effectiveness in similar scenarios.

Table 1: Hebbian Learning Configurations and Features

Configuration	Architecture	Competition Mechanism	Learning Rule	Additional Features
SoftHebb-SoftWTA-Instar	SoftHebb	Soft WTA	GI	Hebb. & anti-Hebb. learning
Lagani-HardWTA/Cos-Instar	Lagani	Hard WTA	GI	Cos. sim. activation
Lagani_Deep-HardWTA/Cos-Instar	Lagani_Deep	Hard WTA	GI	Cos. sim. activation
SoftHebb-HardWTA/Cos-Instar	SoftHebb	Hard WTA	GI	Cos. sim. activation
SoftHebb-HardWTA/Cos-BCM	SoftHebb	Hard WTA	BCM	Cos. sim. activation
SoftHebb-None-Instar	SoftHebb	None	GI	-
SoftHebb-Pre/HardWTA/Cos-Instar	SoftHebb	Hard WTA+Pre	GI	Cos. sim. activation
SoftHebb-Temp/HardWTA/Cos-Instar	SoftHebb	Hard WTA+Temp	GI	Cos. sim. activation
SoftHebb-Stats/HardWTA/Cos-Instar	SoftHebb	Hard WTA + Stats	GI	Cos. sim. activation
SoftHebb-Sur/HardWTA/Cos-Instar	SoftHebb	Hard WTA+Sur	GI	Cos. sim., Sur lat. inhib.
Depthwise_SoftHebb-Sur/HardWTA/Cos-Instar	SoftHebb	Hard WTA+Sur	GI	Cos. sim., Depthwise sep. conv.
Residual_SoftHebb-Sur/HardWTA/Cos-Instar	SoftHebb	Hard WTA+Sur	GI	Cos. sim., Res. connections
Dale_SoftHebb-Sur/HardWTA/Cos-Instar	SoftHebb	Hard WTA+Sur	GI	Cos. sim., Dale's Principle
SoftHebb-Backpropagation	SoftHebb	N/A	BP	-

Note: WTA: Winner-Take-All, GI: Grossberg Instar, BCM: Bienenstock-Cooper-Munro, BP: Backpropagation, Pre: Presynaptic, Temp: Temporal, Stats: Statistical, Sur: Surround, Cos. sim.: Cosine similarity, sep. conv.: separable convolutions, Res.: Residual, lat. inhib.: lateral inhibition, Hebb.: Hebbian

Table 2: Purposes of Hebbian Learning Configurations

Configuration	Purpose
SoftHebb-SoftWTA-Instar	Evaluate SoftHebb research
Lagani-HardWTA/Cos-Instar	Evaluate Lagani Hard-WTA research
Lagani_Deep-HardWTA/Cos-Instar	Evaluate Lagani Hard-WTA 4-layer research
SoftHebb-HardWTA/Cos-Instar	Evaluate Hard-WTA in network with padding
SoftHebb-HardWTA/Cos-BCM	Compare BCM to Grossberg-Instar rule
SoftHebb-None-Instar	Evaluate basic learning without competition
SoftHebb-Pre/HardWTA/Cos-Instar	Assess impact of presynaptic learning
SoftHebb-Temp/HardWTA/Cos-Instar	Evaluate temporal competition with/without WTA
SoftHebb-Stats/HardWTA/Cos-Instar	Evaluate statistical competition with/without WTA
SoftHebb-Sur/HardWTA/Cos-Instar	Assess effect of surround modulation kernel
Depthwise_SoftHebb-Sur/HardWTA/Cos-Instar	Evaluate depthwise equivalent of hard-WTA research
Residual_SoftHebb-Sur/HardWTA/Cos-Instar	Evaluate residual equivalent of hard-WTA research
Dale_SoftHebb-Sur/HardWTA/Cos-Instar	Evaluate Dale's Principle in hard-WTA research
SoftHebb-Backpropagation	Evaluate backpropagation in SoftHebb models

Configurations following `SoftHebb-Sur/HardWTA/Cos-Instar` and onward incorporated additional features to enhance network performance or investigate specific neural phenomena. The cosine similarity activation, present in most configurations, was used as advantages were noted by Lagani, Falchi, Gennaro & Amato (2022). Some configurations explored architectural modifications, such as depthwise separable convolutions (`Depthwise_SoftHebb-Sur/HardWTA/Cos-Instar`) and residual connections (`Residual_SoftHebb-Surr/HardWTA/Cos-Instar`), implemented after the first layer while maintaining input and output channel equivalence to standard CNN layers. The `Dale_SoftHebb-Surr/HardWTA/Cos-Instar` configuration introduced biological constraints by adhering to Dale’s Principle, offering insights into the impact of neuron-type restrictions on network performance.

This systematic experimental setup provided a foundation for understanding how various components interact and contribute to network performance, potentially leading to the development of more efficient and biologically plausible neural network models.

3.4.2 Training and Testing Details

The training stage for Hebbian configurations followed a sequential approach comprising three distinct phases: feature learning, feature extraction and analysis, and classifier training. In the initial feature learning phase, only the convolutional layers were trained for a single epoch to develop feature representations. For all configurations, a batch size of 64 datapoints was used to train the models.

For Hard-WTA configurations, a learning rate of 0.1 was applied to it, while for Soft-WTA the custom learning rate schedule detailed in Journé et al. (2022) and in Appendix A.2 was used instead. It is noted that configurations employing `SoftWTA` and those using `HardWTA` or no competition mechanism (`None`) utilised slightly different parameters. These variations included the application of ZCA-Whitening. The specific details of these training configurations are described in Appendix A.2 and A.3.

Throughout the Hebbian training process, we employed our custom-visualisation tools to store and visualise layer data. This method allowed for a comprehensive analysis of the network’s developing representational abilities. Once the feature learning phase was completed, the convolutional layer weights were frozen to retain the learned features. At this stage, we gathered more information to obtain a deeper understanding of the representations learnt at different layers.

The final phase of our training procedure involved the classifier head, which was trained on top of the frozen convolutional layers. This classifier was trained through backpropagation for five epochs through an Adam optimiser with a learning rate of 0.001. To ensure rigorous evaluation, we assessed both the training and test sets at each epoch, storing the results for comparative analysis between different configurations.

The whole backpropagation model was trained through an Adam optimiser with a learning rate of 0.01 and no data-whitening. All other hyperparameters remained the same as in Hebbian configurations.

3.4.3 Evaluation Metrics

We employed the following evaluation metrics to facilitate a quantitative comparison between Hebbian learning configurations: accuracy, precision, recall, and F1 score.

In classification tasks, True Positives (TP) refer to instances where the model correctly identifies a positive class accurately. False Negatives (FN) occur when the model incorrectly predicts a positive class as negative. On the other hand, True Negatives (TN) are cases where the model

correctly identifies negative instances correctly. Lastly, False Positives (FP) represent cases where the model incorrectly predicts a negative instance as positive. These four terms are essential for calculating performance metrics .

Accuracy serves as a fundamental measure of the model's overall correctness across all classes. While accuracy provides a good indication of performance, it may yield misleading results for imbalanced datasets. Equation 43 expresses this mathematically:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (43)$$

Precision quantifies the model's ability to avoid labelling negative samples as positive. This metric is particularly crucial in applications where false positives incur significant costs or risks. Equation 44 describes this mathematically:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (44)$$

Recall reflects the model's capacity to identify all positive instances within the dataset. This metric is important in scenarios where false negatives are particularly undesirable or costly. Equation 43 expresses this mathematically:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (45)$$

The **F1 score** provides a balanced measure of the model's performance. By combining precision and recall into a single metric, the F1 score offers a holistic view of the model's performance. Equation 46 represents this mathematically:

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (46)$$

4 Results

4.1 Initialisation of Weights

The weights for each synaptic connection were initialised using the same weight distribution, as both Soft-WTA and Hard-WTA functioned well with this scheme. A normal random distribution with a large radius range was used, as low radius range did not produce learning in the SoftHebb model, detailed in Journé et al. (2022).

$$\text{weight_range} = \frac{25}{\sqrt{C_{in} \cdot K_h \cdot K_w}}$$

$$W = \text{weight_range} \cdot X$$

$$X \sim \mathcal{N}(0, 1)$$

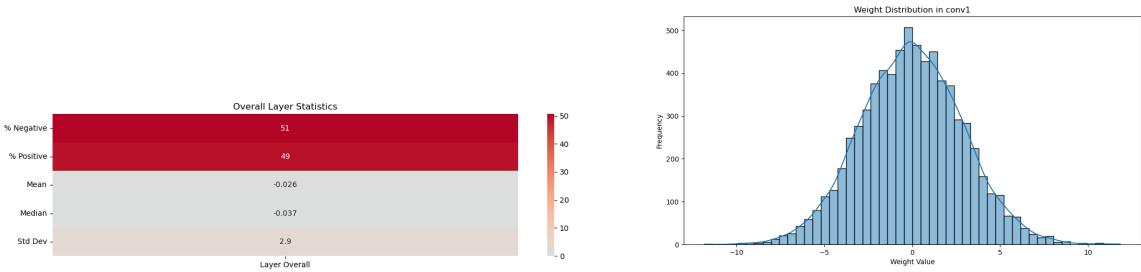


Figure 8: Results of the initial statistics and normal distribution of weights in the first layer. Deeper layers also present similar initial statistics and distribution.

The obtained standard normal distribution was centred at 0, so approximately half of the generated values were negative and half positive, before being scaled by weight_range to create a large range of weight values. Statistical results for the initial weights and the initial distribution for the first convolutional layer can be visualised in Figure 8, visualised using the tools we developed for the project.

4.2 Configurations Quantitative Results

Each experimental configuration detailed in Table 1 was evaluated with a fixed random seed for equal experimental conditions, and the mean Accuracy and F1-score metrics for the test set are presented in Table 3.

SoftHebb-SoftWTA-Instar and **SoftHebb-Surround/HardWTA/Cos-Instar** demonstrated superior performance, with 73.6% and 73% accuracy respectively. With these results, we achieved a new milestone and state-of-the-art performance for Hard-WTA competition, significantly advancing the field of Hebbian learning in neural networks. State-of-the-art (SOTA) includes all results surpassing 70% accuracy.

In addition, the new depthwise architecture in **Depthwise_SoftHebb-Surround/HardWTA/Cos-Instar** achieved competitive results with notably less parameters, reducing parameters by 6.6x from 5,893,162 to 883,978 while achieving 69% accuracy. Other configurations including **SoftHebb-HardWTA/Cos-BCM**, **SoftHebb-Temp/HardWTA/Cos-Instar**, **Residual_SoftHebb-Surround/HardWTA/Cos-Instar** and **SoftHebb-HardWTA/Cos-Instar** surpassed 70% accuracy through a different learning rule, competition mechanism, architectural change or dependence on a sole competition mechanism respectively.

Table 3: Performance Metrics for Experimental Configurations on same Fixed Seed

Configuration	Accuracy (%)	F1 Score
SoftHebb-SoftWTA-Instar	73.6	0.738
Lagani-HardWTA/Cos-Instar	57.7	0.571
Lagani_Deep-HardWTA/Cos-Instar	53.2	0.527
SoftHebb-HardWTA/Cos-Instar	72.1	0.72
SoftHebb-HardWTA/Cos-BCM	71.1	0.711
SoftHebb-None-Instar	19.9	0.6
SoftHebb-Presynaptic/HardWTA/Cos-Instar	66.5	0.664
SoftHebb-Temp/HardWTA/Cos-Instar	71.5	0.715
SoftHebb-Stats/HardWTA/Cos-Instar	59.8	0.582
SoftHebb-Surround/HardWTA/Cos-Instar	73	0.73
Depthwise_SoftHebb-Surround/HardWTA/Cos-Instar	69	0.689
Residual_SoftHebb-Surround/HardWTA/Cos-Instar	71	0.708
Dale_Depthwise_HardHebb-Surround/HardWTA/Cos-Instar	62	0.615
SoftHebb-Backpropagation	70.8	0.704

The remaining configurations attained results closer to previous Hebbian-CNN literature, around 60% accuracy, with the low accuracy of 19% of **SoftHebb-None-Instar** demonstrating the necessity of competition mechanisms.

This high accuracy level in our best configurations was achieved by training for only 5 epochs rather than the 50-20 epochs employed in research, and further improvements can be obtained when allowing the classifier to converge its weights. Training a model for longer epochs can improve performance, but it can also conceal the expected accuracy for shorter training regimes. For comparison, Journé et al. (2022) achieved 80.3% accuracy after training the classifier for 50 epochs, Miconi (2021) achieved 64.6% accuracy after training the classifier for 20 epochs, and Lagani, Falchi, Gennaro & Amato (2022) achieved 64.43% under 20 epochs.

To ensure robustness and stability against randomness, the best models in terms of performance and parameters were evaluated in two additional runs with different random seeds. One of the runs was evaluated for 10 epochs to showcase improved performance could be achieved with longer training times.

Table 4: Accuracy for Experimental Configurations on different Fixed Seeds per Run

Configuration	Run 1 (%)	Run 2 (%)	Run 3 (%)	Mean (%)	10 Epoch Acc.(%)
SoftHebb-SoftWTA-Instar	73.6	75	72.6	73.73	75.2
SoftHebb-Surr/HardWTA/Cos-Instar	73	73.35	73.48	73.27	74.29
Depthwise_SoftHebb-Surround/HardWTA/Cos-Instar	69	68.9	69	68.97	70.5

Table 4 presents these robustness results. **SoftHebb-SoftWTA-Instar**, **SoftHebb-Surround/HardWTA/Cos-Instar** and **Depthwise_SoftHebb-Surround/HardWTA/Cos-Instar** maintained their high performance, and across these runs the largest difference was 1.27% accuracy from the mean for the **SoftHebb-SoftWTA-Instar** configuration. When allowing the classifier to train for 10 epochs, improvements of up to 2% in the accuracy were reached, indicating training the classifier for longer epochs can further improve the final results.

Figure 9 shows the results obtained from configuration **SoftHebb-SoftWTA-Instar** (blue line) and **SoftHebb-Surround/HardWTA/Cos-Instar** (red line). In both cases, similar performance was observed both during training (left panel) and testing (right panel), with a fast accuracy increase at each epoch. Learning was remarkably faster in our Hebbian configurations than

with backpropagation (yellow line) on both train and test data under the shorter training times, with the Hebbian configurations presenting steep slope increases at each epoch compared to the consistently gradual increases in backpropagation, a sign of fast learning.

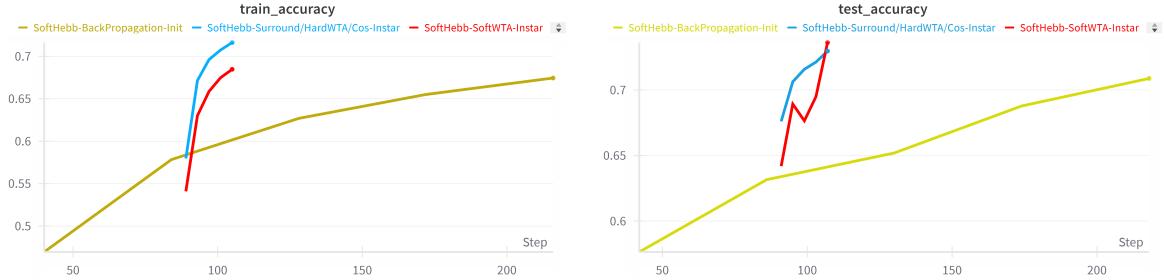


Figure 9: Train and Test accuracy of **SoftHebb-SoftWTA-Instar** and **SoftHebb-Surround/HardWTA/Cos-Instar** configurations compared to Backpropagation during 5 epochs. Epochs are marked incorrectly as WanDB handles steps/epochs inaccurately.

The distribution of weights varied depending on the competition mode and learning rule used, as plotted in Figures 10-12. SoftWTA (Figure 10) produced a distribution with a pronounced central peak and relatively narrow range, closely adhering to its initial Gaussian distribution, indicating only small but distributed changes were made to the original weights. Hard-WTA with the Grossberg Instar rule (Figure 11) resulted in a similar shape but with a wider spread of values. The BCM learning rule (Figure 12) led to a distinct distribution characterised by a sharp peak near zero and heavy tails, indicating a tendency towards weight sparsity while also allowing for some large positive and negative weight values. Deeper layers for these configurations aligned with their respective distribution tendencies.

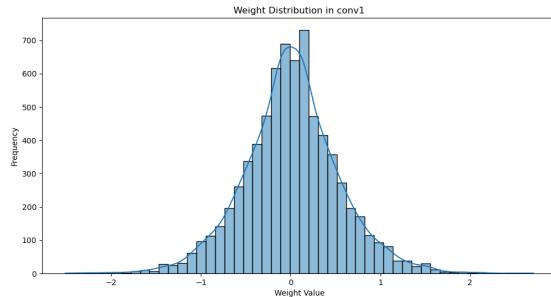


Figure 10: Weight distribution at first layer characteristic of Soft-WTA.

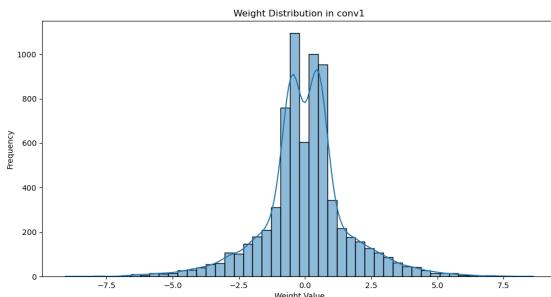


Figure 11: Weight distribution at first layer characteristic of Hard-WTA competition with Grossberg Instar rule.

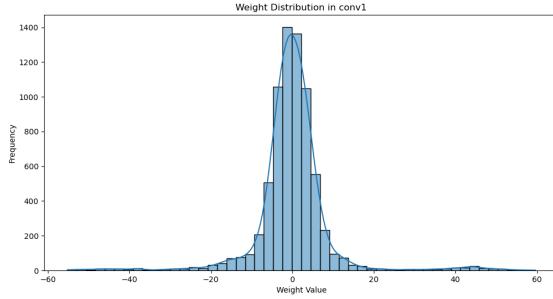


Figure 12: Weight distribution at first layer characteristic of Hard-WTA with BCM learning rule.

4.3 Configurations Qualitative Results

This section emphasises the most relevant information gleaned from our visualisation tools, highlighting the differences in hierarchical representations learnt by our 3-layer configurations.

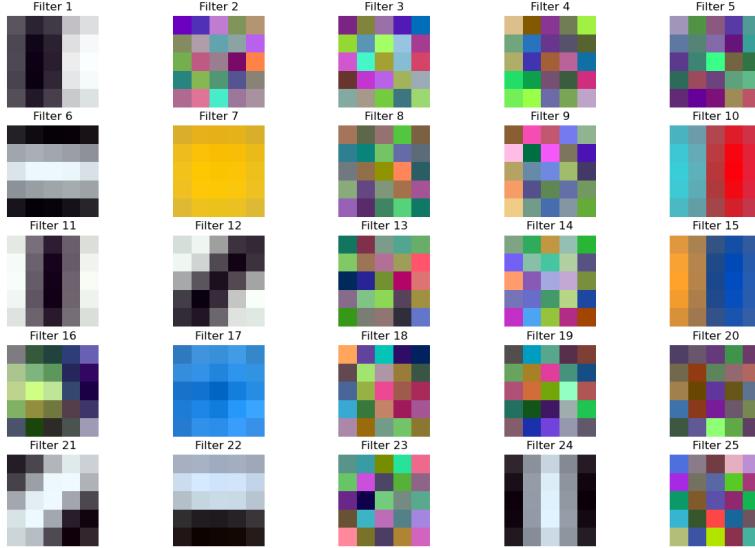


Figure 13: Plots of filters from first 25 neurons directly plotted for first layer of the **SoftHebb-Surround/HardWTA/Cos-Instar** configuration.

Figures 13 and 14 showcase the direct filters learnt at the first layer for the highest-accuracy experiments (**SoftHebb-Surr/HardWTA/Cos-Instar** and **SoftHebb-SoftWTA-Instar**), presented as 5x5 grids where each square represents a different filter or neuron, labelled from *Filter 1* to *Filter 25*. The remaining filters or neurons results per layers were also presented with this structure.

The **SoftHebb-Surr/HardWTA/Cos-Instar** configuration yielded kernels with edges or colour combinations (Filters 12 and 15 respectively in Figure 13), and the **SoftHebb-SoftWTA-Instar** configuration primarily strengthened or diminished initial weight intensities (Filters 11 and 20 respectively in Figure 14), with occasional colour-dominant kernels.

Although there is no single criterion for optimal representations since they depend on each dataset, learning rule, and architecture, we generally expected robust representations in the first convolutional layer will merge colour and edges, resulting in distributed representations across various filters. Similar results were obtained with a backpropagation model (Appendix A.1).

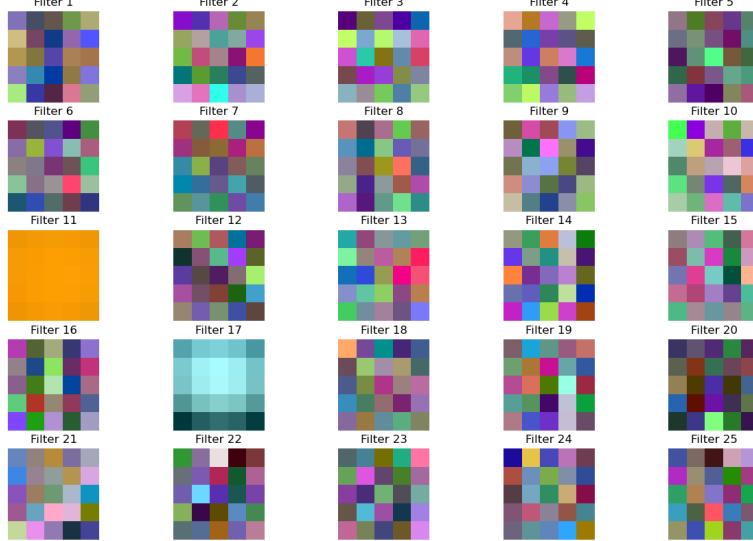


Figure 14: Plots of filters from first 25 neurons directly plotted for first layer of the **SoftHebb-SoftWTA-Instar** configuration.

The direct filter plots focused on the patterns learnt by the filters. The filters in the first convolutional layer of configurations with Hard-WTA and Soft-WTA competition, as seen by comparing Figures 13 and 14 with Figure 22, matched the representations in the receptive fields. However in deeper layers (all layers after the first convolution), as seen by comparing Figures 17 to Figure 21, the direct plots were only capable of identifying patterns like centre-surround kernels or repeating patterns. It could not capture the representations to which neurons were responsive, as the mean across channels reduces these direct kernel plots to patterns.

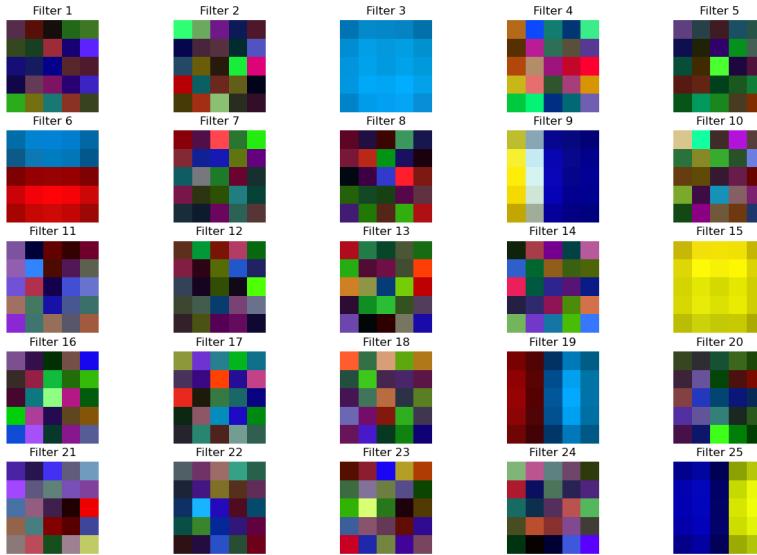


Figure 15: Plots of filters from first 25 neurons directly plotted for first layer of **Dale_HardHebb-Surr/HardWTA/Cos-Instar** configuration.

Figures 15 and 16 highlight the unique filters for **Dale_SoftHebb-SoftWTA-Instar** and **SoftHebb-HardWTA/Cos-BCM** configurations that set them apart from other configurations at the first layer. For the Dale network operating under biological constraints, specialised colour filters emerged (Filter 25 in Figure 15). Models trained using the BCM rule started to combine colours and edge orientations together in its filters (Filter 2 in Figure 16).

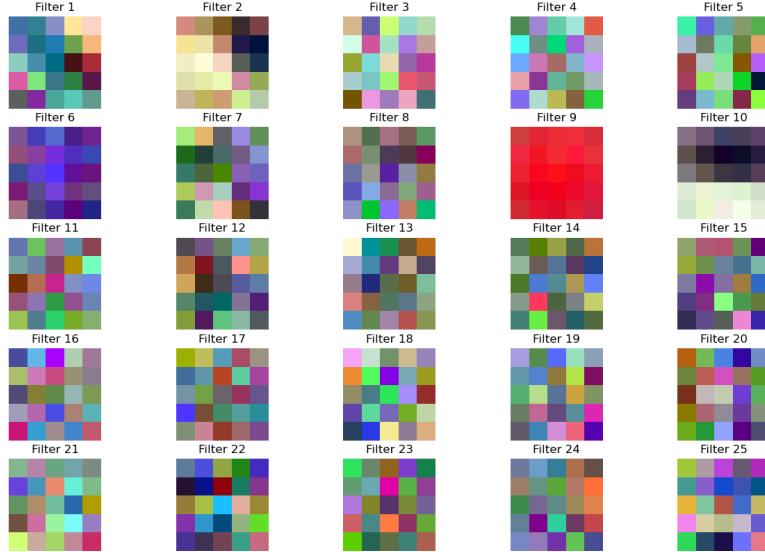


Figure 16: Plots of filters from first 25 neurons directly plotted for first layer of the **SoftHebb-HardWTA/Cos-BCM** configuration.

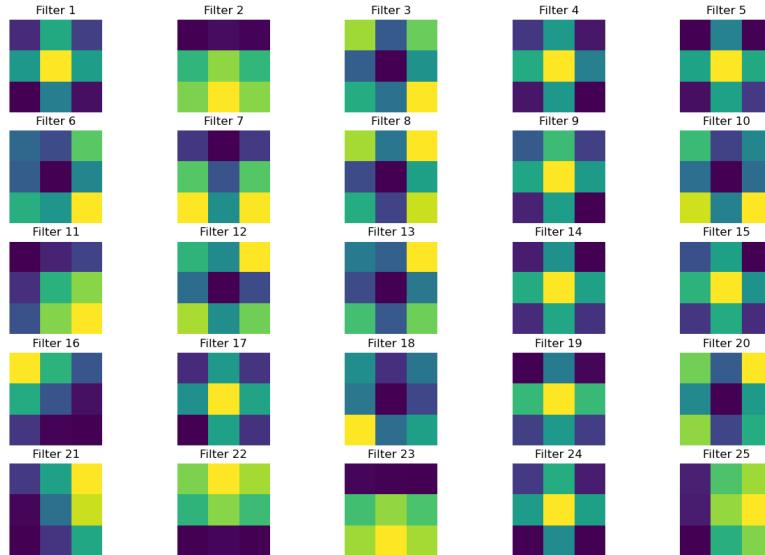


Figure 17: Filter patterns from first 25 neurons for last layer of the **Depthwise_SoftHebb-Surround/HardWTA/Cos-Instar** configuration

Figures 17 and 18 present the filter patterns learnt at the last layer for **Depthwise_SoftHebb-Surround/HardWTA/Cos-Instar** and **Dale_SoftHebb-Surround/HardWTA/Cos-Instar** configurations. The Depthwise model produced different kernel patterns, with a dominance of centre-surround kernels (Filters 11-13 and 4-5 respectively in Figure 17), while the Dale Principle experiment formed exclusively centre-surround kernels.

This natural emergence of centre surround kernels with both **Depthwise_SoftHebb-Surround/HardWTA/Cos-Instar** and **Dale_SoftHebb-Surround/HardWTA/Cos-Instar** reinforces the biological plausibility of combining Hebbian learning and depthwise convolutions. Both present different approaches to learning in deeper layers, either presenting a diverse configuration of patterns or a dominance of biologically realistic centre-surround patterns.

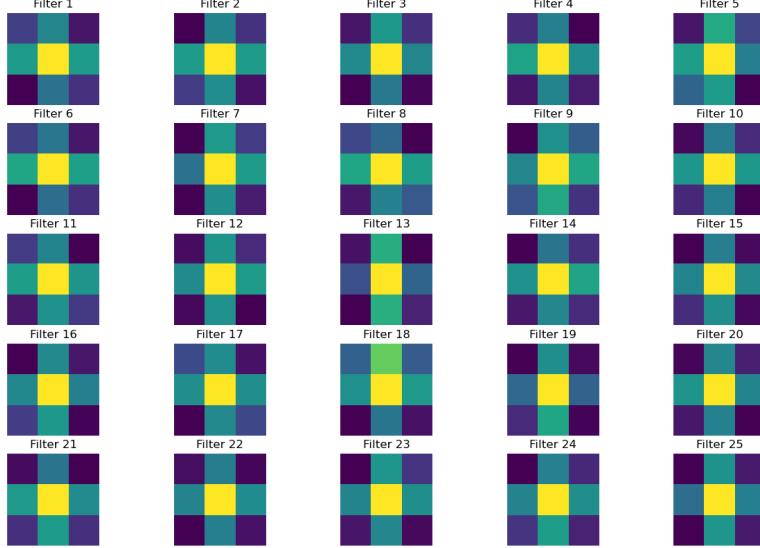


Figure 18: Filter patterns from first 25 neurons for last layer of the **Dale_SoftHebb-Surround/HardWTA/Cos-Instar** configuration.

The class clustering and separation focused on the learnt final representations from a neural network of classes in the dataset. A successfully learnt representation will cluster instances of a same class together and will separate them from instances of other classes. Classes which are similar to each other, such as species of a same animal, will be clustered together given their similarity in features. Colours are used to label classes and help identify the formation of clusters.

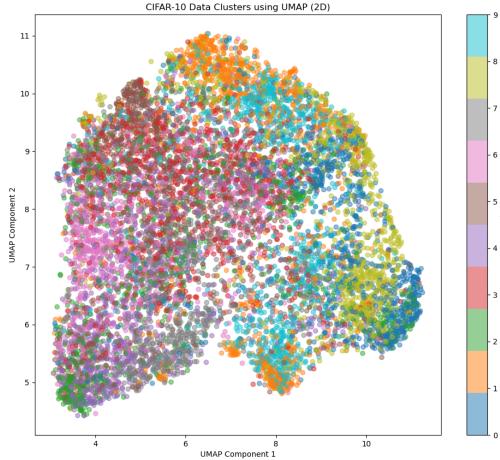


Figure 19: Plots of class clustering/separation of the **SoftHebb-Surround/HardWTA/Cos-Instar** configuration.

Figures 19 and 20 display the class clustering obtained from our highest accuracy **SoftHebb-Surround/HardWTA/Cos-Instar** and **SoftHebb-SoftWTA-Instar** configurations. Each class formed its own cluster or different clusters, and attempts at separation between clusters was found in both models, with **SoftHebb-SoftWTA-Instar** displaying more concentrated clusters. These findings suggest hierarchical learning is occurring, with the model recognising features shared by instances within the same classes and distinguishing them from those in different classes.

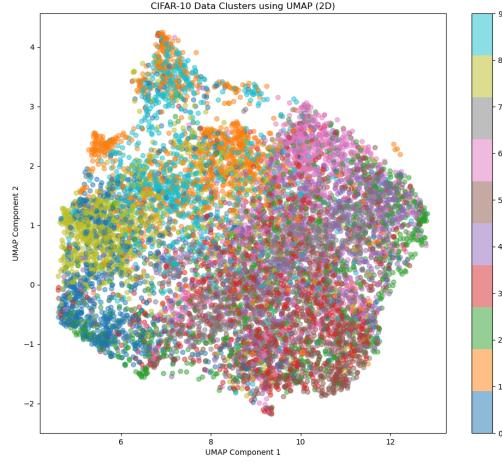
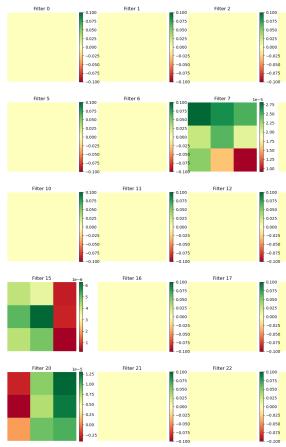


Figure 20: Plots of class clustering/separation of the **SoftHebb-SoftWTA-Instar** configuration.

Our new plots, LTP/LTD per pixel of a kernel, indicated whether the mean of a pixel in a filter grows or decreases in strength as well as when weights have converged and no further updates will change their behaviour. The stored change in weights in the temporary memory buffer was used to determine the values of LTP/LTD. Pixels shaded in green denoted Long-Term Potentiation (LTP) or an increase in weight strength, while those shaded in red represented Long-Term Depression (LTD) or a decrease in weight strength. The neutral beige tone meant no change in weight strength, indicating weight convergence.

Temporal Thresholding and Hard-WTA

LTP/LTD per Weight (Mean across channels)



Hard-WTA

LTP/LTD per Weight (Mean across channels)

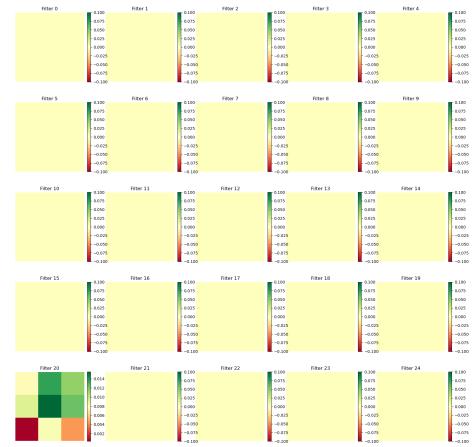


Figure 21: New LTP/LTD plots for first 25 neurons of last layer of **SoftHebb-Temp/HardWTA/Cos-Instar** and **SoftHebb-HardWTA/Cos-Instar** configurations.

Figure 21 illustrates an example of our novel visualisation method to understand convergence of Hebbian learning in the **SoftHebb-Temp/HardWTA/Cos-Instar** and **SoftHebb-HardWTA/Cos-Instar** configurations. Temporal selectivity permitted more neurons to update over extended timesteps, whereas in the absence of this competition mechanisms, most weights had converged.

Receptive fields captured the inputs which maximally activate neurons in a specific layer. The expected representations in the first layer should be simple patterns such as edges and colours which match the direct visualisations from Figure 13 and 14, and increasingly abstract shapes in deeper layers (Figure 24). This hierarchical process was seen in the receptive fields of the same configurations at the first and last layers in Figures 22 and 24.

The deeper layers' receptive fields lacked recognisable object shapes (anatomical features of a cat as an example), as we processed 32x32 pixel images through small 3-layer models. This constraint resulted in abstract feature extraction rather than object-specific features from classes.

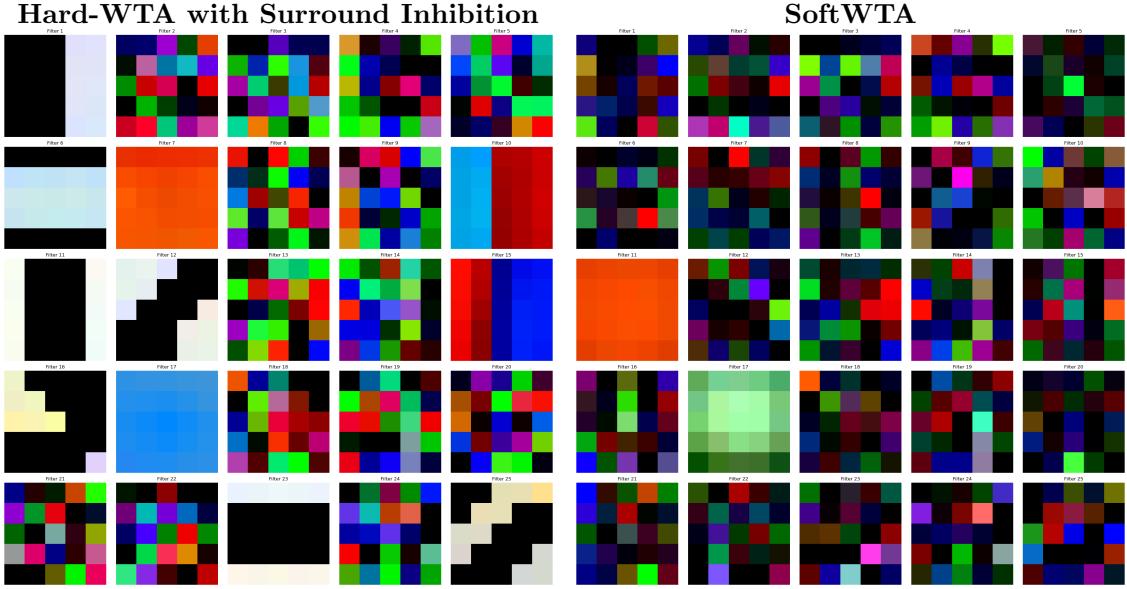


Figure 22: Plots of receptive fields for first 25 neurons of first layer of **SoftHebb-Surround/HardWTA/Cos-Instar** and **SoftHebb-SoftWTA-Instar** configurations.

Figure 22 displays the receptive fields for neurons at the first convolutional layer for our **SoftHebb-Surround/HardWTA/Cos-Instar** and **SoftHebb-SoftWTA-Instar** configurations. Edges with orientations, selective uniform colours or opposing colours (orange-blue) activated neurons trained with Hard-WTA competition, matching the direct visualisation of filters at this layer from Figures 13-14. On the other hand, Soft-WTA retained the initial random patterns from initialisation, and adjusted the intensity of these patterns, causing receptive fields at the first layer to focus on distributed representations.

Figure 23 displays the results of **SoftHebb-None-Instar** and **SoftHebb-Stats/HardWTA/Cos-Instar** configurations. Without competition, neither meaningful edge orientations nor distributed patterns were captured by the receptive fields, highlighting the importance of competition mechanisms to capture discriminative features. Statistical thresholding allowed all neurons to learn different information, proving its homeostatic capabilities. However, the neurons could not form distinct shapes and only colour sensitive receptive fields, indicating hard-WTA competition was not functioning strongly and homeostasis was instead dominating competition.

Figure 24 displays the receptive fields at the deepest layer of our best configurations, **SoftHebb-Surround/HardWTA/Cos-Instar** and **SoftHebb-SoftWTA-Instar**. Both displayed abstract shapes and patterns, which combined different colours with increasing complexity. The shapes of patterns with Surround Lateral Inhibition were smoother when compared to both **SoftHebb-SoftWTA-Instar** and its counterpart without this competition mechanism in Figure 25, indicating lateral inhibition was providing an additional smoothing property on the learnt representations, which was not captured in previous literature.

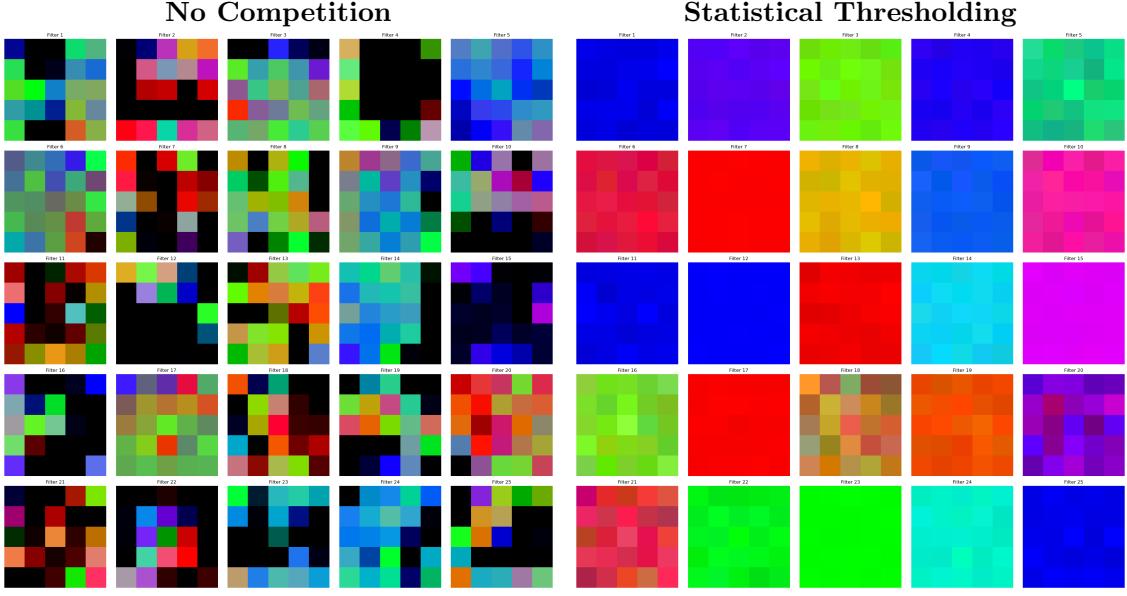


Figure 23: Plots of receptive fields for first 25 neurons of first layer of **SoftHebb-None-Instar** and **SoftHebb-Stats/HardWTA/Cos-Instar** configurations.

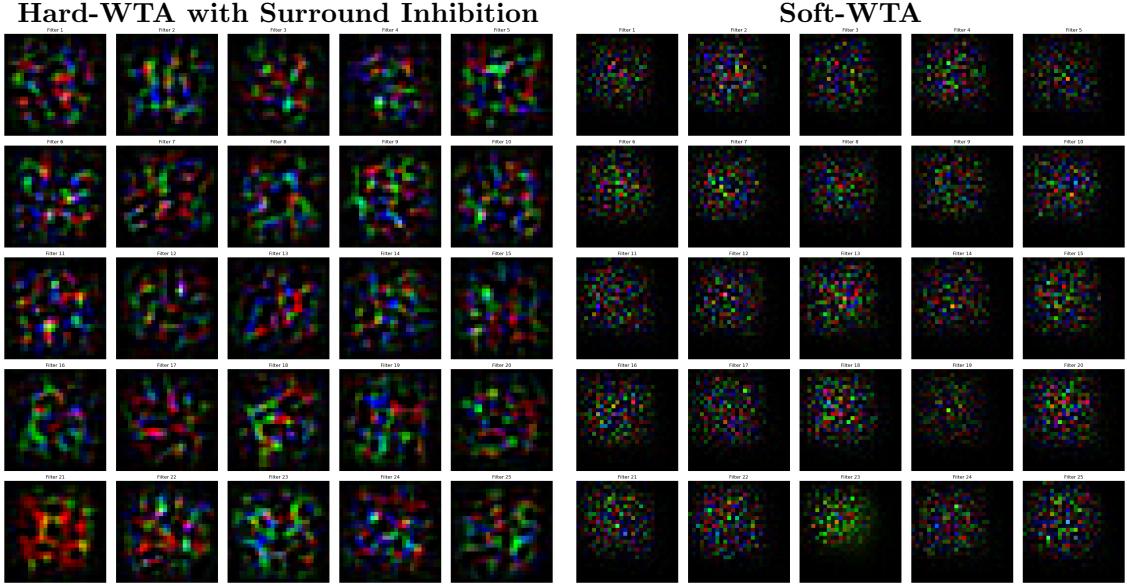


Figure 24: Plots of receptive fields for first 25 neurons of last layer of **SoftHebb-Surround/HardWTA/Cos-Instar** and **SoftHebb-SoftWTA-Instar** configurations.

Figure 25 presents insightful results for the receptive fields at the deepest layer of **SoftHebb-HardWTA/Cos-Instar** and **SoftHebb-None-Instar** configurations. The receptive fields of **SoftHebb-HardWTA/Cos-Instar** were comparable to **SoftHebb-SoftWTA-Instar**, with distributed representations of abstract patterns with diverse colours. The receptive fields without any form of competition mechanisms presented repeated patterns of single colours, which indicated unsuccessful representational learning.

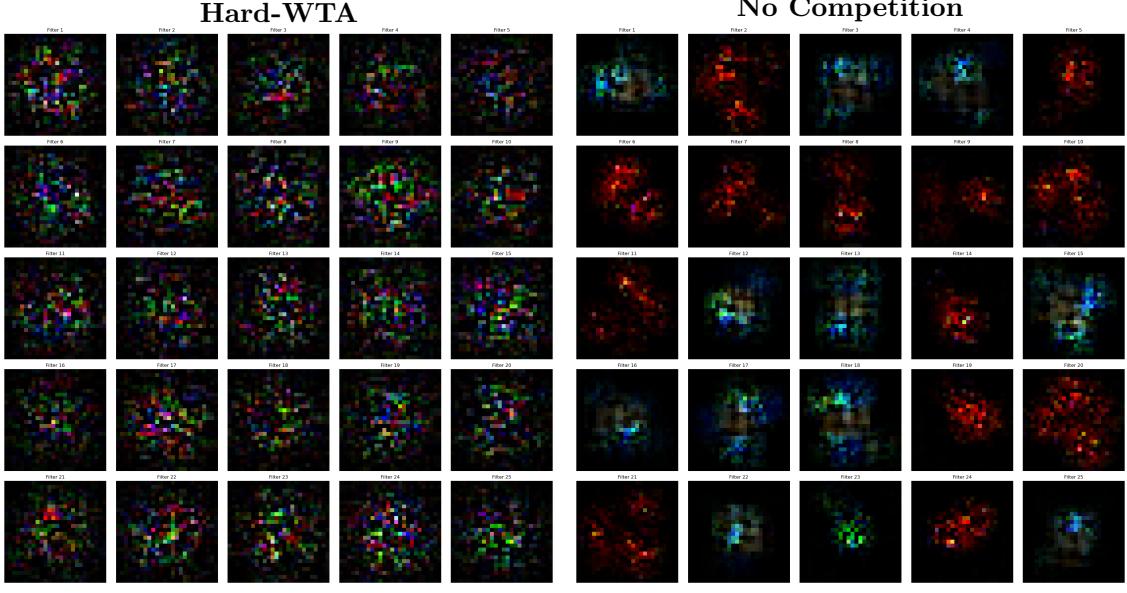


Figure 25: Plots of receptive fields for first 25 neurons of last layer of **SoftHebb-HardWTA/Cos-Instar** and **SoftHebb-None-Instar** configurations.

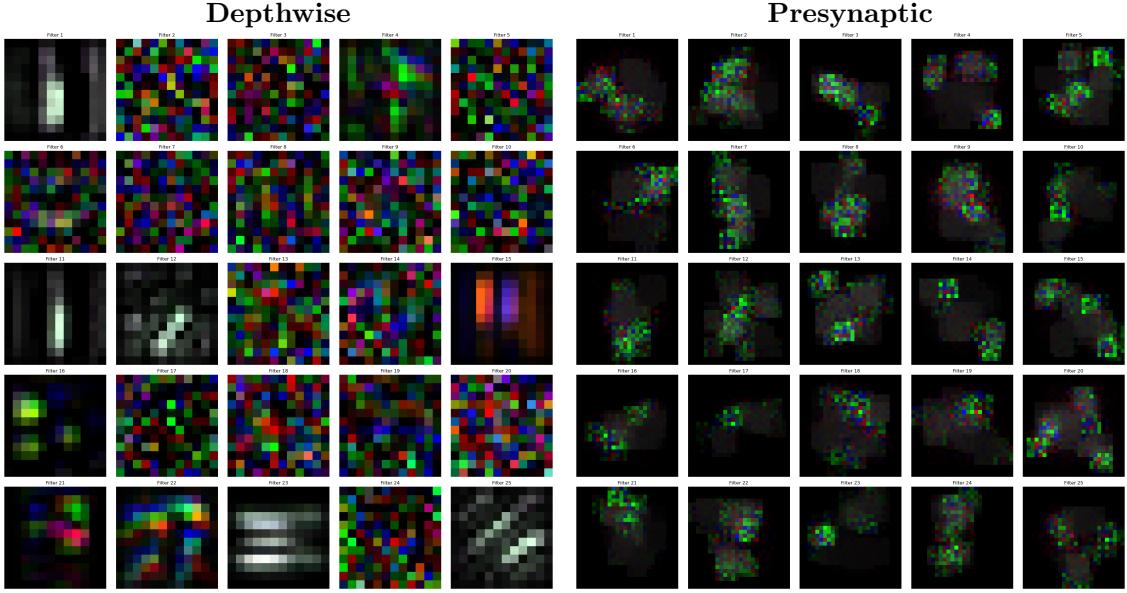


Figure 26: Plots of receptive fields for first 25 neurons of second separable layer of **Depthwise_SoftHebb-HardWTA/Cos-Instar** and last layer of **SoftHebb-Presynaptic/HardWTA/Cos-Instar** configurations.

Figure 26 illustrates distinctive results obtained from the receptive fields of the second layer of the depthwise architecture with **Depthwise_SoftHebb-HardWTA/Cos-Instar** and the unique results obtained at the last layer with **SoftHebb-Presynaptic/HardWTA/Cos-Instar**. For the depthwise configuration, the receptive fields of the depthwise convolution, which focuses solely on spatial features, correctly displayed spatial information, with a focus on edges and colour combinations even at deeper layers, as well as the emergence of more complex patterns and shapes. For the presynaptic competition, unique and diverse clusters dominated by green intensities were captured by neurons.

5 Discussion

5.1 Analysis of results

The experimental results demonstrate varying efficacy across different Hebbian learning approaches in CNNs. **SoftHebb-SoftWTA-Instar** emerges as a superior method, achieving over 70% accuracy across all benchmarks. The weight distribution in Figure 11 closely approximates the original normal distribution 8, suggesting stability and distributed representations within the network. Notably, **SoftHebb-SoftWTA-Instar** exhibits robust class clustering and separation, indicative of effective hierarchical learning in Figure 20. This performance can be attributed to its capacity to form incrementally complex and abstract representations, as demonstrated by the receptive fields of neurons at each layer in Figure 24.

The HardWTA mode in **Lagani-HardWTA/Cos-Instar**, implemented with the original Lagani Architecture and cosine similarity activation during both forward propagation and weight updates as proposed by Amato et al. (2019), achieves an accuracy exceeding 55%. This approach forms basic representations in its receptive fields, progressing from edges and colours in earlier layers as seen in Figure 13 to simple shapes and diverse patterns in deeper layers. While class clustering is observed in Figure, it demonstrates limited efficacy in separating similar classes, indicating basic hierarchical learning. The simple representations can be seen in Appendix A.1 Figure 30.

Further experiments were conducted to assess the impact of different architectures on the performance of hard-WTA competition in **Lagani_Deep-HardWTA/Cos-Instar**. Notably, the original 4-layer architecture by (Amato et al. 2019) attained 53% accuracy, supporting the reported performance decline with the addition of deeper layers. Significant improvements were realised when we adopted an architecture similar to SoftHebb without padding, denominated as **HardHebb** model. This modified approach yielded substantial accuracy improvements, surpassing 65% accuracy.

A critical result emerged from our implementation of HardWTA with cosine similarity activation function and the exact SoftHebb Architecture. This **SoftHebb-HardWTA/Cos-Instar** configuration achieved state-of-the-art performance with an impressive 72% accuracy, representing the highest performance attained in the deepest network yet for a fully Hard-WTA Hebbian-CNN model. The hierarchical representations formed are comparable to those of SoftHebb, with each layer exhibiting abstract and distributed learned representations as seen in Figure 25.

Given its superior performance, the SoftHebb architecture was adopted as the default CNN model for subsequent experiments. With the absence of competition in **SoftHebb-None-Instar**, the network's performance deteriorates significantly, reaching only 19% accuracy due to its inability to extract meaningful features. This is demonstrated by a focus on colours in the receptive fields of neurons across all layers as seen in Figure 23 and 25, representational collapse in class clusters, and signs of weight convergence as indicated by a negatively-skewed distribution with a sharp peak.

The BCM learning rule in **SoftHebb-HardWTA/Cos-BCM** also achieved state-of-the-art performance with 71% accuracy but behaved vastly differently than its Grossberg Instar counterpart. While the deeper representations remain similar, the initial layer displays filters which combine distributed representations with edges. The distribution of weights are also significantly distinct with a sparse distribution as seen in Figure 12, where the mean is close to zero but presents a large standard deviation. This suggests sparse connections could be removed for further optimisation and reduction in parameters.

Novel temporal and statistical competition modes were evaluated alongside Hard-WTA competition using the SoftHebb architecture. The temporal competition in **SoftHebb-Temp/HardWTA/Cos-Instar** obtained comparable results to **SoftHebb-HardWTA/Cos-Instar**, with 71% accuracy.

Weights had also not all converged yet, as seen by Figure 21. Notably, the absence of WTA specialisation led to representational collapse in clusters and weight convergence, highlighting its necessity.

The statistical threshold method in **SoftHebb-Stats/HardWTA/Cos-Instar**, when evaluated with WTA competition, yielded an accuracy of 60%. Interestingly, this approach demonstrated a stronger response to colours rather than shapes and patterns as seen in Figure 23. Despite the simplicity of these filters, the model successfully separated class clusters.

Various iterations of presynaptic competition were tested in **SoftHebb-Presynaptic/HardWTA/Cos-Instar**, with relative success shown only in the input channels competition dimension. Other versions exhibited representational collapse, characterised by minor changes to weights in the first layer but significant alterations in deeper layer patterns. Spatial and output channel-based competition resulted in repeated simple patterns and shapes in the receptive fields and filters, showcasing convergence of weights and redundant representations. These receptive fields were dominated by either green colour sensitivity as seen in Figure 26 or exclusively RGB colour dominance in the case of output and spatial competition modes, while global presynaptic competition produced noisy receptive fields. Examples of these representations can be seen in Appendix A.1 Figure 32.

The simple surround inhibition kernel between neurons of a same filter in **SoftHebb-Surround/HardWTA/Cos-Instar** provided a refinement over **SoftHebb-HardWTA/Cos-Instar**, enhancing edges and smoothing uniform features achieving a light improvement with an accuracy of 73%, a new milestone in Hard-WTA competition. The receptive fields of deeper layers are smoother abstract shapes than the previous best hard-WTA configuration, which could permit increased generalisability, as seen in Figure 24. Given this improvement, it is used in all subsequent Hard-WTA configurations.

A depthwise architecture equivalent to the SoftHebb model with **Depthwise_SoftHebb-Surr/HardWTA/Cos-Instar**, achieved near state-of-the-art performance with 69% accuracy while reducing parameters by a factor of 6.6 (from 5,893,162 to 883,978 parameters). Visualisation of the receptive fields in depthwise convolutions provided unique insights into the learning process of this novel architecture. The second layer demonstrated a larger focus on edges and shapes, with first-layer filters producing edges also generating edges in their counterparts in the following layer as seen in Figure 26, while noisier distributed representations in the first layer led to similar patterns in the subsequent filter counterparts.

The network partially formed by residual blocks of depthwise separable convolutions in **Residual_SoftHebb-Surr/HardWTA/Cos-Instar**, similar to MobileNet, performed comparably to **Depthwise_SoftHebb-Surr/HardWTA/Cos-Instar** in terms of class separation and weight distribution across layers. A notable difference was observed in the features detected by the receptive fields of the second layer, where residual blocks focused on specific spots of colour rather than edges and lines. This difference can be noted in Appendix A.1 Figure 31.

Lastly, a depthwise network fully respecting Dale’s principle with only excitatory synapses was evaluated to enhance the biological plausibility of the model in **Dale_SoftHebb-Surr/HardWTA/Cos-Instar**. The SoftHebb architecture with padding produced inferior results, forming only basic filters of colours and edges. In contrast, the HardHebb architecture without padding improved performance, allowing for the formation of edge and colour filters.

Significantly, in the depthwise convolutions of the Dale variant, most filters exhibited centre-surround patterns as seen in Figure 18, further accentuating the biological plausibility of this network. This difference may help explain the biological implausibility of padding, as neurons cannot create artificial regional extensions to inputs from beyond their receptive fields, further proving the biological plausibility of this variant.

Comparative analysis of backpropagation using an identical 3-layer architecture to SoftHebb in [SoftHebb-Backpropagation](#) revealed comparable performance outcomes, irrespective of whether the weight initialisation replicated our Hebbian models or employed the default PyTorch initialisation. A notable observation emerged when utilising our initialisation scheme as the initial filter patterns remained largely unaltered, with only subtle, imperceptible modifications occurring. In contrast, the default weight initialisation led to the emergence of filters in the first layer that combined edges and colours, while deeper layers exhibited more significant alterations to the original patterns, occasionally producing centre-surround kernels. These results can be found in Appendix A.1 Figures 27, 28, and 29.

Accuracy metrics for both initialisation methods converged at approximately 70%, highlighting the competitive efficacy of Hebbian learning in shallow networks. However, it is important to note that the strengths of backpropagation become more pronounced with the incorporation of additional deep layers and extended training epochs. These augmentations facilitate further optimisation, surpassing the performance ceiling of Hebbian learning approaches in more complex network architectures.

5.2 Objectives Evaluation

This research successfully addressed the established objectives, contributing significantly to the field of Hebbian learning in Convolutional Neural Networks and benefiting AI researchers, neuroscientists and future industry applications. The outcomes for each objective are as follows:

The **first objective** of developing a baseline model with a simple Hebbian learning rule, utilising datasets and parameters consistent with existing research for fair comparisons was successfully attained. This achievement facilitated a fundamental understanding of basic Hebbian learning and its integration into CNNs, while also enabling rigorous examination of previous implementations.

The **second objective** of implementing state-of-the-art research and ensuring similar results was completed despite encountering unanticipated complexities that extended the development timeline. It was the most essential step, as through understanding these implementations and their inner workings we successfully replicated results from literature. To complete this objective, some implementation details were extracted from the code as they did not appear in the published papers, with the aim of facilitating future implementation.

The **third objective** of understanding limitations in these implementations and researching extensions to enhance performance was mostly completed. As this is a novel field inside Artificial Intelligence, research regarding Hebbian-CNN literature was limited to specific papers, and additional interdisciplinary insights from biological-inspired methodologies and neuroscience-based research were used to find some extensions to state-of-the-art methods. Optimisation of these methods was a major focus, from a computational and structural perspective.

Finally, the **last objective** of comprehensive experiments and comparative analysis between experiments was performed successfully. This was accomplished with correct evaluation metrics and by devising a suite of visualisation tools new to this field and to the machine learning community. These new tools allow an interpretable and explainable understanding of learning processes in neural networks, with a focus on Hebbian learning but applicable to backpropagation as well.

5.3 Key Findings

The key finding in this project has been the success of combining Hard-WTA competition and Hebbian learning to achieve state-of-the-art performance through an architecture with the

deepest number of Hebbian layers in literature. The scaling factor of number of filters per layer from SoftHebb also has a large impact on the performance of Hard-WTA competition, undocumented in literature before this project.

Extensions to competition, learning rules and lateral interactions were evaluated, with noted success in the incorporation of a surround modulation kernel to simulate lateral inhibition, improving by 1% accuracy of our best Hard-WTA configuration and creating smoother representations which can aid with generalisability. Temporal competition alongside hard-WTA allowed updates to occur in deeper layers during longer timesteps before weight convergence, indicating more selectivity in the updates of neurons, while statistical thresholds provided a type of indirect homeostasis, albeit an overly aggressive form.

BCM learning rule was integrated into Hebbian-CNN and behaved differently from all other experiments. Rather than retaining the basic form of a Gaussian distribution in its weights, it creates a very sparse distribution with most values close to 0 and incrementally larger heavy tails. Even with sparse weights, it can form efficient representations which combine edge and colour information and achieve state-of-the-art performance, indicating its efficiency and adequacy in neural networks.

The first implementation of Hebbian Depthwise Separable Convolutions in this field substantially decreases the parameter count by a factor of 6 while maintaining state-of-the-art performance. This reduction in parameters from a computational standpoint can alleviate the increased number of filters required to achieve high precision, while increasing its biological plausibility. This architecture is beneficial for both Soft and Hard-WTA competition.

Residual blocks increased the depth of the network by including additional pointwise convolutions, achieving high accuracy on a deeper network, reinforcing the importance of additional connections such as skip connections to allow deeper training in Hebbian networks. A Depthwise Separable Convolutional Hebbian network trained following Dale's principle naturally produces a type of centre-surround kernels in deeper layers, supporting the biological plausibility of this convolution with Hebbian learning and strengthening the research of Babaiee et al. (2024).

These findings exceeded our initial expectations for the objectives, yielding outstanding outcomes through a learning dynamic that is both computationally efficient and biologically plausible, capable of competing with the widely-used backpropagation algorithm, offering a new optimised training paradigm for Artificial Intelligence.

5.4 Answer to Research Question

The main question for this research project was: **Can Hebbian learning algorithms achieve representational learning comparable to backpropagation while maintaining biological plausibility?**

This question can be answered affirmatively with many different Hebbian configurations under certain constraints. In shallow networks of 3 layers and under training time constraints, Hebbian Learning in CNN with competition mechanisms can achieve or surpass backpropagation at a faster rate for both training and test sets. This result is encouraging for the application of this biologically plausible learning rule in more complex environments and tasks, and may possess stronger generalisation capabilities given its unsupervised foundation.

Further sub-questions where also posed:

- How do different Hebbian learning rules affect the network's ability to form hierarchical representations?

This question can be answered with the differences in behaviour found between the

Grossberg Instar and the BCM learning rule under the same conditions, where the first rule supports a broader spread of synaptic strengths while BCM learning rule focuses on selectivity through sparsity.

- What role does neuronal competition play in Hebbian learning?

Competition between neurons is crucial for creating meaningful representations which extract discriminative features from data to achieve high performance. Hard-WTA competition can attain comparable performance to Soft-WTA competition, and lateral mechanisms can help reinforce competition and selectivity between neurons.

6 Evaluation, Reflection and Conclusion

6.1 Evaluation and Reflection of Objectives and Findings

The objectives laid out permitted a structured approach to the project, starting from a solid foundation with a baseline and improving this starting point with state-of-the-art Hebbian-CNN research. Once results from literature could be obtained, extensions and understanding of the limitations of these methods became the intuitive step to improve on these results. Finally, carefully controlled and rigorous experiments facilitated an analysis of emergent behaviours from competition mechanisms between neurons and how they interact with different learning rules and architectures.

This bottom-up approach provided flexibility and focus at each step, and these properties permitted the successful completion of these objectives. The structured approach facilitated by our objectives provided a robust framework for the project’s progression, systematically incorporating state-of-the-art Hebbian-CNN research and extending our understanding of these methods’ limitations.

Regarding our main findings, we are pleased by the results obtained through our thorough analysis of hard-WTA competition. Achieving state-of-the-art performance through architectural changes is a simple yet effective modification, emphasising the importance of increasing diversity in neurons to allow correct winner selection. This demonstrates that simple yet effective changes can yield substantial improvements.

Further architectural changes to preserve more biological plausibility while reducing parameters also maintained SOTA accuracy, a surprising result alongside the natural emergence of centre-surround kernels when following Dale’s Principle. This unexpected outcome highlights the potential for biologically inspired approaches to yield competitive performance.

Residual networks also achieve close to SOTA performance, achieving the best performance for the deepest Hebbian network yet. This represents a significant advancement in the depth of Hebbian networks while maintaining high accuracy.

Finally, temporal and statistical competition thresholds were designed for the first time, and can be expanded upon to improve the generalisability of the model. Pre-synaptic competition completely changes the behaviour of the model, and further experiments to integrate it with post-synaptic competition must be studied. These novel designs open new avenues for improving model generalisability and understanding the complex interactions between different types of neuronal competition.

6.2 Evaluation of Limitations

Our main limitation has been the lack of full customisation of a model layer by layer, and the control of learning rules and competition mechanisms through command arguments. Currently, all layers in a model have the same competition and learning rules applied to them. To change some of these learning mechanisms, manual changes must be made in their respective files to allow it to work. This constraint hinders the flexibility of experimental setups and the ability to fine-tune individual layers.

Another limitation has been the results obtained from the extended competition mechanisms. While temporal competition has met SOTA performance, without Hard-WTA competition it cannot form representations, indicating Hard-WTA is still the main competition mechanism. Statistical competition only creates simple representations, and further tests on homeostatic mechanisms could be beneficial. These findings suggest that while our extended mechanisms

show promise, there is still room for improvement and more comprehensive understanding of their interactions and individual contributions to model performance.

6.3 Future Work

Our goal, driven by these remarkable findings, is to write this project for publication, sharing the insights gained from merging Hebbian learning and local competition mechanisms with modern deep learning structures. This will advocate the application of Hebbian learning and its advantages to a broader scientific audience, fostering additional research in this innovative learning paradigm.

Ensuring full network customisability layer by layer and control over mechanisms through command arguments would facilitate the speed of setting up experiments and provide finer control of each layer, allowing the combination of different mechanisms in the same experiment. This enhancement would streamline the research process and enable more nuanced investigations into the interactions between different learning rules and competition mechanisms.

Evaluating these findings on different datasets and with deeper networks is paramount to understand the representational capabilities of these networks. STL-10 dataset (Coates et al. 2011) and more real-world application datasets are a good starting point. This expansion would provide a more comprehensive assessment of the model’s generalisability and practical applicability.

Recurrent connections between layers could improve both performance and the biological plausibility of the network. Similarly, feedback connections from deeper layers to lower layers could improve representational formation. These additions could potentially enhance the model’s ability to capture complex temporal dependencies and hierarchical relationships in the data.

To enhance the biological inspiration, distinct excitatory and inhibitory neural populations which adhere to Dale’s Principle could improve formation of representation in a manner approved by neuroscientists. This approach would align our model more closely with current neuroscientific understanding, potentially yielding insights relevant to both artificial intelligence and neuroscience.

6.4 Conclusion

This project provides a competitive alternative to backpropagation through a biologically plausible learning rule, reducing the computation resources required to create meaningful representations for classification tasks. It is able to achieve state-of-the-art performance in a Hebbian-CNN architecture, with over 73% accuracy on CIFAR-10.

Hard-WTA competition was proven to be as effective as Soft-WTA competition, rivalling the performance from the benchmark set by Journé et al. (2022). New extensions including learning rules and competition mechanisms helped improved performance and laid the foundation for further research in these mechanisms.

The code offers a synthesis and optimisation of previous work, facilitating research in this intersection between neuroscience and Artificial Intelligence through an accessible platform designed with modularity and extensibility in mind. This codebase serves as a catalyst for continued improvements in both performance and biological plausibility.

In conclusion, this research contributes significantly to the field by demonstrating the viability of Hebbian learning in deep neural networks, opening new pathways for the development of more efficient and biologically inspired artificial intelligence systems. We will continue building up

from this project to improve performance and biological plausibility, and aim for applying this new technology in real-world tasks.

Ultimately, this will bridge the gap between neuroscientific principles and practical artificial intelligence applications, bringing the biological benefits of rapid training, reduced computational resources and energy efficiency minimising the effect on the ecosystem, and increased generalisability without reliance of target values in a fully unsupervised fashion.

References

- Amari, S.-i. (1993), ‘Backpropagation and stochastic gradient descent method’, *Neurocomputing* **5**(4-5), 185–196.
- Amato, G., Carrara, F., Falchi, F., Gennaro, C. & Lagani, G. (2019), Hebbian learning meets deep convolutional neural networks, in ‘Image Analysis and Processing–ICIAP 2019: 20th International Conference, Trento, Italy, September 9–13, 2019, Proceedings, Part I 20’, Springer, pp. 324–334.
- Babaiee, Z., Kiasari, P. M., Rus, D. & Grosu, R. (2024), Neural echos: Depthwise convolutional filters replicate biological receptive fields, in ‘Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision’, pp. 8216–8225.
- Becker, S. & Plumley, M. (1996), ‘Unsupervised neural network learning procedures for feature extraction and classification’, *Applied Intelligence* **6**(3), 185–203.
- Bengio, Y. & Frasconi, P. (1993), ‘Credit assignment through time: Alternatives to backpropagation’, *Advances in neural information processing systems* **6**.
- Bienenstock, E. L., Cooper, L. N. & Munro, P. W. (1982), ‘Theory for the development of neuron selectivity: orientation specificity and binocular interaction in visual cortex’, *Journal of Neuroscience* **2**(1), 32–48.
- Campbell, J. et al. (2022), ‘Considerations of biological plausibility in deep learning’, *Cornell Undergraduate Research Journal* **1**(1), 4–12.
- Cekic, M., Bakiskan, C. & Madhow, U. (2022), ‘Towards robust, interpretable neural networks via hebbian/anti-hebbian learning: A software framework for training with feature-based costs’, *Software Impacts* **13**, 100347.
- Choe, Y. (2022), Anti-hebbian learning, in ‘Encyclopedia of Computational Neuroscience’, Springer, pp. 213–216.
- Chollet, F. (2017), Xception: Deep learning with depthwise separable convolutions, in ‘Proceedings of the IEEE conference on computer vision and pattern recognition’, pp. 1251–1258.
- Churhe, P. (2024), ‘Unraveling deep learning: A comprehensive guide to neural networks’. Accessed: 2024-09-28.
URL: <https://prathmeshchurhe.medium.com/unraveling-deep-learning-a-comprehensive-guide-to-neural-networks-2623351da219>
- Coates, A., Ng, A. & Lee, H. (2011), An analysis of single-layer networks in unsupervised feature learning, in ‘Proceedings of the fourteenth international conference on artificial intelligence and statistics’, JMLR Workshop and Conference Proceedings, pp. 215–223.
- Cornford, J., Kalajdzievski, D., Leite, M., Lamarquette, A., Kullmann, D. M. & Richards, B. (2020), ‘Learning to live with dale’s principle: Anns with separate excitatory and inhibitory units’, *bioRxiv* pp. 2020–11.
- Fukushima, K. (1980), ‘Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position’, *Biological cybernetics* **36**(4), 193–202.
- Gabbott, P. & Somogyi, P. (1986), ‘Quantitative distribution of gaba-immunoreactive neurons in the visual cortex (area 17) of the cat’, *Experimental brain research* **61**, 323–331.
- Grossberg, S. (1976), ‘Adaptive pattern classification and universal recoding: I. parallel development and coding of neural feature detectors’, *Biological cybernetics* **23**(3), 121–134.

- Gupta, M., Modi, S. K., Zhang, H., Lee, J. H. & Lim, J. H. (2022), ‘Is bio-inspired learning better than backprop? benchmarking bio learning vs. backprop’, *arXiv preprint arXiv:2212.04614* .
- Hasani, H., Soleymani, M. & Aghajan, H. (2019), ‘Surround modulation: A bio-inspired connectivity structure for convolutional neural networks’, *Advances in neural information processing systems* **32**.
- He, K., Zhang, X., Ren, S. & Sun, J. (2016), Deep residual learning for image recognition, in ‘Proceedings of the IEEE conference on computer vision and pattern recognition’, pp. 770–778.
- Journé, A., Rodriguez, H. G., Guo, Q. & Moraitis, T. (2022), ‘Hebbian deep learning without feedback’, *arXiv preprint arXiv:2209.11883* .
- Kerr, D. (2024), ‘Ai brings soaring emissions for google and microsoft, a major contributor to climate change’, *NPR* .
URL: <https://www.npr.org/2024/07/12/g-s1-9545/ai-brings-soaring-emissions-for-google-and-microsoft-a-major-contributor-to-climate-change>
- Kheradpisheh, S. R., Ganjtabesh, M., Thorpe, S. J. & Masquelier, T. (2018), ‘Stdp-based spiking deep convolutional neural networks for object recognition’, *Neural Networks* **99**, 56–67.
- Krizhevsky, A., Hinton, G. et al. (2009), ‘Learning multiple layers of features from tiny images’.
- Kucharlapati, A. (2024), ‘Gradient descent: Gradient descent is an iterative optimization algorithm used to find local minima’. Accessed: 2024-09-28.
URL: <https://medium.com/@kucharlapatiaparna123/gradient-descent-gradient-descent-is-an-iterative-optimization-algorithm-used-to-find-local-minima-874d56a2ead9>
- Lagani, G. (2024), ‘Bio-inspired deep learning research directions’, *Google Docs* .
URL: <https://docs.google.com/document/d/1mRNC0AEkrW7cztbNuXdoz428ZCwbBW8kA2FZKZBh2o/edit>
- Lagani, G., Falchi, F., Gennaro, C. & Amato, G. (2022), ‘Comparing the performance of hebbian against backpropagation learning using convolutional neural networks’, *Neural Computing and Applications* **34**(8), 6503–6519.
- Lagani, G., Falchi, F., Gennaro, C. & Amato, G. (2023), ‘Synaptic plasticity models and bio-inspired unsupervised deep learning: A survey’, *arXiv preprint arXiv:2307.16236* .
- Lagani, G., Gennaro, C., Fassold, H. & Amato, G. (2022), Fasthebb: Scaling hebbian training of deep neural networks to imagenet level, in ‘International Conference on Similarity Search and Applications’, Springer, pp. 251–264.
- LeCun, Y., Bottou, L., Bengio, Y. & Haffner, P. (1998), ‘Gradient-based learning applied to document recognition’, *Proceedings of the IEEE* **86**(11), 2278–2324.
- Lillicrap, T. P., Santoro, A., Marris, L., Akerman, C. J. & Hinton, G. (2020), ‘Backpropagation and the brain’, *Nature Reviews Neuroscience* **21**(6), 335–346.
- Lindsay, G. W. (2018), ‘Deep convolutional neural networks as models of the visual system: Qa’. Accessed: 2024-09-28.
URL: <https://gracewlindsay.com/2018/05/17/deep-convolutional-neural-networks-as-models-of-the-visual-system-qa/>
- Magee, J. C. & Johnston, D. (1997), ‘A synaptically controlled, associative signal for hebbian plasticity in hippocampal neurons’, *Science* **275**(5297), 209–213.
- Mayur, I. (2024), ‘Simple convolutional neural network (cnn) for dummies in pytorch: A step-by-step guide’. Accessed: 2024-09-28.

URL: <https://medium.com/@myringoleMLGOD/simple-convolutional-neural-network-cnn-for-dummies-in-pytorch-a-step-by-step-guide-6f4109f6df80>

- Miconi, T. (2017), ‘Biologically plausible learning in recurrent neural networks reproduces neural dynamics observed during cognitive tasks’, *Elife* **6**, e20899.
- Miconi, T. (2021), ‘Hebbian learning with gradients: Hebbian convolutional neural networks with modern deep learning frameworks’, *arXiv preprint arXiv:2107.01729*.
- Nowlan, S. (1989), ‘Maximum likelihood competitive learning’, *Advances in neural information processing systems* **2**.
- Oja, E. (1982), ‘Simplified neuron model as a principal component analyzer’, *Journal of mathematical biology* **15**, 267–273.
- Pogodin, R., Mehta, Y., Lillicrap, T. & Latham, P. E. (2021), ‘Towards biologically plausible convolutional networks’, *Advances in Neural Information Processing Systems* **34**, 13924–13936.
- Pulvermüller, F., Tomasello, R., Henningsen-Schomers, M. R. & Wennekers, T. (2021), ‘Biological constraints on neural network models of cognitive function’, *Nature Reviews Neuroscience* **22**(8), 488–502.
- Rojas, R. & Rojas, R. (1996), ‘The backpropagation algorithm’, *Neural networks: a systematic introduction* pp. 149–182.
- Ruder, S. (2016), ‘An overview of gradient descent optimization algorithms’, *arXiv preprint arXiv:1609.04747*.
- Rumelhart, D. E., Hinton, G. E. & Williams, R. J. (1986), ‘Learning representations by back-propagating errors’, *nature* **323**(6088), 533–536.
- Rumelhart, D. E. & Zipser, D. (1985), ‘Feature discovery by competitive learning’, *Cognitive science* **9**(1), 75–112.
- Schuman, C. D., Kulkarni, S. R., Parsa, M., Mitchell, J. P., Kay, B. et al. (2022), ‘Opportunities for neuromorphic computing algorithms and applications’, *Nature Computational Science* **2**(1), 10–19.
- Sejnowski, T. J. & Tesauro, G. (1989), ‘Building network learning algorithms from hebbian synapses’, *Brain organization and memory: cells, systems, and circuits* pp. 338–355.
- Shepherd, G. M. & Grillner, S. (2018), *Handbook of brain microcircuits*, Oxford University Press.
- Song, Y., Lukasiewicz, T., Xu, Z. & Bogacz, R. (2020), ‘Can the brain do backpropagation?—exact implementation of backpropagation in predictive coding networks’, *Advances in neural information processing systems* **33**, 22566–22579.
- Spoerer, C. J., McClure, P. & Kriegeskorte, N. (2017), ‘Recurrent convolutional neural networks: a better model of biological object recognition’, *Frontiers in psychology* **8**, 1551.
- Sultonov, F., Park, J.-H., Yun, S., Lim, D.-W. & Kang, J.-M. (2022), ‘Mixer u-net: An improved automatic road extraction from uav imagery’, *Applied Sciences* **12**, 1953.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V. & Rabinovich, A. (2015), Going deeper with convolutions, in ‘Proceedings of the IEEE conference on computer vision and pattern recognition’, pp. 1–9.
- Wu, J. (2017), ‘Introduction to convolutional neural networks’, *National Key Lab for Novel Software Technology. Nanjing University. China* **5**(23), 495.

A Appendix

A.1 Further Results

Additional results for the backpropagation models, both trained using our weight initialisation scheme or the default PyTorch scheme. Similar results to our Hebbian configurations were obtained when trained on the same initial weights, while the default Pytorch weights combined colours and shapes together more effectively. Class separation was also similar, displaying no meaningful differences with our Hebbian methodology. Figures 27-29 all belong to backpropagation experiments.

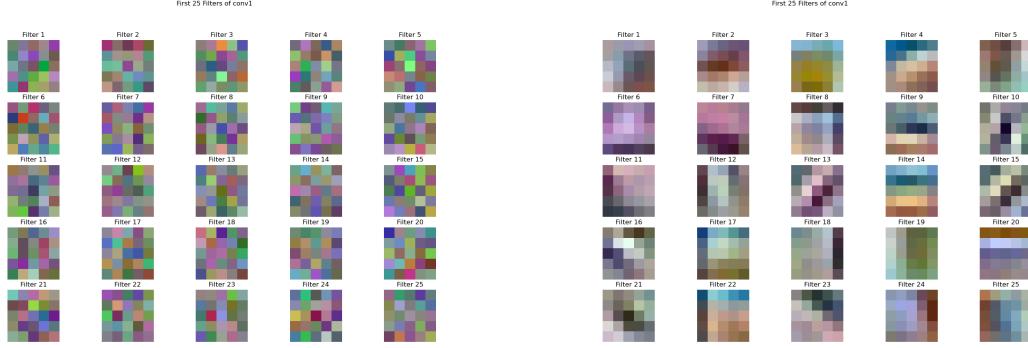


Figure 27: Direct filter plots of kernels in first layer of **SoftHebb-Backpropagation** model, using the same weight initialisation as Hebbian experiments or default Pytorch initialisation.

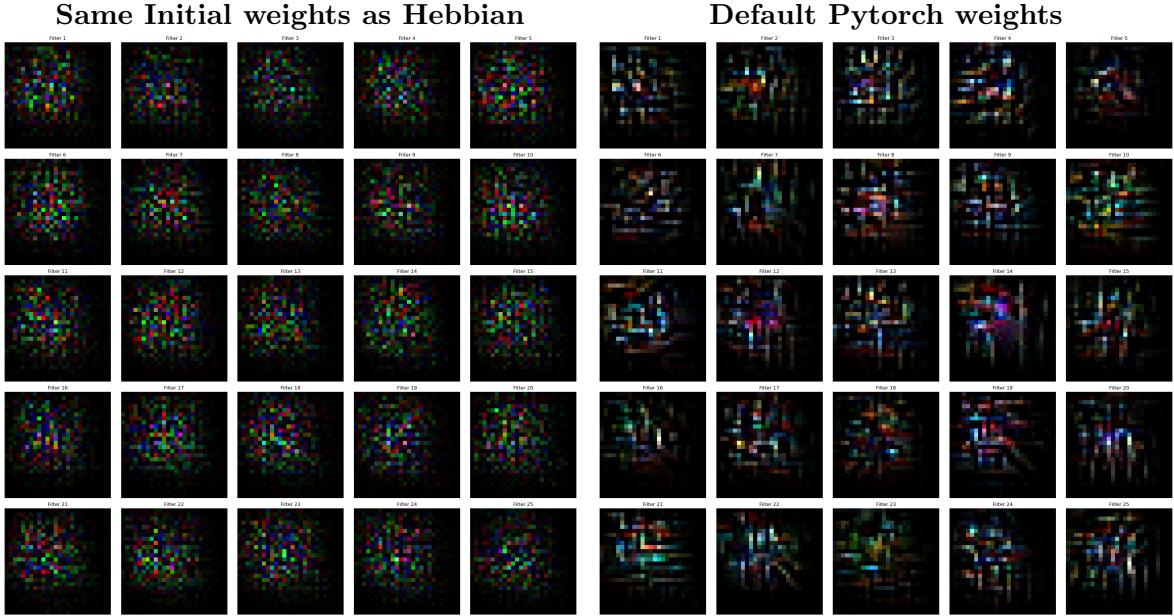


Figure 28: Receptive fields of kernels in last layer of **SoftHebb-Backpropagation** model, using the same weight initialisation as Hebbian experiments or default Pytorch initialisation.

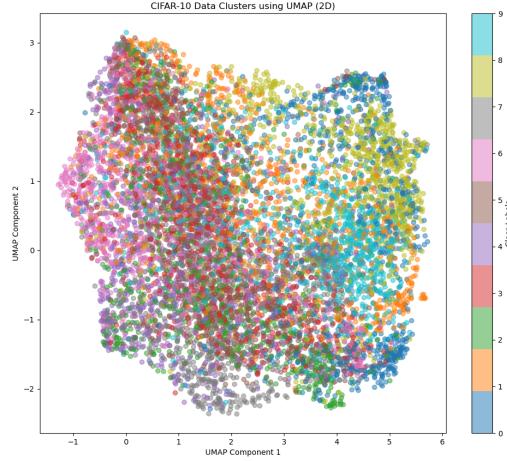


Figure 29: Plots of class clustering/separation of **SoftHebb-Backpropagation** configuration, with same initial weights as Hebbian configurations.

Results using the **Lagani-HardWTA/Cos-Instar** configuration are provided in Figure 30, to showcase how only simple patterns and shapes were formed when using a smaller architecture with smaller receptive fields and less filters per layer.

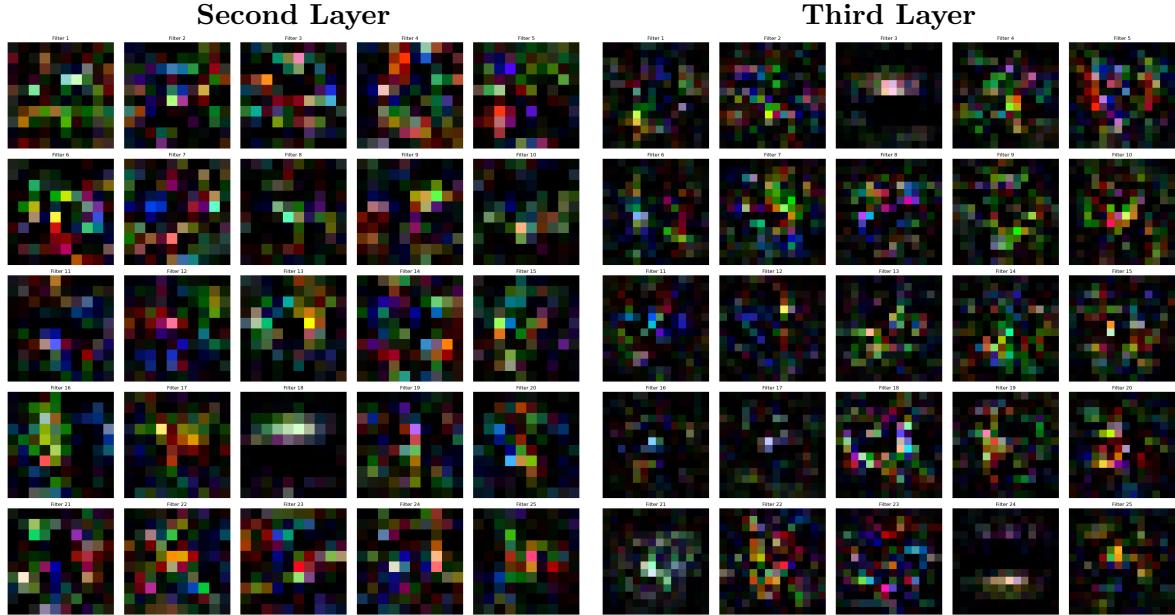


Figure 30: Receptive fields of neurons in second and last layer of **Lagani-HardWTA-Instar** configuration.

Supplementary results to aid our discussion have been provided for the residual and presynaptic configurations. Figure 31 displays how Residual depthwise convolutions focused on colours rather than edge orientations at its deeper layers, unlike the depthwise architecture. Figure 32 presents results from the presynaptic competition at a global and Output level. Global offers only noisy random representations, while output competition displays repeated focused patterns. Neither are informative for forming hierarchical representations.

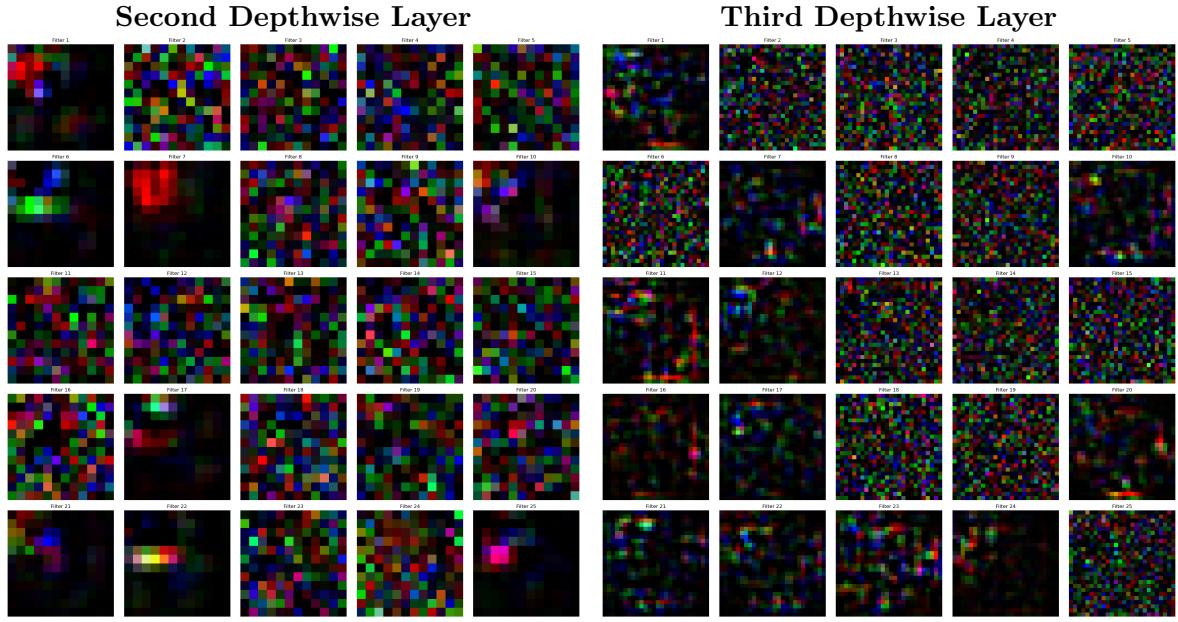


Figure 31: Receptive fields plots of kernels in second and last layer of **Residual_SoftHebb-Surr/HardWTA/Cos-Instar** configuration.

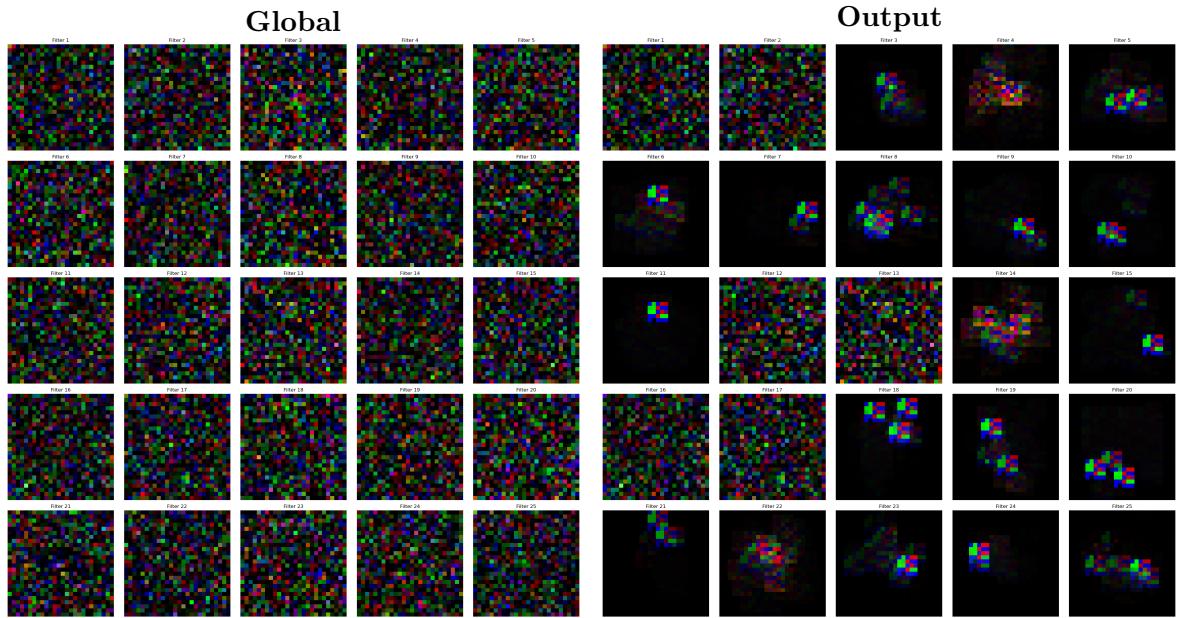


Figure 32: Receptive fields plots of kernels in last layer of configuration using Presynaptic Global or Output competition

A.2 Code Implementation

A.2.1 Data.py

```

1 import os
2 import torch
3 import torch.nn.functional as F
4 from torch.utils.data import DataLoader
5 from torchvision.datasets import CIFAR10, STL10
6 import torchvision.transforms as T

```

```

7 import matplotlib.pyplot as plt
8 import numpy as np
9
10 import params as P
11
12
13 class ZCAWhitening:
14     def __init__(self, epsilon=1e-1):
15         self.epsilon = epsilon
16         self.zca_matrix = None
17         self.mean = None
18         self.std = None
19
20     def fit(self, x: torch.Tensor, transpose=True, dataset: str = "CIFAR10"):
21         path = os.path.join("zca_data", dataset, f"{dataset}_zca.pt")
22         os.makedirs(os.path.dirname(path), exist_ok=True)
23
24     try:
25         saved_data = torch.load(path, map_location='cpu')
26         self.zca_matrix = saved_data['zca']
27         self.mean = saved_data['mean']
28         self.std = saved_data['std']
29         print(f"Loaded pre-computed ZCA matrix for {dataset}")
30     except FileNotFoundError:
31         print(f"Computing ZCA matrix for {dataset}")
32         if transpose and x.dim() == 4:
33             x = x.permute(0, 3, 1, 2)
34
35         x = x.reshape(x.shape[0], -1)
36         self.mean = x.mean(dim=0, keepdim=True)
37         self.std = x.std(dim=0, keepdim=True)
38         x = (x - self.mean) / (self.std + self.epsilon)
39
40         print(f"Data shape after reshaping: {x.shape}")
41         print(
42             f"Data statistics - Min: {x.min().item():.4f}, Max: {x.max().item():.4f}, Mean: {x.mean().item():.4f}, Std: {x.std().item():.4f}")
43
44         cov = torch.mm(x.T, x) / (x.shape[0] - 1)
45
46     # Debug print
47     print(f"Covariance matrix shape: {cov.shape}")
48     print(f"Covariance matrix statistics - Min: {cov.min().item():.4f}, Max: {cov.max().item():.4f}")
49
50     # Check for NaN or inf values
51     if torch.isnan(cov).any() or torch.isinf(cov).any():
52         print("Warning: NaN or inf values detected in covariance matrix")
53         cov = torch.where(torch.isnan(cov) | torch.isinf(cov), torch.zeros_like(cov), cov)
54
55     if dataset == "CIFAR10":
56         u, s, v = torch.svd(cov)

```

```

57         inv_sqrt_s = torch.diag(1.0 / torch.sqrt(s + self.
58                                     epsilon))
59         self.zca_matrix = torch.mm(torch.mm(u, inv_sqrt_s),
60                                     u.T)
61
60     else:
61         # Use eigendecomposition instead of SVD
62         eigenvalues, eigenvectors = torch.linalg.eigh(cov)
63         inv_sqrt_eigenvalues = torch.diag(1.0 / torch.sqrt(
64             eigenvalues + self.epsilon))
65         self.zca_matrix = torch.mm(torch.mm(eigenvectors,
66             inv_sqrt_eigenvalues), eigenvectors.T)
67
66     torch.save({'zca': self.zca_matrix, 'mean': self.mean,
67                 'std': self.std}, path)
68     print(f"Saved computed ZCA matrix for {dataset}")
69
69 def transform(self, x: torch.Tensor):
70     if self.zca_matrix is None:
71         raise ValueError("ZCA matrix not computed. Call fit()
72                         first.")
73
73     original_shape = x.shape
74     x = x.reshape(x.shape[0], -1)
75     x = (x - self.mean) / (self.std + self.epsilon)
76     x_whitened = torch.mm(x, self.zca_matrix)
77     return x_whitened.reshape(original_shape)
78
79
80 def whitening_zca(x: torch.Tensor, transpose=True, dataset: str = "CIFAR10"):
81     zca = ZCAWhitening()
82     zca.fit(x, transpose, dataset)
83     return zca.transform(x)
84
85
86 class ZCATransformation:
87     def __init__(self, zca):
88         self.zca = zca
89
90     def __call__(self, x):
91         if x.dim() == 3:
92             x = x.unsqueeze(0)
93         x_whitened = self.zca.transform(x)
94         return x_whitened.squeeze(0) if x_whitened.shape[0] == 1
95             else x_whitened
96
96
97 def visualize_zca_effect(original_data, whitened_data, num_samples
98 =5):
99     fig, axes = plt.subplots(2, num_samples, figsize=(15, 6))
100    for i in range(num_samples):
101        # Normalize data to [0, 1] range
102        orig = original_data[i].permute(1, 2, 0)
103        orig = (orig - orig.min()) / (orig.max() - orig.min())
104
104        whit = whitened_data[i].permute(1, 2, 0)
105        whit = (whit - whit.min()) / (whit.max() - whit.min())

```

```

106     axes[0, i].imshow(orig)
107     axes[0, i].axis('off')
108     axes[0, i].set_title('Original')
109
110     axes[1, i].imshow(whit)
111     axes[1, i].axis('off')
112     axes[1, i].set_title('Whitened')
113
114 plt.tight_layout()
115 plt.show()
116
117
118
119 def check_covariance(data):
120     data_flat = data.reshape(data.shape[0], -1)
121     cov_matrix = torch.cov(data_flat.T)
122     diag_mean = cov_matrix.diag().mean().item()
123     off_diag_mean = (cov_matrix - torch.diag(cov_matrix.diag())).abs_()
124     print(f"Mean of diagonal elements: {diag_mean:.6f}")
125     print(f"Mean of off-diagonal elements: {off_diag_mean:.6f}")
126     print(f"Ratio of off-diagonal to diagonal: {off_diag_mean / diag_mean:.6f}")
127
128 def check_normalization(data):
129     data_flat = data.reshape(data.shape[0], -1)
130     mean = data_flat.mean().item()
131     std = data_flat.std().item()
132     print(f"Mean: {mean:.6f}")
133     print(f"Standard deviation: {std:.6f}")
134     print(f"Min: {data_flat.min().item():.6f}")
135     print(f"Max: {data_flat.max().item():.6f}")
136     print()
137
138
139 def get_data(dataset='cifar10', root='datasets', batch_size=32,
140             num_workers=0, whiten_lvl=None):
141     trn_set, tst_set = None, None
142     if dataset == 'cifar10':
143         trn_set = CIFAR10(root=os.path.join(root, dataset), train=True,
144                           download=True, transform=T.ToTensor())
145         tst_set = CIFAR10(root=os.path.join(root, dataset), train=False,
146                           download=True, transform=T.ToTensor())
147         all_data = torch.cat([torch.tensor(trn_set.data), torch.
148                               tensor(tst_set.data)], dim=0)
149         all_data = all_data.float() / 255.0 # Normalize data to [0,
150                                           1]
151
152     elif dataset == 'stl10':
153         trn_set = STL10(root=os.path.join(root, dataset), split='train',
154                         download=True, transform=T.ToTensor())
155         tst_set = STL10(root=os.path.join(root, dataset), split='test',
156                         download=True, transform=T.ToTensor())
157         all_data = torch.cat([trn_sample[0].unsqueeze(0) for
158                               trn_sample in trn_set] +
159                               [tst_sample[0].unsqueeze(0) for
160                               tst_sample in tst_set], dim=0)
161
162 else:

```

```

153     raise NotImplementedError("Dataset {} not supported.".format
154         (dataset))
155
156     zca = None
157     if whiten_lvl is not None:
158         print("Data whitening")
159         zca = ZCAWhitening(epsilon=whiten_lvl)
160         zca.fit(all_data, transpose=False, dataset=dataset)
161
162     # Create a temporary loader with only ToTensor transform for
163     # visualization
164     if dataset=="cifar10":
165         temp_transform = T.Compose([T.Resize(32), T.ToTensor()])
166         temp_dataset = CIFAR10(root=os.path.join(root, dataset),
167                               train=True, download=False,
168                               transform=temp_transform)
169
170         full_transform = T.Compose([
171             T.RandomHorizontalFlip(),
172             T.Resize(32),
173             T.ToTensor(),
174             ZCATransformation(zca),
175         ])
176     else:
177         temp_transform = T.Compose([T.Resize(96), T.ToTensor()])
178         temp_dataset = STL10(root=os.path.join(root, dataset),
179                               split='train', download=False,
180                               transform=temp_transform)
181
182         full_transform = T.Compose([
183             T.RandomHorizontalFlip(),
184             T.Resize(96),
185             T.ToTensor(),
186             ZCATransformation(zca),
187         ])
188
189     # temp_loader = DataLoader(temp_dataset, batch_size=
190     #     batch_size, shuffle=True, num_workers=num_workers)
191
192     # # Get a batch for visualization
193     # original_batch, _ = next(iter(temp_loader))
194     # whitened_batch = zca.transform(original_batch)
195     # print("\nVisualization of ZCA effect:")
196     # visualize_zca_effect(original_batch, whitened_batch)
197     # print("\nCovariance check:")
198     # print("Before whitening:")
199     # check_covariance(original_batch)
200     # print("\nAfter whitening:")
201     # check_covariance(whitened_batch)
202
203     #
204     # print("\nNormalization check:")
205     # print("Before whitening:")
206     # check_normalization(original_batch)
207     # print("\nAfter whitening:")
208     # check_normalization(whitened_batch)
209
210
211     # Now apply the full transform including ZCA to the datasets
212
213
214     trn_set.transform = full_transform
215     tst_set.transform = full_transform
216
217     else:

```

```

206     print("No ZCA")
207     if dataset == 'cifar10':
208         mean = [0.4914, 0.4822, 0.4465]
209         std = [0.2023, 0.1994, 0.2010]
210         temp_transform = T.Compose([
211             T.RandomHorizontalFlip(),
212             T.Resize(32),
213             T.ToTensor(),
214             T.Normalize(mean=mean, std=std)])
215     elif dataset == 'stl10':
216         mean = [0.4467, 0.4398, 0.4066]
217         std = [0.2603, 0.2566, 0.2713]
218         temp_transform = T.Compose([
219             T.RandomHorizontalFlip(),
220             T.Resize(96),
221             T.ToTensor(),
222             T.Normalize(mean=mean, std=std)])
223     trn_set.transform = temp_transform
224     tst_set.transform = temp_transform
225     trn_loader = DataLoader(trn_set, batch_size=batch_size, shuffle=
226                             False, num_workers=P.NUM_WORKERS)
227     tst_loader = DataLoader(tst_set, batch_size=batch_size, shuffle=
228                             False, num_workers=P.NUM_WORKERS)
229
230
231 def visualize_whitening_effect(dataset, zca, num_images=5):
232     """
233     Visualize the effect of ZCA whitening on a few sample images.
234
235     :param dataset: The dataset to sample images from
236     :param zca: The ZCAWhitening object
237     :param num_images: Number of images to visualize
238     """
239
240     # Get a few sample images
241     sampler = torch.utils.data.RandomSampler(dataset, num_samples=
242                                               num_images)
243     loader = DataLoader(dataset, batch_size=1, sampler=sampler)
244
245     fig, axes = plt.subplots(num_images, 2, figsize=(10, 5 *
246                                                 num_images))
247     fig.suptitle("ZCA Whitening Effect", fontsize=16)
248
249     for i, (img, _) in enumerate(loader):
250         # Original image
251         orig_img = img.squeeze().permute(1, 2, 0).numpy()
252         axes[i, 0].imshow(orig_img)
253         axes[i, 0].set_title(f"Original Image {i + 1}")
254         axes[i, 0].axis('off')
255
256         # Whitened image
257         whitened_img = zca.transform(img).squeeze().permute(1, 2, 0)
258             .numpy()
259         axes[i, 1].imshow(whitened_img)
260         axes[i, 1].set_title(f"Whitened Image {i + 1}")
261         axes[i, 1].axis('off')

```

```

259     plt.tight_layout()
260     plt.show()
261
262 # Example usage:
263 if __name__ == "__main__":
264     # Use CIFAR-10
265     cifar_trn_loader, cifar_tst_loader, cifar_zca = get_data(dataset=
266         ='cifar10', batch_size=64, whiten_lvl=1e-3)
267     print("CIFAR-10 data loaders and ZCA whitening prepared.")
268
269
270 # Use STL-10
271 stl_trn_loader, stl_tst_loader, stl_zca = get_data(dataset='
272     stl10', batch_size=64, whiten_lvl=1e-3)
273     print("STL-10 data loaders and ZCA whitening prepared.")

```

Listing 1: Centering the Data

A.2.2 Hebb.py

```

1 import math
2
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 from torch.nn.modules.utils import _pair
7 import matplotlib.pyplot as plt
8
9 # Code uses elements from {https://github.com/GabrieleLagani/
10 # hebbdemo} and {
11 # https://github.com/NeuromorphicComputing/SoftHebb}
12 import torch.nn.init as init
13
14 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
15 torch.manual_seed(0)
16
17
18 def normalize(x, dim=None):
19     nrm = (x ** 2).sum(dim=dim, keepdim=True) ** 0.5
20     nrm[nrm == 0] = 1.
21     return x / nrm
22
23
24 def symmetric_pad(x, padding):
25     if padding == 0:
26         return x
27     return F.pad(x, (padding,) * 4, mode='reflect')
28
29
30 def create_sm_kernel(kernel_size=5, sigma_e=1.2, sigma_i=1.4):
31 """
32     Create a surround modulation kernel.
33     :param kernel_size: Size of the SM kernel.
34     :param sigma_e: Standard deviation for the excitatory Gaussian.
35     :param sigma_i: Standard deviation for the inhibitory Gaussian.

```

```

36     :return: A normalized SM kernel.
37     """
38     center = kernel_size // 2
39     x, y = torch.meshgrid(torch.arange(kernel_size), torch.arange(
40         kernel_size), indexing="ij")
41     x = x.float() - center
42     y = y.float() - center
43     # Compute the excitatory and inhibitory Gaussians
44     gaussian_e = torch.exp(-(x ** 2 + y ** 2) / (2 * sigma_e ** 2))
45     gaussian_i = torch.exp(-(x ** 2 + y ** 2) / (2 * sigma_i ** 2))
46     # Compute the Difference of Gaussians (DoG)
47     dog = gaussian_e / (2 * math.pi * sigma_e ** 2) - gaussian_i / (
48         2 * math.pi * sigma_i ** 2)
49     # Normalize the DoG so that the center value is 1
50     sm_kernel = dog / dog[center, center]
51     return sm_kernel.unsqueeze(0).unsqueeze(0).to(device)

52 class HebbianConv2d(nn.Module):
53     """
54     A 2d convolutional layer that learns through Hebbian plasticity
55     """
56
57     MODE_HPCA = 'h pca'
58     MODE_BASIC_HEBBIAN = 'basic'
59     MODE_WTA = 'wta'
60     MODE_SOFTWTA = 'soft'
61     MODE_BCM = 'bcm'
62     MODE_HARDWT = "hard"
63     MODE_PRESYNAPTIC_COMPETITION = "pre"
64     MODE_TEMPORAL_COMPETITION = "temp"
65     MODE_ADAPTIVE_THRESHOLD = "thresh"
66     MODE_ANTIHARDWT = "antihard"

67     def __init__(self, in_channels, out_channels, kernel_size,
68                  stride=1, dilation=1, padding=0, groups=1,
69                  w_nrm=False, bias=False, act=nn.Identity(),
70                  mode=MODE_SOFTWTA, k=1, patchwise=True,
71                  contrast=1., uniformity=False, alpha=1.,
72                  wta_competition='similarity_spatial',
73                  lateral_competition="combined",
74                  lateral_inhibition_strength=0.01, top_k=1,
75                  prune_rate=0, t_invert=1.):
76         """
77             :param out_channels: output channels of the convolutional
78                 kernel
79             :param in_channels: input channels of the convolutional
80                 kernel
81             :param kernel_size: size of the convolutional kernel (int or
82                 tuple)
83             :param stride: stride of the convolutional kernel (int or
84                 tuple)
85             :param w_nrm: whether to normalize the weight vectors before
86                 computing outputs
87             :param act: the nonlinear activation function after
88                 convolution
89             :param mode: the learning mode, either 'swta' or 'h pca'
```

```

83     :param k: softmax inverse temperature parameter used for
84         swta-type learning
85     :param patchwise: whether updates for each convolutional
86         patch should be computed separately,
87         and then aggregated
88     :param contrast: coefficient that rescales negative compared
89         to positive updates in contrastive-type learning
90     :param uniformity: whether to use uniformity weighting in
91         contrastive-type learning.
92     :param alpha: weighting coefficient between hebbian and
93         backprop updates (0 means fully backprop, 1 means fully
94         hebbian).
95     """
96
97     super(HebbianConv2d, self).__init__()
98     self.mode = mode
99     self.out_channels = out_channels
100    self.in_channels = in_channels
101    self.kernel = kernel_size
102    self.kernel_size = _pair(kernel_size)
103    self.stride = _pair(stride)
104    self.dilation = _pair(dilation)
105    self.padding = padding
106    self.padding_mode = 'reflect'
107    if mode == "hard":
108        self.padding_mode = 'symmetric'
109
110    self.F_padding = (padding, padding, padding, padding)
111    self.groups = 1 # in_channels for depthwise
112
113    weight_range = 25 / math.sqrt(in_channels * kernel_size *
114        kernel_size)
115    self.weight = nn.Parameter(
116        weight_range * torch.randn((out_channels, in_channels // 
117            self.groups, *self.kernel_size)))
118
119    # Different initialisation schemes for future work
120
121    # self.weight = nn.Parameter(torch.empty(out_channels,
122        in_channels // self.groups, *self.kernel_size))
123    # init.kaiming_uniform_(self.weight)
124    # self.weight = center_surround_init(out_channels,
125        in_channels, kernel_size, 1)
126
127    print(self.weight.shape)
128    self.w_nrm = w_nrm
129    self.act = act
130
131    # BCM parameters
132    if mode == "bcm":
133        self.theta_decay = 0.5
134        self.theta = nn.Parameter(torch.ones(out_channels),
135            requires_grad=False)
136
137        self.register_buffer('delta_w', torch.zeros_like(self.weight
138            ))
139        self.top_k = top_k
140        self.patchwise = patchwise
141        self.contrast = contrast

```

```

129     self.uniformity = uniformity
130     self.alpha = alpha
131     self.lebesgue_p = 2
132     self.prune_rate = prune_rate # 99% of connections are
133         pruned
134     self.t_invert = torch.tensor(t_invert)
135
136     #For presynaptic competition
137     self.presynaptic_competition_type = "softmax"
138     self.presynaptic_weights = False # presynaptic competition
139         in forward pass
140
141     # For temporal and statistical thresholds
142     self.activation_history = None
143     self.temporal_window = 500
144     self.competition_k = 2
145     self.competition_type = "hard"
146
147     # For surround lateral inhibition kernel
148     if self.kernel != 1:
149         self.sm_kernel = create_sm_kernel()
150         self.register_buffer('surround_kernel', self.sm_kernel)
151         # Visualise what the kernel will look like
152         self.visualize_surround_modulation_kernel()
153
154     # For direct homeostasis effect
155     self.target_activity = 0.1
156     self.scaling_rate = 0.001
157     self.register_buffer('average_activity', torch.zeros(
158         out_channels))
159
160     self.growth_probability = 0.1
161     self.new_synapse_strength = 1
162     self.prune_threshold_percentile = 10 # Prune bottom 10% of
163         weights
164
165     def visualize_surround_modulation_kernel(self):
166         """
167             Visualizes the surround modulation kernel using matplotlib.
168         """
169         sm_kernel = self.sm_kernel.squeeze().cpu().detach().numpy()
170         plt.figure(figsize=(5, 5))
171         plt.imshow(sm_kernel, cmap='jet')
172         plt.colorbar()
173         plt.title('Surround Modulation Kernel')
174         plt.show()
175
176     def apply_lebesgue_norm(self, w):
177         """Apply Lebesgue norm to weights."""
178         return torch.sign(w) * torch.abs(w) ** (self.lebesgue_p - 1)
179
180     def apply_weights(self, x, w):
181         """
182             Apply convolutional operation with weights to input.
183         """
184         return F.conv2d(x, w, None, self.stride, 0, self.dilation,
185             groups=self.groups)

```

```

182     def update_average_activity(self, y):
183         """Update the average activity of neurons."""
184         current_activity = y.mean(dim=(0, 2, 3))
185         self.average_activity = 0.9 * self.average_activity + 0.1 *
186             current_activity
187
188     def synaptic_scaling(self):
189         """Apply synaptic scaling to maintain target activity."""
190         scale_factor = self.target_activity / (self.average_activity
191             + 1e-6)
192         self.weight.data *= (1 + self.scaling_rate * (scale_factor -
193             1)).view(-1, 1, 1, 1)
194
195     def structural_plasticity(self):
196         """Prune sparse weights and create new ones."""
197         with torch.no_grad():
198             # Pruning step
199             prune_threshold = torch.quantile(torch.abs(self.weight),
200                 self.prune_threshold_percentile / 100)
201             weak_synapses = torch.abs(self.weight) < prune_threshold
202             self.weight.data[weak_synapses] = 0
203             # Growth step
204             zero_weights = self.weight.data == 0
205             new_synapses = torch.rand_like(self.weight) < self.
206                 growth_probability
207             new_synapses &= zero_weights
208             self.weight.data[new_synapses] = torch.randn_like(self.
209                 weight)[new_synapses] * self.new_synapse_strength
210
211     def cosine(self, x, w):
212         """Compute cosine similarity between input and weights."""
213         w_normalized = F.normalize(w, p=2, dim=1)
214         # conv_output = symmetric_pad(x, self.padding)
215         conv_output = F.conv2d(x, w_normalized, None, self.stride,
216             0, self.dilation, groups=self.groups)
217         x_squared = x.pow(2)
218         x_squared_sum = F.conv2d(x_squared, torch.ones_like(w), None
219             , self.stride, 0, self.dilation,
220                 self.groups)
221         x_norm = torch.sqrt(x_squared_sum + 1e-8)
222         cosine_sim = conv_output / x_norm
223         return cosine_sim
224
225     def apply_surround_modulation(self, y):
226         """Apply surround modulation to the output."""
227         return F.conv2d(y, self.sm_kernel.repeat(self.out_channels,
228             1, 1, 1),
229                         padding=self.sm_kernel.size(-1) // 2, groups
230                         =self.out_channels)
231
232     def compute_activation(self, x):
233         """Compute the activation of the layer."""
234         x = symmetric_pad(x, self.padding)
235         w = self.weight
236         if self.w_nrm: w = normalize(w, dim=(1, 2, 3))
237         if self.presynaptic_weights: w = self.
238             compute_presynaptic_competition_global(w)
239         # No activation required:

```

```

229     # y = self.act(self.apply_weights(x, w))
230     # For cosine similarity activation:
231     # Must manually change to allow/disable cosine similarity
232     y = self.cosine(x, w)
233     return x, y, w
234
235 def forward(self, x):
236     x,y, w = self.compute_activation(x)
237     # Only apply lateral inhibition if spatial dimensions exists
238     # Must manually change to allow/disable surround inhibition
239     if self.kernel != 1:
240         y = self.apply_surround_modulation(y)
241     if self.training:
242         # Optional mechanisms
243         # self.update_average_activity(y)
244         # self.synaptic_scaling()
245         self.compute_update(x, y, w)
246     return y
247
248 def compute_update(self, x, y, weight):
249     """Compute weight updates based on the chosen learning mode.
250     """
251     if self.mode == self.MODE_BASIC_HEBBIAN:
252         update = self.update_basic_hebbian(x, y, weight)
253     elif self.mode == self.MODE_HARDWT:
254         update = self.update_hardwt(x, y, weight)
255     elif self.mode == self.MODE_SOFTWTA:
256         update = self.update_softwta(x, y, weight)
257     elif self.mode == self.MODE_ANTIHARDWT:
258         update = self.update_antihardwt(x, y, weight)
259     elif self.mode == self.MODE_BCM:
260         update = self.update_bcm(x, y, weight)
261     elif self.mode == self.MODE_TEMPORAL_COMPETITION:
262         update = self.update_temporal_competition(x, y, weight)
263     elif self.mode == self.MODE_ADAPTIVE_THRESHOLD:
264         update = self.update_adaptive_threshold(x, y, weight)
265     else:
266         raise NotImplementedError(f"Learning mode {self.mode} "
267             f"unavailable for {self.__class__.__name__} layer")
268     # Weight Normalization and added to weight change buffer
269     update.div_(torch.abs(update).amax() + 1e-30)
270     self.delta_w += update
271
272 def update_basic_hebbian(self, x, y, weight):
273     """Implement basic Hebbian learning (Grossberg Instar rule).
274     """
275     yx = self.compute_yx(x, y)
276     y_sum = y.sum(dim=(0, 2, 3)).view(-1, 1, 1, 1)
277     yw = y_sum * weight
278     update = yx - yw
279     return update
280
281 def update_hardwt(self, x, y, weight):
282     """Implement hard Winner-Take-All (WTA) Hebbian learning."""
283     y_wta = y * self.compute_wta_mask(y)
284     yx = self.compute_yx(x, y_wta)
285     yu = torch.sum(y_wta, dim=(0, 2, 3))
286     update = yx - yu.view(-1, 1, 1, 1) * weight

```

```

284     return update
285
286     def update_softwta(self, x, y, weight):
287         """Implement soft Winner-Take-All Hebbian learning (SoftHebb)
288         """
289         softwta_activs = self.compute_softwta_activations(y)
290         yx = self.compute_yx(x, softwta_activs)
291         yu = torch.sum(torch.mul(softwta_activs, y), dim=(0, 2, 3))
292         update = yx - yu.view(-1, 1, 1, 1) * weight
293         return update
294
295     def update_antihardwt(self, x, y, weight):
296         """Implement anti-Hebbian hard Winner-Take-All learning."""
297         hardwta_activs = self.compute_antihardwta_activations(y)
298         yx = self.compute_yx(x, hardwta_activs)
299         yu = torch.sum(torch.mul(hardwta_activs, y), dim=(0, 2, 3))
300         update = yx - yu.view(-1, 1, 1, 1) * weight
301         return update
302
303     def update_bcm(self, x, y, weight):
304         """Implement BCM (Bienenstock-Cooper-Munro) learning rule.
305         """
306         y_wta = y * self.compute_wta_mask(y)
307         y_squared = y_wta.pow(2).mean(dim=(0, 2, 3))
308         self.theta.data = (1 - self.theta_decay) * self.theta + self
309             .theta_decay * y_squared
310         y_minus_theta = y_wta - self.theta.view(1, -1, 1, 1)
311         bcm_factor = y_wta * y_minus_theta
312         yx = self.compute_yx(x, bcm_factor)
313         update = yx.view(weight.shape)
314         return update
315
316     def update_temporal_competition(self, x, y, weight):
317         """Implement temporal competition-based learning."""
318         self.update_activation_history(y)
319         temporal_winners = self.compute_temporal_winners(y)
320         y_winners = temporal_winners * y
321         y_winners = y_winners * self.apply_competition(y_winners)
322         yx = self.compute_yx(x, y_winners)
323         y_sum = y_winners.sum(dim=(0, 2, 3)).view(-1, 1, 1, 1)
324         update = yx - y_sum * weight
325         return update
326
327     def update_adaptive_threshold(self, x, y, weight):
328         """Implement adaptive/statistics threshold-based learning.
329         """
330         batch_size, out_channels, height_out, width_out = y.shape
331         similarities = F.conv2d(x, weight, stride=self.stride,
332             padding=self.padding, groups=self.groups)
333         similarities = similarities / (torch.norm(weight.view(
334             out_channels, -1), dim=1).view(1, -1, 1, 1) + 1e-10)
335         threshold = self.compute_adaptive_threshold(similarities)
336         winners = (similarities > threshold).float()
337         y_winners = winners * similarities
338         y_winners = y_winners * self.apply_competition(y_winners)
339         yx = self.compute_yx(x, y_winners)
340         y_sum = y_winners.sum(dim=(0, 2, 3)).view(-1, 1, 1, 1)
341         update = yx - y_sum * weight

```

```

336         return update
337
338     def update_activation_history(self, y):
339         """Update the activation history for temporal competition.
340         """
341         if self.activation_history is None:
342             self.activation_history = y.detach().clone()
343         else:
344             self.activation_history = torch.cat([self.
345                 activation_history, y.detach()], dim=0)
346         if self.activation_history.size(0)>self.temporal_window:
347             self.activation_history = self.activation_history[-
348                 self.temporal_window:]
349
350     def compute_temporal_winners(self, y):
351         """Select winners of temporal competition."""
352         batch_size, out_channels, height_out, width_out = y.shape
353         history_spatial = self.activation_history.view(-1,
354             out_channels, height_out, width_out)
355         median_activations = torch.median(history_spatial, dim=0)[0]
356         # Compute threshold for each spatial location
357         temporal_threshold = torch.mean(median_activations, dim=0,
358             keepdim=True)
359         return (y > temporal_threshold).float()
360
361     def compute_adaptive_threshold(self, similarities):
362         """Create threshold for statistical thresholding."""
363         mean_sim = similarities.mean(dim=1, keepdim=True)
364         std_sim = similarities.std(dim=1, keepdim=True)
365         return mean_sim + self.competition_k * std_sim
366
367     def apply_competition(self, y):
368         # Hard-WTA or Soft-WTA competition mask for temporal/
369         # statistical thresholding
370         batch_size, out_channels, height, width = y.shape
371         if self.mode in [self.MODE_TEMPORAL_COMPETITION, self.
372             MODE_ADAPTIVE_THRESHOLD]:
373             if self.competition_type == 'hard':
374                 y = y.view(batch_size, out_channels, -1)
375                 top_k_values, top_k_indices = torch.topk(y, self.
376                     top_k, dim=1, largest=True, sorted=False)
377                 y_compete = torch.zeros_like(y)
378                 y_compete.scatter_(1, top_k_indices, top_k_values)
379                 return y_compete.view(batch_size, out_channels,
380                     height, width)
381             elif self.competition_type == 'soft':
382                 y_flat = y.view(batch_size, out_channels, -1)
383                 y_soft = torch.softmax(self.t_invert * y_flat, dim
384                     =1)
385                 return y_soft.view(batch_size, out_channels, height,
386                     width)
387         return y
388
389     def compute_yx(self, x, y):
390         """Compute yx term from Grossberg Instar rule"""
391         yx = F.conv2d(x.transpose(0, 1), y.transpose(0, 1), padding
392             =0,

```

```

381                     stride=self.dilation, dilation=self.stride).
382                     transpose(0, 1)
383             if self.groups != 1:
384                 yx = yx.mean(dim=1, keepdim=True)
385             return yx
386
387     def compute_wta_mask(self, y):
388         """Compute Hard-WTA mask."""
389         batch_size, out_channels, height_out, width_out = y.shape
390         y_flat = y.transpose(0, 1).reshape(out_channels, -1)
391         win_neurons = torch.argmax(y_flat, dim=0)
392         wta_mask = F.one_hot(win_neurons, num_classes=out_channels).
393                     float()
394         return wta_mask.transpose(0, 1).view(out_channels,
395                                             batch_size, height_out, width_out).transpose(0, 1)
396
397     def compute_softwta_activations(self, y):
398         """Computes SoftHebb mask and Hebb/AntiHebb implementation.
399         """
400         batch_size, out_channels, height_out, width_out = y.shape
401         flat_weighted_inputs = y.transpose(0, 1).reshape(
402             out_channels, -1)
403         flat_softwta_activs = torch.softmax(self.t_invert *
404                                              flat_weighted_inputs, dim=0)
405         flat_softwta_activs = -flat_softwta_activs
406         win_neurons = torch.argmax(flat_weighted_inputs, dim=0)
407         competing_idx = torch.arange(flat_weighted_inputs.size(1))
408         flat_softwta_activs[win_neurons, competing_idx] = -
409             flat_softwta_activs[win_neurons, competing_idx]
410         return flat_softwta_activs.view(out_channels, batch_size,
411                                         height_out, width_out).transpose(0, 1)
412
413     def compute_antihardwta_activations(self, y):
414         # Computes Hard-WTA mask and Hebb/AntiHebb implementation.
415         batch_size, out_channels, height_out, width_out = y.shape
416         flat_weighted_inputs = y.transpose(0, 1).reshape(
417             out_channels, -1)
418         win_neurons = torch.argmax(flat_weighted_inputs, dim=0)
419         competing_idx = torch.arange(flat_weighted_inputs.size(1))
420         anti_hebbian_mask = torch.ones_like(flat_weighted_inputs) *
421             -1
422         anti_hebbian_mask[win_neurons, competing_idx] = 1
423         flat_hardwta_activs = flat_weighted_inputs *
424             anti_hebbian_mask
425         return flat_hardwta_activs.view(out_channels, batch_size,
426                                         height_out, width_out).transpose(0, 1)
427
428     def compute_presynaptic_competition(self, m):
429         # This presynaptic competition promotes diversity among
430         # output channels.
431         m = 1 / (torch.abs(m) + 1e-6)
432         if self.presynaptic_competition_type == 'linear':
433             return m / (m.sum(dim=0, keepdim=True) + 1e-6)
434         elif self.presynaptic_competition_type == 'softmax':
435             return F.softmax(m, dim=0)
436         elif self.presynaptic_competition_type == 'lp_norm':
437             return F.normalize(m, p=2, dim=0)
438         else:

```

```

426     raise ValueError(f"Unknown competition type: {self.
427                     competition_type}")
428
429     def compute_presynaptic_competition_spatial(self, m):
430         # The presynaptic spatial competition encourages each input-
431         # output channel pair.
432         m = 1 / (torch.abs(m) + 1e-6)
433         if self.presynaptic_competition_type == 'linear':
434             # Sum across spatial dimensions (last two dimensions)
435             return m / (m.sum(dim=(-2, -1), keepdim=True) + 1e-6)
436         elif self.presynaptic_competition_type == 'softmax':
437             # Reshape to combine spatial dimensions
438             shape = m.shape
439             m_flat = m.view(*shape[:-2], -1)
440             # Apply softmax across spatial dimensions
441             m_comp = F.softmax(m_flat, dim=-1)
442             # Reshape back to original shape
443             return m_comp.view(*shape)
444         elif self.presynaptic_competition_type == 'lp_norm':
445             # Normalize across spatial dimensions
446             return F.normalize(m, p=2, dim=(-2, -1))
447         else:
448             raise ValueError(f"Unknown competition type: {self.
449                             presynaptic_competition_type}")
450
451     def compute_presynaptic_competition_input_channels(self, m):
452         # The presynaptic input competition promotes specialisation
453         # of each output channel across input features.
454         m = 1 / (torch.abs(m) + 1e-6)
455         if self.presynaptic_competition_type == 'linear':
456             # Sum across input channel dimension
457             return m / (m.sum(dim=1, keepdim=True) + 1e-6)
458         elif self.presynaptic_competition_type == 'softmax':
459             # Apply softmax across input channels
460             return F.softmax(m, dim=1)
461         elif self.presynaptic_competition_type == 'lp_norm':
462             # Normalize across input channels
463             return F.normalize(m, p=2, dim=1)
464         else:
465             raise ValueError(f"Unknown competition type: {self.
466                             presynaptic_competition_type}")
467
468     def compute_presynaptic_competition_global(self, m):
469         # The global competition creates a more intense competition
470         # where every weight competes with all others,
471         # potentially leading to very sparse but highly specialised
472         # connections.
473         m = 1 / (torch.abs(m) + 1e-6)
474         if self.presynaptic_competition_type == 'linear':
475             # Global sum across all dimensions
476             return m / (m.sum() + 1e-6)
477         elif self.presynaptic_competition_type == 'softmax':
478             # Flatten and apply softmax globally
479             m_flat = m.view(-1)
480             return F.softmax(m_flat, dim=0).view(m.shape)
481         elif self.presynaptic_competition_type == 'lp_norm':
482             # Global normalization
483             return F.normalize(m.view(-1), p=2).view(m.shape)

```

```

477         else:
478             raise ValueError(f"Unknown competition type: {self.
479                             presynaptic_competition_type}")
480
481     @torch.no_grad()
482     def local_update(self):
483         """
484             This function transfers a previously computed weight update,
485             stored in buffer self.delta_w, to the gradient
486             self.weight.grad of the weight parameter.
487
488             This function should be called before optimizer.step(), so
489             that the optimizer will use the locally computed
490             update as optimization direction. Local updates can also be
491             combined with end-to-end updates by calling this
492             function between loss.backward() and optimizer.step(). loss.
493             backward will store the end-to-end gradient in
494             self.weight.grad, and this function combines this value with
495             self.delta_w as
496             self.weight.grad = (1 - alpha) * self.weight.grad - alpha *
497             self.delta_w
498             Parameter alpha determines the scale of the local update
499             compared to the end-to-end gradient in the combination.
500             """
501
502     if self.weight.grad is None:
503         self.weight.grad = -self.alpha * self.delta_w
504     else:
505         self.weight.grad = (1 - self.alpha) * self.weight.grad -
506                             self.alpha * self.delta_w
507         self.delta_w.zero_()
508
509     # Update method to modify weights without requiring an optimiser
510     # @torch.no_grad()
511     # # Weight Update
512     # def local_update(self):
513     #     new_weight = self.weight + 0.1 * self.alpha * self.delta_w
514     #     # Update weights
515     #     self.weight.copy_(new_weight)
516     #     # self.structural_plasticity()
517     #     self.delta_w.zero_()

```

Listing 2: Calculation of post-synaptic activity in a CNN

A.2.3 Hebb_depthwise.py

```

1 import math
2
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 from torch.nn.modules.utils import _pair
7 import torch.nn.init as init
8
9
10 """
11 Uses almost identical code to hebb.py. Please refer to hebb.py for
12 additional explanations on code functionality

```

```

12 Only changes are the competition modes, which apply competition
    across spatial neurons in a filter for channel independence
13 """
14
15 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
16 torch.manual_seed(36)
17
18
19 def normalize(x, dim=None):
20     nrm = (x ** 2).sum(dim=dim, keepdim=True) ** 0.5
21     nrm[nrm == 0] = 1.
22     return x / nrm
23
24
25 def symmetric_pad(x, padding):
26     if padding == 0:
27         return x
28     return F.pad(x, (padding,) * 4, mode='reflect')
29
30
31 def center_surround_init(out_channels, in_channels, kernel_size,
32     groups=1):
33     # Calculate weight range
34     weight_range = 25 / math.sqrt(in_channels * kernel_size *
35         kernel_size)
36     # Calculate sigma based on kernel size (using equation 3 from
37     # the paper)
38     gamma = torch.empty(out_channels).uniform_(0, 0.5)
39     sigma = (kernel_size / 4) * torch.sqrt((1 - gamma ** 2) / (-
40         torch.log(gamma)))
41     # Create meshgrid for x and y coordinates
42     x = torch.linspace(-(kernel_size - 1) / 2, (kernel_size - 1) /
43         2, kernel_size)
44     y = torch.linspace(-(kernel_size - 1) / 2, (kernel_size - 1) /
45         2, kernel_size)
46     xx, yy = torch.meshgrid(x, y, indexing='ij')
47     # Calculate center and surround Gaussians
48     center = torch.exp(-(xx ** 2 + yy ** 2) / (2 * (gamma.view(-1,
49         1, 1) * sigma.view(-1, 1, 1)) ** 2))
50     surround = torch.exp(-(xx ** 2 + yy ** 2) / (2 * sigma.view(-1,
51         1, 1) ** 2))
52     # Calculate DoG (Difference of Gaussians)
53     dog = center - surround
54     # Normalize DoG
55     ac = torch.sum(torch.clamp(dog, min=0))
56     as_ = torch.sum(-torch.clamp(dog, max=0))
57     dog = weight_range * 0.5 * dog / (ac + as_)
58     # Assign excitatory (positive) or inhibitory (negative) centers
59     center_type = torch.cat([torch.ones(out_channels // 2), -torch.
60         ones(out_channels - out_channels // 2)])
61     center_type = center_type[torch.randperm(out_channels)].view(-1,
62         1, 1)
63     dog = dog * center_type
64     # Repeat for in_channels and reshape to match conv2d weight
65     # shape
66     dog = dog.unsqueeze(1).repeat(1, in_channels // groups, 1, 1)

```

```

56     dog = dog.reshape(out_channels, in_channels // groups,
57                         kernel_size, kernel_size)
58     return nn.Parameter(dog)
59
60 def create_sm_kernel(kernel_size=5, sigma_e=1.2, sigma_i=1.4):
61     center = kernel_size // 2
62     x, y = torch.meshgrid(torch.arange(kernel_size), torch.arange(
63                           kernel_size))
64     x = x.float() - center
65     y = y.float() - center
66     gaussian_e = torch.exp(-(x ** 2 + y ** 2) / (2 * sigma_e ** 2))
67     gaussian_i = torch.exp(-(x ** 2 + y ** 2) / (2 * sigma_i ** 2))
68     dog = gaussian_e / (2 * math.pi * sigma_e ** 2) - gaussian_i /
69         (2 * math.pi * sigma_i ** 2)
70     sm_kernel = dog / dog[center, center]
71     return sm_kernel.unsqueeze(0).unsqueeze(0).to(device)
72
73 # Doubts:
74 # Visualizing weights, as separated between channels and spatial
75 #
76
77 class HebbianDepthConv2d(nn.Module):
78     """
79     A 2d convolutional layer that learns through Hebbian plasticity
80     """
81
82     MODE_HPCA = 'h pca'
83     MODE_BASIC_HEBBIAN = 'basic'
84     MODE_WTA = 'wta'
85     MODE_SOFTWTA = 'soft'
86     MODE_BCM = 'bcm'
87     MODE_HARDWT = "hard"
88     MODE_PRESYNAPTIC_COMPETITION = "pre"
89     MODE_TEMPORAL_COMPETITION = "temp"
90     MODE_ADAPTIVE_THRESHOLD = "thresh"
91     MODE_ANTIHARDWT = "antihard"
92
93     def __init__(self, in_channels, out_channels, kernel_size,
94                  stride=1, dilation=1, padding=0, groups = 1,
95                               w_nrm=False, bias=False, act=nn.Identity(),
96                               mode=MODE_SOFTWTA, k=1, patchwise=True,
97                               contrast=1., uniformity=False, alpha=1.,
98                               wta_competition='similarity_spatial',
99                               lateral_competition="filter",
100                              lateral_inhibition_strength=0.01, top_k=1,
101                              prune_rate=0, t_invert=1.):
102
103     """
104
105         :param out_channels: output channels of the convolutional
106             kernel
107         :param in_channels: input channels of the convolutional
108             kernel
109         :param kernel_size: size of the convolutional kernel (int or
110             tuple)
111         :param stride: stride of the convolutional kernel (int or
112             tuple)
113         :param w_nrm: whether to normalize the weight vectors before
114             computing outputs

```

```

103     :param act: the nonlinear activation function after
104         convolution
105     :param mode: the learning mode, either 'swta' or 'h pca'
106     :param k: softmax inverse temperature parameter used for
107         swta-type learning
108     :param patchwise: whether updates for each convolutional
109         patch should be computed separately,
110         and then aggregated
111     :param contrast: coefficient that rescales negative compared
112         to positive updates in contrastive-type learning
113     :param uniformity: whether to use uniformity weighting in
114         contrastive-type learning.
115     :param alpha: weighting coefficient between hebbian and
116         backprop updates (0 means fully backprop, 1 means fully
117         hebbian).
118     """
119     super(HebbianDepthConv2d, self).__init__()
120     # spatial competition is the appropriate comp modes
121     self.mode = mode
122     self.out_channels = out_channels
123     self.in_channels = in_channels
124     self.kernel = kernel_size
125     self.kernel_size = _pair(kernel_size)
126     self.stride = _pair(stride)
127     self.dilation = _pair(dilation)
128     self.padding = padding
129     self.padding_mode = 'reflect'
130     if mode == "hard":
131         self.padding_mode = 'symmetric'
132     self.F_padding = (padding, padding, padding, padding)
133     self.groups = in_channels # in_channels for depthwise
134
135     # # Depthwise separable weights
136     weight_range = 25 / math.sqrt(in_channels * kernel_size *
137         kernel_size)
138     self.weight = nn.Parameter(weight_range * torch.randn(
139         in_channels, 1, *self.kernel_size))
140
141     # self.weight = nn.Parameter(torch.empty(in_channels, 1, *
142         # self.kernel_size))
143     # init.kaiming_uniform_(self.weight)
144     # self.weight = center_surround_init(in_channels, 1,
145         kernel_size, 1)
146
147     self.w_nrm = w_nrm
148     self.act = act
149     self.theta_decay = 0.5
150     if mode == "bcm":
151         self.theta = nn.Parameter(torch.ones(out_channels))
152
153     self.register_buffer('delta_w', torch.zeros_like(self.weight
154         ))
155     self.top_k = top_k
156     self.patchwise = patchwise
157     self.contrast = contrast
158     self.uniformity = uniformity
159     self.alpha = alpha
160     self.wta_competition = wta_competition

```

```

149     self.lateral_inhibition_mode = lateral_competition
150     self.lateral_learning_rate = lateral_inhibition_strength # 
151         Adjust as needed
152     self.lebesgue_p = 2
153
154     self.prune_rate = prune_rate # 99% of connections are
155         pruned
156     self.t_invert = torch.tensor(t_invert)
157
158     self.presynaptic_competition_type = "softmax"
159     self.presynaptic_weights = False # presynaptic competition
160         in forward pass
161
162     self.activation_history = None
163     self.temporal_window = 500
164     self.competition_k = 2
165     self.competition_type = "hard"
166
167     if kernel_size !=1:
168         self.sm_kernel = create_sm_kernel()
169         self.register_buffer('surround_kernel', self.sm_kernel)
170
171     def apply_lebesgue_norm(self, w):
172         return torch.sign(w) * torch.abs(w) ** (self.lebesgue_p - 1)
173
174     def cosine(self, x, w):
175         w_normalized = F.normalize(w, p=2, dim=1)
176         conv_output = F.conv2d(x, w_normalized, None, self.stride,
177             0, self.dilation, groups=self.groups)
178         x_squared = x.pow(2)
179         x_squared_sum = F.conv2d(x_squared, torch.ones_like(w), None
180             , self.stride, 0, self.dilation,
181                 self.groups)
182         x_norm = torch.sqrt(x_squared_sum + 1e-8)
183         cosine_sim = conv_output / x_norm
184         return cosine_sim
185
186     def apply_weights(self, x, w):
187         """
188             This function provides the logic for combining input x and
189             weight w
190         """
191
192         # w = self.apply_lebesgue_norm(self.weight)
193         return F.conv2d(x, w, None, self.stride, 0, self.dilation,
194             groups=self.groups)
195
196     def apply_surround_modulation(self, y):
197         return F.conv2d(y, self.sm_kernel.repeat(self.out_channels,
198             1, 1, 1),
199                         padding=self.sm_kernel.size(-1) // 2, groups
200                             =self.out_channels)
201
202     def compute_activation(self, x):
203         x = symmetric_pad(x, self.padding)
204         w = self.weight
205         if self.w_nrm: w = normalize(w, dim=(1, 2, 3))
206         if self.presynaptic_weights: w = self.
207             compute_presynaptic_competition(w)

```

```

197     # y_depthwise = self.act(self.apply_weights(x, w))
198     # For cosine similarity activation if cosine is to be used
199     # for next layer
200     y_depthwise = self.cosine(x, w)
201     return x, y_depthwise, w
202
203 def forward(self, x):
204     x, y_depthwise, w = self.compute_activation(x)
205     if self.kernel != 1:
206         y_depthwise = self.apply_surround_modulation(y_depthwise)
207     if self.training:
208         self.compute_update(x, y_depthwise, w)
209     return y_depthwise
210
211 def compute_update(self, x, y, weight):
212     if self.mode == self.MODE_BASIC_HEBBIAN:
213         update = self.update_basic_hebbian(x, y, weight)
214     elif self.mode == self.MODE_HARDWT:
215         update = self.update_hardwt(x, y, weight)
216     elif self.mode == self.MODE_SOFTWTA:
217         update = self.update_softwta(x, y, weight)
218     elif self.mode == self.MODE_BCM:
219         update = self.update_bcm(x, y, weight)
220     elif self.mode == self.MODE_TEMPORAL_COMPETITION:
221         update = self.update_temporal_competition(x, y, weight)
222     elif self.mode == self.MODE_ADAPTIVE_THRESHOLD:
223         update = self.update_adaptive_threshold(x, y, weight)
224     else:
225         raise NotImplementedError(f"Learning mode {self.mode} unavailable for {self.__class__.__name__} layer")
226
227     # Weight Normalization and added to weight change buffer
228     update.div_(torch.abs(update).amax() + 1e-30)
229     self.delta_w += update
230
231 def update_basic_hebbian(self, x, y, weight):
232     yx = self.compute_yx(x, y)
233     y_sum = y.sum(dim=(0, 2, 3)).view(self.in_channels, 1, 1, 1)
234     yw = y_sum * weight
235     return yx - yw
236
237 def update_hardwt(self, x, y, weight):
238     y_wta = y * self.compute_wta_mask(y)
239     yx = self.compute_yx(x, y_wta)
240     yu = torch.sum(y_wta, dim=(0, 2, 3)).view(self.in_channels,
241         1, 1, 1)
242     return yx - yu * weight
243
244 def update_softwta(self, x, y, weight):
245     softwta_actvns = self.compute_softwta_activations(y)
246     yx = self.compute_yx(x, softwta_actvns)
247     yu = torch.sum(torch.mul(softwta_actvns, y), dim=(0, 2, 3)).
248         view(self.in_channels, 1, 1, 1)
249     return yx - yu * weight
250
251 def update_bcm(self, x, y, weight):
252     batch_size, out_channels, height, width = y.shape

```

```

250     y_wta = y * self.compute_wta_mask(y)
251     # Compute squared activation for each spatial location and
252     # channel
253     y_squared = y_wta.pow(2)
254     # Reshape theta to match spatial dimensions
255     theta_spatial = self.theta.view(1, -1, 1, 1).expand(1,
256         out_channels, height, width)
257     # Update theta for each spatial location
258     theta_update = self.theta_decay * (y_squared - theta_spatial
259         )
260     self.theta.data += theta_update.mean(dim=(0, 2, 3))
261     # Compute BCM updates for each spatial location
262     y_minus_theta = y_wta - theta_spatial
263     bcm_factor = y_wta * y_minus_theta
264     yx = self.compute_yx(x, bcm_factor)
265     update = yx.view(weight.shape)
266     return update
267
268     def update_temporal_competition(self, x, y, weight):
269         batch_size, out_channels, height_out, width_out = y.shape
270         self.update_activation_history(y)
271         temporal_winners = self.compute_temporal_winners(y)
272         y_winners = temporal_winners * y
273         y_winners = self.apply_competition(y_winners, batch_size,
274             out_channels)
275         yx = self.compute_yx(x, y_winners)
276         y_sum = y_winners.sum(dim=(0, 2, 3)).view(self.in_channels,
277             1, 1, 1)
278         update = yx - y_sum * weight
279         return update
280
281     def update_adaptive_threshold(self, x, y, weight):
282         batch_size, out_channels, height_out, width_out = y.shape
283         similarities = F.conv2d(x, weight, stride=self.stride,
284             padding=self.padding, groups=self.groups)
285         similarities = similarities / (torch.norm(weight.view(
286             out_channels, -1), dim=1).view(1, -1, 1, 1) + 1e-10)
287         threshold = self.compute_adaptive_threshold(similarities)
288         winners = (similarities > threshold).float()
289         y_winners = winners * similarities
290         y_winners = self.apply_competition(y_winners, batch_size,
291             out_channels)
292         yx = self.compute_yx(x, y_winners)
293         y_sum = y_winners.sum(dim=(0, 2, 3)).view(self.in_channels,
294             1, 1, 1)
295         update = yx - y_sum * weight
296         return update
297
298     def compute_yx(self, x, y):
299         yx = F.conv2d(x.transpose(0, 1), y.transpose(0, 1), padding
300             =0,
301                 stride=self.dilation, dilation=self.stride).
302                 transpose(0, 1)
303         yx = yx.diagonal(dim1=0, dim2=1).permute(2, 0, 1).unsqueeze
304             (1)
305         return yx
306
307     def compute_wta_mask(self, y):

```

```

296     batch_size, out_channels, height_out, width_out = y.shape
297     # WTA competition within each channel
298     y_flat = y.view(batch_size, out_channels, -1)
299     win_neurons = torch.argmax(y_flat, dim=2)
300     wta_mask = F.one_hot(win_neurons, num_classes=height_out *
301                           width_out).float()
302     return wta_mask.view(batch_size, out_channels, height_out,
303                           width_out)
304
305 def compute_softwta_activations(self, y):
306     # Competition and anti-Hebbian learning for y_depthwise
307     batch_size, in_channels, height_depthwise, width_depthwise =
308         y.shape
309     # Reshape to apply softmax within each channel
310     y_depthwise_reshaped = y.view(batch_size, in_channels, -1)
311     # Apply softmax within each channel
312     flat_softwta_actvts_depthwise = torch.softmax(self.t_invert
313             * y_depthwise_reshaped, dim=2)
314     # Turn all postsynaptic activations into anti-Hebbian
315     flat_softwta_actvts_depthwise = -
316         flat_softwta_actvts_depthwise
317     # Find winners within each channel
318     win_neurons_depthwise = torch.argmax(y_depthwise_reshaped,
319                                         dim=2)
320     # Create a mask to flip the sign of winning neurons
321     mask = torch.zeros_like(flat_softwta_actvts_depthwise)
322     mask.scatter_(2, win_neurons_depthwise.unsqueeze(2), 1)
323     # Flip the sign of winning neurons
324     flat_softwta_actvts_depthwise =
325         flat_softwta_actvts_depthwise * (1 - 2 * mask)
326     # Reshape back to original shape
327     return flat_softwta_actvts_depthwise.view(batch_size,
328                                               in_channels, height_depthwise,
329                                               width_depthwise)
330
331 def update_activation_history(self, y):
332     if self.activation_history is None:
333         self.activation_history = y.detach().clone()
334     else:
335         self.activation_history = torch.cat([self.
336             activation_history, y.detach()], dim=0)
337     if self.activation_history.size(0) > self.
338         temporal_window:
339         self.activation_history = self.activation_history[-
340             self.temporal_window:]
341
342 def compute_temporal_winners(self, y):
343     batch_size, out_channels, height_out, width_out = y.shape
344     history_spatial = self.activation_history.view(-1,
345             out_channels, height_out, width_out)
346     median_activations = torch.median(history_spatial, dim=0)[0]
347     temporal_threshold = torch.mean(median_activations, dim=(1,
348                     2), keepdim=True)
349     return (median_activations > temporal_threshold).float()
350
351 def compute_adaptive_threshold(self, similarities):

```

```

339     mean_sim = similarities.mean(dim=(2, 3), keepdim=True)
340     std_sim = similarities.std(dim=(2, 3), keepdim=True)
341     return mean_sim + self.competition_k * std_sim
342
343     def apply_competition(self, y, batch_size, out_channels):
344         if self.mode in [self.MODE_TEMPORAL_COMPETITION, self.
345             MODE_ADAPTIVE_THRESHOLD]:
346             if self.competition_type == 'hard':
347                 y = y.view(batch_size, out_channels, -1)
348                 top_k_indices = torch.topk(y, self.top_k, dim=2,
349                     largest=True, sorted=False).indices
350                 y_compete = torch.zeros_like(y)
351                 y_compete.scatter_(2, top_k_indices, y.gather(2,
352                     top_k_indices))
353                 return y_compete.view_as(y)
354             elif self.competition_type == 'soft':
355                 return torch.softmax(self.t_invert * y.view(
356                     batch_size, out_channels, -1), dim=2).view_as(y)
357         return y
358
359     @torch.no_grad()
360     def local_update(self):
361         """
362             This function transfers a previously computed weight update,
363             stored in buffer self.delta_w, to the gradient
364             self.weight.grad of the weight parameter.
365
366             This function should be called before optimizer.step(), so
367             that the optimizer will use the locally computed
368             update as optimization direction. Local updates can also be
369             combined with end-to-end updates by calling this
370             function between loss.backward() and optimizer.step(). loss.
371             backward will store the end-to-end gradient in
372             self.weight.grad, and this function combines this value with
373             self.delta_w as
374             self.weight.grad = (1 - alpha) * self.weight.grad - alpha *
375             self.delta_w
376             Parameter alpha determines the scale of the local update
377             compared to the end-to-end gradient in the combination.
378             """
379             if self.weight.grad is None:
380                 self.weight.grad = -self.alpha * self.delta_w
381             else:
382                 self.weight.grad = (1 - self.alpha) * self.weight.grad -
383                     self.alpha * self.delta_w
384             self.delta_w.zero_()
385
386             # @torch.no_grad()
387             # # Weight Update
388             # def local_update(self):
389             #     new_weight = self.weight + 0.1 * self.alpha * self.delta_w
390             #     # Update weights
391             #     self.weight.copy_(new_weight)
392             #     # self.structural_plasticity()
393             #     self.delta_w.zero_()

```

A.2.4 Hebb_abs.py

```

1 import math
2
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 from torch.nn.modules.utils import _pair
7 import matplotlib.pyplot as plt
8
9 import torch.nn.init as init
10
11 """
12 Uses almost identical code to hebb.py. Please refer to hebb.py for
13 additional explanations on code functionality
14 Only changes are the initialisation and update of weights, to ensure
15 adherence to Dale's Principle
16 """
17
18
19
20 def normalize(x, dim=None):
21     nrm = (x ** 2).sum(dim=dim, keepdim=True) ** 0.5
22     nrm[nrm == 0] = 1.
23     return x / nrm
24
25
26 def symmetric_pad(x, padding):
27     if padding == 0:
28         return x
29     return F.pad(x, (padding,) * 4, mode='reflect')
30
31
32 def create_sm_kernel(kernel_size=5, sigma_e=1.2, sigma_i=1.4):
33 """
34 Create a surround modulation kernel that matches the paper's
35 specifications.
36 :param kernel_size: Size of the SM kernel.
37 :param sigma_e: Standard deviation for the excitatory Gaussian.
38 :param sigma_i: Standard deviation for the inhibitory Gaussian.
39 :return: A normalized SM kernel.
40 """
41     center = kernel_size // 2
42     x, y = torch.meshgrid(torch.arange(kernel_size), torch.arange(
43         kernel_size), indexing="ij")
44     x = x.float() - center
45     y = y.float() - center
46     # Compute the excitatory and inhibitory Gaussians
47     gaussian_e = torch.exp(-(x ** 2 + y ** 2) / (2 * sigma_e ** 2))
48     gaussian_i = torch.exp(-(x ** 2 + y ** 2) / (2 * sigma_i ** 2))
49     # Compute the Difference of Gaussians (DoG)
50     dog = gaussian_e / (2 * math.pi * sigma_e ** 2) - gaussian_i / (
51         2 * math.pi * sigma_i ** 2)
52     # Normalize the DoG so that the center value is 1
53     sm_kernel = dog / dog[center, center]
54     return sm_kernel.unsqueeze(0).unsqueeze(0).to(device)

```

```

52
53
54 class HebbianConv2d(nn.Module):
55     """
56     A 2d convolutional layer that learns through Hebbian plasticity
57     """
58
59     MODE_HPCA = 'h pca'
60     MODE_BASIC_HEBBIAN = 'basic'
61     MODE_WTA = 'wta'
62     MODE_SOFTWTA = 'soft'
63     MODE_BCM = 'bcm'
64     MODE_HARDWT = "hard"
65     MODE_PRESYNAPTIC_COMPETITION = "pre"
66     MODE_TEMPORAL_COMPETITION = "temp"
67     MODE_ADAPTIVE_THRESHOLD = "thresh"
68     MODE_ANTIHARDWT = "antihard"
69
70     def __init__(self, in_channels, out_channels, kernel_size,
71                  stride=1, dilation=1, padding=0, groups=1,
72                  w_nrm=False, bias=False, act=nn.Identity(),
73                  mode=MODE_SOFTWTA, k=1, patchwise=True,
74                  contrast=1., uniformity=False, alpha=1.,
75                  wta_competition='similarity_spatial',
76                  lateral_competition="combined",
77                  lateral_inhibition_strength=0.01, top_k=1,
78                  prune_rate=0, t_invert=1.):
79         """
80
81         :param out_channels: output channels of the convolutional
82             kernel
83         :param in_channels: input channels of the convolutional
84             kernel
85         :param kernel_size: size of the convolutional kernel (int or
86             tuple)
87         :param stride: stride of the convolutional kernel (int or
88             tuple)
89         :param w_nrm: whether to normalize the weight vectors before
90             computing outputs
91         :param act: the nonlinear activation function after
92             convolution
93         :param mode: the learning mode, either 'swta' or 'h pca'
94         :param k: softmax inverse temperature parameter used for
             swta-type learning
         :param patchwise: whether updates for each convolutional
             patch should be computed separately,
             and then aggregated
         :param contrast: coefficient that rescales negative compared
             to positive updates in contrastive-type learning
         :param uniformity: whether to use uniformity weighting in
             contrastive-type learning.
         :param alpha: weighting coefficient between hebbian and
             backprop updates (0 means fully backprop, 1 means fully
             hebbian).
95
96         super(HebbianConv2d, self).__init__()
97         # spatial competition is the appropriate comp modes
98         self.mode = mode

```

```

95     self.out_channels = out_channels
96     self.in_channels = in_channels
97     self.kernel = kernel_size
98     self.kernel_size = _pair(kernel_size)
99     self.stride = _pair(stride)
100    self.dilation = _pair(dilation)
101    self.padding = padding
102    self.padding_mode = 'reflect'
103    if mode == "hard":
104        self.padding_mode = 'symmetric'
105
106    self.F_padding = (padding, padding, padding, padding)
107    self.groups = 1 # in_channels for depthwise
108
109    weight_range = 25 / math.sqrt(in_channels * kernel_size *
110        kernel_size)
111    self.weight = nn.Parameter(weight_range * torch.abs(torch.
112        randn((out_channels, in_channels // self.groups, *self.
113            kernel_size))))
114
115    print(self.weight.shape)
116    self.w_nrm = w_nrm
117    self.act = act
118    self.theta_decay = 0.5
119    if mode == "bcm":
120        self.theta = nn.Parameter(torch.ones(out_channels),
121            requires_grad=False)
122
123    self.register_buffer('delta_w', torch.zeros_like(self.weight
124        ))
125    self.top_k = top_k
126    self.patchwise = patchwise
127    self.contrast = contrast
128    self.uniformity = uniformity
129    self.alpha = alpha
130    self.wta_competition = wta_competition
131    self.lateral_inhibition_mode = lateral_competition
132    self.lateral_learning_rate = lateral_inhibition_strength # #
133    Adjust as needed
134    self.lebesgue_p = 2
135
136    self.prune_rate = prune_rate # 99% of connections are
137        pruned
138    self.t_invert = torch.tensor(t_invert)
139
140    self.presynaptic_competition_type = "softmax"
141    self.presynaptic_weights = False # presynaptic competition
142        in forward pass
143
144    self.activation_history = None
145    self.temporal_window = 500
146    self.competition_k = 2
147    self.competition_type = "hard"
148
149    if self.kernel != 1:
150        self.sm_kernel = create_sm_kernel()
151        self.register_buffer('surround_kernel', self.sm_kernel)
152        self.visualize_surround_modulation_kernel()

```

```

145
146     self.target_activity = 0.1
147     self.scaling_rate = 0.001
148     self.register_buffer('average_activity', torch.zeros(
149         out_channels))
150
151     self.growth_probability = 0.1
152     self.new_synapse_strength = 1
153     self.prune_threshold_percentile = 10 # Prune bottom 10% of
154     weights
155
156     def visualize_surround_modulation_kernel(self):
157         """
158             Visualizes the surround modulation kernel using matplotlib.
159         """
160         sm_kernel = self.sm_kernel.squeeze().cpu().detach().numpy()
161         plt.figure(figsize=(5, 5))
162         plt.imshow(sm_kernel, cmap='jet')
163         plt.colorbar()
164         plt.title('Surround Modulation Kernel')
165         plt.show()
166
167     def apply_lebesgue_norm(self, w):
168         return torch.sign(w) * torch.abs(w)**(self.lebesgue_p - 1)
169
170     def apply_weights(self, x, w):
171         """
172             This function combines input x and weight w
173         """
174         return F.conv2d(x, w, None, self.stride, 0, self.dilation,
175                         groups=self.groups)
176
177     def update_average_activity(self, y):
178         current_activity = y.mean(dim=(0, 2, 3))
179         self.average_activity = 0.9 * self.average_activity + 0.1 *
180             current_activity
181
182     def synaptic_scaling(self):
183         scale_factor = self.target_activity / (self.average_activity
184             + 1e-6)
185         self.weight.data *= (1 + self.scaling_rate * (scale_factor -
186             1)).view(-1, 1, 1, 1)
187
188     def structural_plasticity(self):
189         with torch.no_grad():
190             # Pruning step
191             prune_threshold = torch.quantile(torch.abs(self.weight),
192                 self.prune_threshold_percentile / 100)
193             weak_synapses = torch.abs(self.weight) < prune_threshold
194             self.weight.data[weak_synapses] = 0
195             # Growth step
196             zero_weights = self.weight.data == 0
197             new_synapses = torch.rand_like(self.weight) < self.
198                 growth_probability
199             new_synapses &= zero_weights
200             self.weight.data[new_synapses] = torch.randn_like(self.
201                 weight)[new_synapses] * self.new_synapse_strength

```

```

194     def cosine(self, x, w):
195         w_normalized = F.normalize(w, p=2, dim=1)
196         conv_output = F.conv2d(x, w_normalized, None, self.stride,
197                               0, self.dilation, groups=self.groups)
198         x_squared = x.pow(2)
199         x_squared_sum = F.conv2d(x_squared, torch.ones_like(w), None
200                               , self.stride, 0, self.dilation,
201                               self.groups)
202         x_norm = torch.sqrt(x_squared_sum + 1e-8)
203         cosine_sim = conv_output / x_norm
204         return cosine_sim
205
206
207     def apply_surround_modulation(self, y):
208         return F.conv2d(y, self.sm_kernel.repeat(self.out_channels,
209                         1, 1, 1),
210                         padding=self.sm_kernel.size(-1) // 2, groups
211                         =self.out_channels)
212
213
214     def compute_activation(self, x):
215         x = symmetric_pad(x, self.padding)
216         w = self.weight.abs()
217         if self.w_nrm: w = normalize(w, dim=(1, 2, 3))
218         if self.presynaptic_weights: w = self.
219             compute_presynaptic_competition_global(w)
220         # y = self.act(self.apply_weights(x, w))
221         # For cosine similarity activation if cosine is to be used
222             for next layer
223         y = self.cosine(x, w)
224         return x, y, w
225
226
227     def forward(self, x):
228         x, y, w = self.compute_activation(x)
229         if self.kernel != 1:
230             y = self.apply_surround_modulation(y)
231         if self.training:
232             # self.update_average_activity(y)
233             # self.synaptic_scaling()
234             self.compute_update(x, y, w)
235         return y
236
237
238     def compute_update(self, x, y, weight):
239         if self.mode == self.MODE_BASIC_HEBBIAN:
240             update = self.update_basic_hebbian(x, y, weight)
241         elif self.mode == self.MODE_HARDWT:
242             update = self.update_hardwt(x, y, weight)
243         elif self.mode == self.MODE_SOFTWTA:
244             update = self.update_softwta(x, y, weight)
245         elif self.mode == self.MODE_BCM:
246             update = self.update_bcm(x, y, weight)
247         else:
248             raise NotImplementedError(f"Learning mode {self.mode}
249                                     unavailable for {self.__class__.__name__} layer")
250         # Weight Normalization and added to weight change buffer
251         update.div_(torch.abs(update).amax() + 1e-30)
252         self.delta_w += update
253
254
255     def update_basic_hebbian(self, x, y, weight):
256         # Grossberg Instar rule

```

```

245     yx = self.compute_yx(x, y)
246     y_sum = y.sum(dim=(0, 2, 3)).view(-1, 1, 1, 1)
247     yw = y_sum * weight
248     update = yx - yw
249     return update
250
251 def update_hardwt(self, x, y, weight):
252     # Grossberg Instar with wta mask
253     y_wta = y * self.compute_wta_mask(y)
254     yx = self.compute_yx(x, y_wta)
255     yu = torch.sum(y_wta, dim=(0, 2, 3))
256     update = yx - yu.view(-1, 1, 1, 1) * weight
257     return update
258
259 def update_softwta(self, x, y, weight):
260     # SoftHebb Grossberg Instar Variation
261     softwta_activs = self.compute_softwta_activations(y)
262     yx = self.compute_yx(x, softwta_activs)
263     yu = torch.sum(torch.mul(softwta_activs, y), dim=(0, 2, 3))
264     update = yx - yu.view(-1, 1, 1, 1) * weight
265     return update
266
267 def update_bcm(self, x, y, weight):
268     # BCM learning rule with WTA mask (remove if not needed)
269     y_wta = y * self.compute_wta_mask(y)
270     y_squared = y_wta.pow(2).mean(dim=(0, 2, 3))
271     self.theta.data = (1 - self.theta_decay) * self.theta + self
272         .theta_decay * y_squared
273     y_minus_theta = y_wta - self.theta.view(1, -1, 1, 1)
274     bcm_factor = y_wta * y_minus_theta
275     yx = self.compute_yx(x, bcm_factor)
276     update = yx.view(weight.shape)
277     return update
278
279 def compute_yx(self, x, y):
280     # Computes common y*w term from Grossberg Instar
281     yx = F.conv2d(x.transpose(0, 1), y.transpose(0, 1), padding
282         =0,
283             stride=self.dilation, dilation=self.stride).
284                 transpose(0, 1)
285     if self.groups != 1:
286         yx = yx.mean(dim=1, keepdim=True)
287     return yx
288
289 def compute_wta_mask(self, y):
290     # Computes WTA Mask
291     batch_size, out_channels, height_out, width_out = y.shape
292     y_flat = y.transpose(0, 1).reshape(out_channels, -1)
293     win_neurons = torch.argmax(y_flat, dim=0)
294     wta_mask = F.one_hot(win_neurons, num_classes=out_channels).
295         float()
296     return wta_mask.transpose(0, 1).view(out_channels,
297         batch_size, height_out, width_out).transpose(0, 1)
298
299 def compute_softwta_activations(self, y):
300     # Computes SoftHebb mask and Hebb/AntiHebb implementation
301     batch_size, out_channels, height_out, width_out = y.shape

```

```

297     flat_weighted_inputs = y.transpose(0, 1).reshape(
298         out_channels, -1)
299     flat_softwta_activs = torch.softmax(self.t_invert *
300         flat_weighted_inputs, dim=0)
301     flat_softwta_activs = -flat_softwta_activs
302     win_neurons = torch.argmax(flat_weighted_inputs, dim=0)
303     competing_idx = torch.arange(flat_weighted_inputs.size(1))
304     flat_softwta_activs[win_neurons, competing_idx] = -
305         flat_softwta_activs[win_neurons, competing_idx]
306     return flat_softwta_activs.view(out_channels, batch_size,
307         height_out, width_out).transpose(0, 1)
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

```

A.2.5 Hebb_abs_depthwise.py

```

1 import math
2
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 from torch.nn.modules.utils import _pair
7 import torch.nn.init as init
8
9
10 """
11 Uses almost identical code to hebb.py. Please refer to hebb.py for
12 additional explanations on code functionality
13 Only changes are the competition modes, which apply competition
14 across spatial neurons in a filter for channel independence
15 Cahnegs also in initialisation and update, similar to hebb_abs.py
16 """
17
18 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
19
20 torch.manual_seed(0)

```

```

18
19
20 def normalize(x, dim=None):
21     nrm = (x ** 2).sum(dim=dim, keepdim=True) ** 0.5
22     nrm[nrm == 0] = 1.
23     return x / nrm
24
25
26 def symmetric_pad(x, padding):
27     if padding == 0:
28         return x
29     return F.pad(x, (padding,) * 4, mode='reflect')
30
31
32 def center_surround_init(out_channels, in_channels, kernel_size,
33                           groups=1):
34     # Calculate weight range
35     weight_range = 25 / math.sqrt(in_channels * kernel_size *
36                                   kernel_size)
37
38     # Calculate sigma based on kernel size (using equation 3 from
39     # the paper)
40     gamma = torch.empty(out_channels).uniform_(0, 0.5)
41     sigma = (kernel_size / 4) * torch.sqrt((1 - gamma ** 2) / (-
42                                             torch.log(gamma)))
43
44     # Create meshgrid for x and y coordinates
45     x = torch.linspace(-(kernel_size - 1) / 2, (kernel_size - 1) /
46                         2, kernel_size)
47     y = torch.linspace(-(kernel_size - 1) / 2, (kernel_size - 1) /
48                         2, kernel_size)
49     xx, yy = torch.meshgrid(x, y, indexing='ij')
50
51     # Calculate center and surround Gaussians
52     center = torch.exp(-(xx ** 2 + yy ** 2) / (2 * (gamma.view(-1,
53                                                 1, 1) * sigma.view(-1, 1, 1)) ** 2))
54     surround = torch.exp(-(xx ** 2 + yy ** 2) / (2 * sigma.view(-1,
55                                                 1, 1) ** 2))
56
57     # Calculate DoG (Difference of Gaussians)
58     dog = center - surround
59
60     # Normalize DoG
61     ac = torch.sum(torch.clamp(dog, min=0))
62     as_ = torch.sum(-torch.clamp(dog, max=0))
63     dog = weight_range * 0.5 * dog / (ac + as_)
64
65     # Assign excitatory (positive) or inhibitory (negative) centers
66     center_type = torch.cat([torch.ones(out_channels // 2), -torch.
67                             ones(out_channels - out_channels // 2)])
68     center_type = center_type[torch.randperm(out_channels)].view(-1,
69                         1, 1)
70     dog = dog * center_type
71
72     # Repeat for in_channels and reshape to match conv2d weight
73     # shape
74     dog = dog.unsqueeze(1).repeat(1, in_channels // groups, 1, 1)

```

```

64     dog = dog.reshape(out_channels, in_channels // groups,
65                       kernel_size, kernel_size)
66
67     return nn.Parameter(dog)
68
69 def create_sm_kernel(kernel_size=5, sigma_e=1.2, sigma_i=1.4):
70     center = kernel_size // 2
71     x, y = torch.meshgrid(torch.arange(kernel_size), torch.arange(
72                           kernel_size))
73     x = x.float() - center
74     y = y.float() - center
75
76     gaussian_e = torch.exp(-(x ** 2 + y ** 2) / (2 * sigma_e ** 2))
77     gaussian_i = torch.exp(-(x ** 2 + y ** 2) / (2 * sigma_i ** 2))
78
79     dog = gaussian_e / (2 * math.pi * sigma_e ** 2) - gaussian_i / (
80                           2 * math.pi * sigma_i ** 2)
81     sm_kernel = dog / dog[center, center]
82
83     return sm_kernel.unsqueeze(0).unsqueeze(0).to(device)
84
85
86 class HebbianDepthConv2d(nn.Module):
87     """
88         A 2d convolutional layer that learns through Hebbian plasticity
89     """
90
91     MODE_HPCA = 'h pca'
92     MODE_BASIC_HEBBIAN = 'basic'
93     MODE_WTA = 'wta'
94     MODE_SOFTWTA = 'soft'
95     MODE_BCM = 'bcm'
96     MODE_HARDWT = "hard"
97     MODE_PRESYNAPTIC_COMPETITION = "pre"
98     MODE_TEMPORAL_COMPETITION = "temp"
99     MODE_ADAPTIVE_THRESHOLD = "thresh"
100    MODE_ANTIHARDWT = "antihard"
101
102    def __init__(self, in_channels, out_channels, kernel_size,
103                 stride=1, dilation=1, padding=0, groups = 1,
104                             w_nrm=False, bias=False, act=nn.Identity(),
105                             mode=MODE_SOFTWTA, k=1, patchwise=True,
106                             contrast=1., uniformity=False, alpha=1.,
107                             wta_competition='similarity_spatial',
108                             lateral_competition="filter",
109                             lateral_inhibition_strength=0.01, top_k=1,
110                             prune_rate=0, t_invert=1.):
111        """
112            :param out_channels: output channels of the convolutional
113                                  kernel
114            :param in_channels: input channels of the convolutional
115                                  kernel
116            :param kernel_size: size of the convolutional kernel (int or
117                               tuple)
118            :param stride: stride of the convolutional kernel (int or
119                          tuple)

```

```

111     :param w_nrm: whether to normalize the weight vectors before
112         computing outputs
113     :param act: the nonlinear activation function after
114         convolution
115     :param mode: the learning mode, either 'swta' or 'hpca'
116     :param k: softmax inverse temperature parameter used for
117         swta-type learning
118     :param patchwise: whether updates for each convolutional
119         patch should be computed separately,
120         and then aggregated
121     :param contrast: coefficient that rescales negative compared
122         to positive updates in contrastive-type learning
123     :param uniformity: whether to use uniformity weighting in
124         contrastive-type learning.
125     :param alpha: weighting coefficient between hebbian and
126         backprop updates (0 means fully backprop, 1 means fully
127         hebbian).
128     """
129     super(HebbianDepthConv2d, self).__init__()
130     # spatial competition is the appropriate comp modes
131     self.mode = mode
132     self.out_channels = out_channels
133     self.in_channels = in_channels
134     self.kernel = kernel_size
135     self.kernel_size = _pair(kernel_size)
136     self.stride = _pair(stride)
137     self.dilation = _pair(dilation)
138     self.padding = padding
139     self.padding_mode = 'reflect'
140     if mode == "hard":
141         self.padding_mode = 'symmetric'
142     self.F_padding = (padding, padding, padding, padding)
143     self.groups = in_channels # in_channels for depthwise
144
145     # Depthwise separable weights
146     weight_range = 25 / math.sqrt(in_channels * kernel_size *
147         kernel_size)
148     self.weight = nn.Parameter(weight_range * torch.abs(torch.
149         randn(in_channels, 1, *self.kernel_size)))
150
151     # self.weight = center_surround_init(in_channels, 1,
152     #     kernel_size, 1)
153
154     self.w_nrm = w_nrm
155     self.act = act
156     self.theta_decay = 0.5
157     if mode == "bcm":
158         self.theta = nn.Parameter(torch.ones(out_channels))
159
160     self.register_buffer('delta_w', torch.zeros_like(self.weight
161         ))
162     self.top_k = top_k
163     self.patchwise = patchwise
164     self.contrast = contrast
165     self.uniformity = uniformity
166     self.alpha = alpha
167     self.wta_competition = wta_competition
168     self.lateral_inhibition_mode = lateral_competition

```

```

157     self.lateral_learning_rate = lateral_inhibition_strength # 
158         Adjust as needed
159     self.lebesgue_p = 2
160
161     self.prune_rate = prune_rate # 99% of connections are
162         pruned
163     self.t_invert = torch.tensor(t_invert)
164
165     self.presynaptic_competition_type = "softmax"
166     self.presynaptic_weights = False # presynaptic competition
167         in forward pass
168
169     self.activation_history = None
170     self.temporal_window = 500
171     self.competition_k = 2
172     self.competition_type = "hard"
173
174     if kernel_size !=1:
175         self.sm_kernel = create_sm_kernel()
176         self.register_buffer('surround_kernel', self.sm_kernel)
177
178     def apply_lebesgue_norm(self, w):
179         return torch.sign(w) * torch.abs(w) ** (self.lebesgue_p - 1)
180
181     def cosine(self, x, w):
182         w_normalized = F.normalize(w, p=2, dim=1)
183         conv_output = F.conv2d(x, w_normalized, None, self.stride,
184             0, self.dilation, groups=self.groups)
185         x_squared = x.pow(2)
186         x_squared_sum = F.conv2d(x_squared, torch.ones_like(w), None
187             , self.stride, 0, self.dilation,
188                 self.groups)
189         x_norm = torch.sqrt(x_squared_sum + 1e-8)
190         cosine_sim = conv_output / x_norm
191         return cosine_sim
192
193     def apply_weights(self, x, w):
194         """
195             This function combines input x and weight w
196         """
197         # w = self.apply_lebesgue_norm(self.weight)
198
199         return F.conv2d(x, w, None, self.stride, 0, self.dilation,
200             groups=self.groups)
201
202     def apply_surround_modulation(self, y):
203         return F.conv2d(y, self.sm_kernel.repeat(self.out_channels,
204             1, 1, 1),
205                         padding=self.sm_kernel.size(-1) // 2, groups
206                         =self.out_channels)
207
208     def compute_activation(self, x):
209         x = symmetric_pad(x, self.padding)
210         w = self.weight.abs()
211         if self.w_nrm: w = normalize(w, dim=(1, 2, 3))
212         if self.presynaptic_weights: w = self.
213             compute_presynaptic_competition(w)
214         # y_depthwise = self.act(self.apply_weights(x, w))

```

```

206     # For cosine similarity activation if cosine is to be used
207     # for next layer
208     y_depthwise = self.cosine(x, w)
209     return x, y_depthwise, w
210
211 def forward(self, x):
212     x, y_depthwise, w = self.compute_activation(x)
213     if self.kernel != 1:
214         y_depthwise = self.apply_surround_modulation(y_depthwise)
215     if self.training:
216         self.compute_update(x, y_depthwise, w)
217     return y_depthwise
218
219 def compute_update(self, x, y, weight):
220     if self.mode == self.MODE_BASIC_HEBBIAN:
221         update = self.update_basic_hebbian(x, y, weight)
222     elif self.mode == self.MODE_HARDWT:
223         update = self.update_hardwt(x, y, weight)
224     elif self.mode == self.MODE_SOFTWTA:
225         update = self.update_softwta(x, y, weight)
226     elif self.mode == self.MODE_ANTIHARDWT:
227         update = self.update_antihardwt(x, y, weight)
228     elif self.mode == self.MODE_BCM:
229         update = self.update_bcm(x, y, weight)
230     elif self.mode == self.MODE_TEMPORAL_COMPETITION:
231         update = self.update_temporal_competition(x, y, weight)
232     elif self.mode == self.MODE_ADAPTIVE_THRESHOLD:
233         update = self.update_adaptive_threshold(x, y, weight)
234     else:
235         raise NotImplementedError(f"Learning mode {self.mode} unavailable for {self.__class__.__name__} layer")
236
237     # Weight Normalization and added to weight change buffer
238     update.div_(torch.abs(update).amax() + 1e-30)
239     self.delta_w += update
240
241 def update_basic_hebbian(self, x, y, weight):
242     yx = self.compute_yx(x, y)
243     y_sum = y.sum(dim=(0, 2, 3)).view(self.in_channels, 1, 1, 1)
244     yw = y_sum * weight
245     return yx - yw
246
247 def update_hardwt(self, x, y, weight):
248     y_wta = y * self.compute_wta_mask(y)
249     yx = self.compute_yx(x, y_wta)
250     yu = torch.sum(y_wta, dim=(0, 2, 3)).view(self.in_channels,
251         1, 1, 1)
252     return yx - yu * weight
253
254 def update_softwta(self, x, y, weight):
255     softwta_activs = self.compute_softwta_activations(y)
256     yx = self.compute_yx(x, softwta_activs)
257     yu = torch.sum(torch.mul(softwta_activs, y), dim=(0, 2, 3)).
258         view(self.in_channels, 1, 1, 1)
259     return yx - yu * weight
260
261 def update_bcm(self, x, y, weight):

```

```

259     y_wta = y * self.compute_wta_mask(y)
260     y_squared = y_wta.pow(2).mean(dim=(0, 2, 3))
261     self.theta.data = (1 - self.theta_decay) * self.theta + self
262         .theta_decay * y_squared
263     y_minus_theta = y_wta - self.theta.view(1, -1, 1, 1)
264     bcm_factor = y_wta * y_minus_theta
265     yx = self.compute_yx(x, bcm_factor)
266     return yx.view(weight.shape)
267
267 def update_temporal_competition(self, x, y, weight):
268     batch_size, out_channels, height_out, width_out = y.shape
269     self.update_activation_history(y)
270     temporal_winners = self.compute_temporal_winners(y)
271     y_winners = temporal_winners * y
272     y_winners = self.apply_competition(y_winners, batch_size,
273         out_channels)
273     yx = self.compute_yx(x, y_winners)
274     y_sum = y_winners.sum(dim=(0, 2, 3)).view(self.in_channels,
275         1, 1, 1)
275     update = yx - y_sum * weight
276     return update
277
278 def update_adaptive_threshold(self, x, y, weight):
279     batch_size, out_channels, height_out, width_out = y.shape
280     similarities = F.conv2d(x, weight, stride=self.stride,
281         padding=self.padding, groups=self.groups)
281     similarities = similarities / (torch.norm(weight.view(
282         out_channels, -1), dim=1).view(1, -1, 1, 1) + 1e-10)
283     threshold = self.compute_adaptive_threshold(similarities)
284     winners = (similarities > threshold).float()
285     y_winners = winners * similarities
286     y_winners = self.apply_competition(y_winners, batch_size,
287         out_channels)
287     yx = self.compute_yx(x, y_winners)
288     y_sum = y_winners.sum(dim=(0, 2, 3)).view(self.in_channels,
289         1, 1, 1)
290     update = yx - y_sum * weight
291     return update
292
291 def compute_yx(self, x, y):
292     yx = F.conv2d(x.transpose(0, 1), y.transpose(0, 1), padding
293         =0,
294             stride=self.dilation, dilation=self.stride).
295                 transpose(0, 1)
295     yx = yx.diagonal(dim1=0, dim2=1).permute(2, 0, 1).unsqueeze
296         (1)
297     return yx
298
297 def compute_wta_mask(self, y):
298     batch_size, out_channels, height_out, width_out = y.shape
299     # WTA competition within each channel
300     y_flat = y.view(batch_size, out_channels, -1)
301     win_neurons = torch.argmax(y_flat, dim=2)
302     wta_mask = F.one_hot(win_neurons, num_classes=height_out *
303         width_out).float()
304     return wta_mask.view(batch_size, out_channels, height_out,
304         width_out)

```

```

305     def compute_softwta_activations(self, y):
306         # Competition and anti-Hebbian learning for y_depthwise
307         batch_size, in_channels, height_depthwise, width_depthwise =
308             y.shape
309         # Reshape to apply softmax within each channel
310         y_depthwise_reshaped = y.view(batch_size, in_channels, -1)
311         # Apply softmax within each channel
312         flat_softwta_actvns_depthwise = torch.softmax(self.t_invert
313             * y_depthwise_reshaped, dim=2)
314         # Turn all postsynaptic activations into anti-Hebbian
315         flat_softwta_actvns_depthwise = -
316             flat_softwta_actvns_depthwise
317         # Find winners within each channel
318         win_neurons_depthwise = torch.argmax(y_depthwise_reshaped,
319             dim=2)
320         # Create a mask to flip the sign of winning neurons
321         mask = torch.zeros_like(flat_softwta_actvns_depthwise)
322         mask.scatter_(2, win_neurons_depthwise.unsqueeze(2), 1)
323         # Flip the sign of winning neurons
324         flat_softwta_actvns_depthwise =
325             flat_softwta_actvns_depthwise * (1 - 2 * mask)
326         # Reshape back to original shape
327         return flat_softwta_actvns_depthwise.view(batch_size,
328             in_channels, height_depthwise,
329             width_depthwise)

330     def update_activation_history(self, y):
331         if self.activation_history is None:
332             self.activation_history = y.detach().clone()
333         else:
334             self.activation_history = torch.cat([self.
335                 activation_history, y.detach()], dim=0)
336         if self.activation_history.size(0) > self.
337             temporal_window:
338             self.activation_history = self.activation_history[-
339                 self.temporal_window:]

340     def compute_temporal_winners(self, y):
341         batch_size, out_channels, height_out, width_out = y.shape
342         history_spatial = self.activation_history.view(-1,
343             out_channels, height_out, width_out)
344         median_activations = torch.median(history_spatial, dim=0)[0]
345         temporal_threshold = torch.mean(median_activations, dim=(1,
346             2), keepdim=True)
347         return (median_activations > temporal_threshold).float()

348     def compute_adaptive_threshold(self, similarities):
349         mean_sim = similarities.mean(dim=(2, 3), keepdim=True)
350         std_sim = similarities.std(dim=(2, 3), keepdim=True)
351         return mean_sim + self.competition_k * std_sim

352     def apply_competition(self, y, batch_size, out_channels):
353         if self.mode in [self.MODE_TEMPORAL_COMPETITION, self.
354             MODE_ADAPTIVE_THRESHOLD]:
355             if self.competition_type == 'hard':
356                 y = y.view(batch_size, out_channels, -1)

```

```

349         top_k_indices = torch.topk(y, self.top_k, dim=2,
350             largest=True, sorted=False).indices
351         y_compete = torch.zeros_like(y)
352         y_compete.scatter_(2, top_k_indices, y.gather(2,
353             top_k_indices))
354         return y_compete.view_as(y)
355     elif self.competition_type == 'soft':
356         return torch.softmax(self.t_invert * y.view(
357             batch_size, out_channels, -1), dim=2).view_as(y)
358     return y
359
360     @torch.no_grad()
361     def local_update(self):
362         new_weight = self.weight + 0.1 * self.alpha * self.delta_w
363         self.weight.copy_(new_weight.abs())
364         self.delta_w.zero_()

```

A.2.6 Model_hebb.py

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4
5 from hebb import HebbianConv2d
6 # To use Dale Principle, manuallly change this line
7 # from hebb_abs import HebbianConv2d
8
9 import matplotlib.pyplot as plt
10 import numpy as np
11 import wandb
12 import seaborn as sns
13
14 """
15 File handles the creation of Hebbian models
16 """
17
18 torch.manual_seed(0)
19 DEFAULT_HEBB_PARAMS = {'mode': HebbianConv2d.MODE_SOFTWTA, 'w_nrm':
    True, 'k': 50, 'act': nn.Identity(), 'alpha': 1.}
20
21 # Activation function, as described in https://github.com/
22 # NeuromorphicComputing/SoftHebb
22 class Triangle(nn.Module):
23     def __init__(self, power: float = 1, inplace: bool = True):
24         super(Triangle, self).__init__()
25         self.inplace = inplace
26         self.power = power
27
28     def forward(self, input: torch.Tensor) -> torch.Tensor:
29         input = input - torch.mean(input.data, axis=1, keepdims=True
30             )
30         return F.relu(input, inplace=self.inplace) ** self.power
31
32 class Net_Hebbian(nn.Module):
33     def __init__(self, hebb_params=None, version="softhebb"):
34         super(Net_Hebbian, self).__init__()
35         self.hebb_params = hebb_params or DEFAULT_HEBB_PARAMS
36         self.version = version

```

```

37         self._build_network()
38
39     def _build_network(self):
40         if self.version == "softhebb":
41             self._build_softhebb_network()
42         elif self.version == "hardhebb":
43             self._build_hardhebb_network()
44         elif self.version == "lagani":
45             self._build_lagani_network()
46         elif self.version == "lagani_short":
47             self._build_lagani_short_network()
48         else:
49             raise ValueError(f"Unknown version: {self.version}")
50
51     # Architecture equivalent to SoftHebb research
52     def _build_softhebb_network(self):
53         # Layer 1
54         self.bn1 = nn.BatchNorm2d(3, affine=False)
55         self.conv1 = HebbianConv2d(in_channels=3, out_channels=96,
56             kernel_size=5, stride=1, **self.hebb_params,
57             padding=2, t_invert=1)
58         self.pool1 = nn.MaxPool2d(kernel_size=4, stride=2, padding
59             =1)
60         self.activ1 = Triangle(power=0.7)
61
62         # Layer 2
63         self.bn2 = nn.BatchNorm2d(96, affine=False)
64         self.conv2 = HebbianConv2d(in_channels=96, out_channels=384,
65             kernel_size=3, stride=1, **self.hebb_params,
66             t_invert=0.65, padding=1)
67         self.pool2 = nn.MaxPool2d(kernel_size=4, stride=2, padding
68             =1)
69         self.activ2 = Triangle(power=1.4)
70
71         # Layer 3
72         self.bn3 = nn.BatchNorm2d(384, affine=False)
73         self.conv3 = HebbianConv2d(in_channels=384, out_channels
74             =1536, kernel_size=3, stride=1, **self.hebb_params,
75             t_invert=0.25, padding=1)
76         self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
77             =0)
78         self.activ3 = Triangle(power=1.)
79
80         # Output layers
81         self.flatten = nn.Flatten()
82         self.fc1 = nn.Linear(24576, 10)
83         self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
84             24576)
85         self.dropout = nn.Dropout(0.5)
86
87     # Architecture equivalent to SoftHebb research, but without
88     # padding
89     def _build_hardhebb_network(self): # Hybrid: Similar to Lagani
90         but with softhebb number of filters and layers
91         # Layer 1
92         self.bn1 = nn.BatchNorm2d(3, affine=False)
93         self.conv1 = HebbianConv2d(in_channels=3, out_channels=96,
94             kernel_size=5, stride=1, **self.hebb_params,

```

```

85                                     padding=0, t_invert=1)
86     self.pool1 = nn.MaxPool2d(kernel_size=2, stride=2, padding
87                               =0)
88     self.activ1 = Triangle(power=0.7)
89
90     # Layer 2
91     self.bn2 = nn.BatchNorm2d(96, affine=False)
92     self.conv2 = HebbianConv2d(in_channels=96, out_channels=384,
93                               kernel_size=3, stride=1, **self.hebb_params,
94                               t_invert=0.65, padding=0)
95     self.activ2 = Triangle(power=1.4)
96
97     # Layer 3
98     self.bn3 = nn.BatchNorm2d(384, affine=False)
99     self.conv3 = HebbianConv2d(in_channels=384, out_channels
100                               =1536, kernel_size=3, stride=1,
101                               **self.hebb_params, t_invert
102                               =0.25, padding=0)
103     self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
104                               =0)
105     self.activ3 = Triangle(power=1.)
106
107     # Output layers
108     self.flatten = nn.Flatten()
109     self.fc1 = nn.Linear(38400, 10)
110     self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
111                               38400)
112     self.dropout = nn.Dropout(0.5)
113
114     # Architecture equivalent to Lagani 3-layer research
115     def _build_lagani_short_network(self):
116         # Layer 1
117         self.bn1 = nn.BatchNorm2d(3, affine=False)
118         self.conv1 = HebbianConv2d(in_channels=3, out_channels=96,
119                                   kernel_size=5, stride=1, **self.hebb_params,
120                                   padding=0, t_invert=1)
121         self.pool1 = nn.MaxPool2d(kernel_size=2, stride=2, padding
122                               =0)
123         self.activ1 = Triangle(power=1.)
124
125         # Layer 2
126         self.bn2 = nn.BatchNorm2d(96, affine=False)
127         self.conv2 = HebbianConv2d(in_channels=96, out_channels=128,
128                                   kernel_size=3, stride=1, **self.hebb_params,
129                                   t_invert=0.65, padding=0)
130         # self.pool2 = nn.MaxPool2d(kernel_size=2, stride=2, padding
131                               =0)
132         self.activ2 = Triangle(power=1.)
133
134         # Layer 3
135         self.bn3 = nn.BatchNorm2d(128, affine=False)
136         self.conv3 = HebbianConv2d(in_channels=128, out_channels
137                                   =192, kernel_size=3, stride=1, **self.hebb_params,
138                                   t_invert=0.25, padding=0)
139         self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
140                               =0)
141         self.activ3 = Triangle(power=1.)

```

```

131     # Output layers
132     self.flatten = nn.Flatten()
133     self.fc1 = nn.Linear(4800, 10)
134     self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
135                                         4800)
136     self.dropout = nn.Dropout(0.5)
137
138     # Architecture equivalent to Lagani 4-layer research
139     def _build_lagani_network(self):
140         # Layer 1
141         self.bn1 = nn.BatchNorm2d(3, affine=False)
142         self.conv1 = HebbianConv2d(in_channels=3, out_channels=96,
143                               kernel_size=5, stride=1, **self.hebb_params,
144                               padding=0, t_invert=1)
145         self.pool1 = nn.MaxPool2d(kernel_size=2, stride=2, padding
146                               =0)
147         self.activ1 = Triangle(power=1.)
148
149         # Layer 2
150         self.bn2 = nn.BatchNorm2d(96, affine=False)
151         self.conv2 = HebbianConv2d(in_channels=96, out_channels=128,
152                               kernel_size=3, stride=1, **self.hebb_params,
153                               t_invert=0.65, padding=0)
154         self.activ2 = Triangle(power=1.)
155
156         # Layer 3
157         self.bn3 = nn.BatchNorm2d(128, affine=False)
158         self.conv3 = HebbianConv2d(in_channels=128, out_channels
159                               =192, kernel_size=3, stride=1, **self.hebb_params,
160                               t_invert=0.25, padding=0)
161         self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
162                               =0)
163         self.activ3 = Triangle(power=1.)
164
165         # Layer 4
166         self.bn4 = nn.BatchNorm2d(192, affine=False)
167         self.conv4 = HebbianConv2d(in_channels=192, out_channels
168                               =256, kernel_size=3, stride=1, **self.hebb_params,
169                               t_invert=0.25, padding=0)
170         self.activ4 = Triangle(power=1.)
171
172         # Output layers
173         self.flatten = nn.Flatten()
174         self.fc1 = nn.Linear(2304, 10)
175         self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
176                                         2304)
177         self.dropout = nn.Dropout(0.5)
178
179     def forward_features(self, x):
180         x = self.pool1(self.activ1(self.conv1(self.bn1(x))))
181         return x
182
183     def features_extract(self, x):
184         x = self.forward_features(x)
185         if self.version == "lagani":
186             x = self.activ2(self.conv2(self.bn2(x)))
187             x = self.pool3(self.activ3(self.conv3(self.bn3(x))))
188             x = self.activ4(self.conv4(self.bn4(x)))

```

```

181     elif self.version == "lagani_short":
182         x = self.activ2(self.conv2(self.bn2(x)))
183         x = self.pool3(self.activ3(self.conv3(self.bn3(x))))
184     elif self.version == "hardhebb":
185         x = self.activ2(self.conv2(self.bn2(x)))
186         x = self.pool3(self.activ3(self.conv3(self.bn3(x))))
187     elif self.version == "softhebb":
188         x = self.pool2(self.activ2(self.conv2(self.bn2(x))))
189         x = self.pool3(self.activ3(self.conv3(self.bn3(x))))
190     return x
191
192 def forward(self, x):
193     x = self.features_extract(x)
194     x = self.flatten(x)
195     x = self.fc1(self.dropout(x))
196     return x
197
198 # Plot neurons/filter of a target layer
199 def plot_grid(self, tensor, path, num_rows=5, num_cols=5,
200               layer_name=""):
201     # Ensure we're working with the first 25 filters (or less if
202     # there are fewer)
203     excitatory = tensor[:20]
204     inhibitory = tensor[-5:]
205     # Symmetric normalization for excitatory weights
206     max_abs_exc = torch.max(torch.abs(excitatory))
207     norm_exc = excitatory / (max_abs_exc + 1e-8)
208     # Symmetric normalization for inhibitory weights
209     max_abs_inh = torch.max(torch.abs(inhibitory))
210     norm_inh = inhibitory / (max_abs_inh + 1e-8)
211     tensor = torch.cat((norm_exc, norm_inh))
212     # Normalize the tensor
213     # Move to CPU and convert to numpy
214     tensor = tensor.cpu().detach().numpy()
215
216     if tensor.shape[2] == 1 and tensor.shape[3] == 1: # 1x1
217         convolution case
218         out_channels, in_channels = tensor.shape[:2]
219         fig = plt.figure(figsize=(14, 10))
220         # Create a gridspec for the layout
221         gs = fig.add_gridspec(2, 2, width_ratios=[20, 1],
222                               height_ratios=[1, 3],
223                               left=0.1, right=0.9, bottom=0.1,
224                               top=0.9, wspace=0.05, hspace
225                               =0.2)
226         ax1 = fig.add_subplot(gs[0, 0])
227         ax2 = fig.add_subplot(gs[1, 0])
228         cbar_ax = fig.add_subplot(gs[:, 1])
229         # Bar plot for average weights per filter
230         avg_weights = tensor.mean(axis=(1, 2, 3))
231         norm = plt.Normalize(vmin=avg_weights.min(), vmax=
232                               avg_weights.max())
233         im1 = ax1.bar(range(out_channels), avg_weights, color=
234                       plt.cm.RdYlGn(norm(avg_weights)))
235         ax1.set_xlabel('Filter Index')
236         ax1.set_ylabel('Average Weight')
237         ax1.set_title(f'Average Weights for 1x1 Kernels in {layer_name}')

```

```

230     ax1.axhline(y=0, color='black', linestyle='--',
231                 linewidth=0.5)
232     # Heatmap for detailed weight distribution
233     im2 = ax2.imshow(tensor.reshape(out_channels,
234                        in_channels), cmap='RdYlGn', aspect='auto', norm=norm
235                         )
236     ax2.set_xlabel('Input Channel')
237     ax2.set_ylabel('Output Channel (Filter)')
238     ax2.set_title('Detailed Weight Distribution')
239     # Add colorbar to the right of both subplots
240     fig.colorbar(im2, cax=cbar_ax, label='Normalized Weight
241                   Value')
242
243     else:
244         fig, axes = plt.subplots(num_rows, num_cols, figsize
245             =(15, 10))
246         fig.suptitle(f'First 25 Filters of {layer_name}')
247         for i, ax in enumerate(axes.flat):
248             if i < tensor.shape[0]:
249                 filter_img = tensor[i]
250                 # Handle different filter shapes
251                 if filter_img.shape[0] == 3: # RGB filter (3, H
252                     , W)
253                     filter_img = np.transpose(filter_img, (1, 2,
254                                         0))
255                     filter_img = (filter_img - filter_img.min())
256                     / (filter_img.max() - filter_img.min() +
257                         1e-8)
258                 elif filter_img.shape[0] == 1: # Grayscale
259                     filter (1, H, W)
260                     filter_img = filter_img.squeeze()
261                 else: # Multi-channel filter (C, H, W), take
262                     mean across channels
263                     filter_img = np.mean(filter_img, axis=0)
264                     ax.imshow(filter_img, cmap='viridis' if
265                               filter_img.ndim == 2 else None)
266                     ax.set_title(f'Filter {i + 1}')
267                     ax.axis('off')
268
269             if path:
270                 fig.savefig(path, bbox_inches='tight')
271                 wandb.log({f'{layer_name} filters': wandb.Image(fig)})
272             plt.close(fig)
273
274
275     def visualize_filters(self, layer_name='conv1', save_path=None):
276         weights = getattr(self, layer_name).weight.data
277         self.plot_grid(weights, save_path, layer_name=layer_name)

```

A.2.7 Model_depthwise

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4
5 from hebb import HebbianConv2d
6 from hebb_depthwise import HebbianDepthConv2d
7

```

```

8 # from hebb_abs import HebbianConv2d
9 # from hebb_abs_depthwise import HebbianDepthConv2d
10
11 import matplotlib.pyplot as plt
12 import numpy as np
13 import wandb
14 import seaborn as sns
15
16 torch.manual_seed(0)
17 DEFAULT_HEBB_PARAMS = { 'mode': HebbianConv2d.MODE_SOFTWTA, 'w_nrm': True, 'k': 50, 'act': nn.Identity(), 'alpha': 1.}
18
19
20 class Triangle(nn.Module):
21     def __init__(self, power: float = 1, inplace: bool = True):
22         super(Triangle, self).__init__()
23         self.inplace = inplace
24         self.power = power
25
26     def forward(self, input: torch.Tensor) -> torch.Tensor:
27         input = input - torch.mean(input.data, axis=1, keepdims=True)
28         return F.relu(input, inplace=self.inplace) ** self.power
29
30 class Net_Depthwise(nn.Module):
31     def __init__(self, hebb_params=None, version="softhebb"):
32         super(Net_Depthwise, self).__init__()
33         self.hebb_params = hebb_params or DEFAULT_HEBB_PARAMS
34         self.version = version
35         self._build_network()
36
37     def _build_network(self):
38         if self.version == "softhebb":
39             self._build_softhebb_network()
40         elif self.version == "hardhebb":
41             self._build_hardhebb_network()
42         elif self.version == "lagani":
43             self._build_lagani_network()
44         else:
45             raise ValueError(f"Unknown version: {self.version}")
46
47     def _build_softhebb_network(self):
48         # Layer 1
49         self.bn1 = nn.BatchNorm2d(3, affine=False)
50         self.conv1 = HebbianConv2d(in_channels=3, out_channels=96,
51             kernel_size=5, stride=1, **self.hebb_params,
52             padding=2, t_invert=1)
53         self.pool1 = nn.MaxPool2d(kernel_size=4, stride=2, padding
54             =1)
55         self.activ1 = Triangle(power=0.7)
56
57         # Layer 2
58         self.bn2 = nn.BatchNorm2d(96, affine=False)
59         self.conv2 = HebbianDepthConv2d(in_channels=96, out_channels
60             =96, kernel_size=3, stride=1, **self.hebb_params,
61             t_invert=0.65, padding=1)
62         self.bn_point2 = nn.BatchNorm2d(96, affine=False)

```

```

60         self.conv_point2 = HebbianConv2d(in_channels=96,
61                                         out_channels=384, kernel_size=1, stride=1, **self.
62                                         hebb_params,
63                                         t_invert=0.65, padding=0)
64         self.pool2 = nn.MaxPool2d(kernel_size=4, stride=2, padding
65                                     =1)
66         self.activ2 = Triangle(power=1.4)
67
68     # Layer 3
69     self.bn3 = nn.BatchNorm2d(384, affine=False)
70     self.conv3 = HebbianDepthConv2d(in_channels=384,
71                                     out_channels=384, kernel_size=3, stride=1, **self.
72                                     hebb_params,
73                                     t_invert=0.25, padding=1)
74     self.bn_point3 = nn.BatchNorm2d(384, affine=False)
75     self.conv_point3 = HebbianConv2d(in_channels=384,
76                                     out_channels=1536, kernel_size=1, stride=1, **self.
77                                     hebb_params,
78                                     t_invert=0.25, padding=0)
79     self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
80                               =0)
81     self.activ3 = Triangle(power=1.)
82
83     # Output layers
84     self.flatten = nn.Flatten()
85     self.fc1 = nn.Linear(24576, 10)
86     self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
87                                              24576)
88     self.dropout = nn.Dropout(0.5)
89
90     def _build_hardhebb_network(self):
91         # Layer 1
92         self.bn1 = nn.BatchNorm2d(3, affine=False)
93         self.conv1 = HebbianConv2d(in_channels=3, out_channels=96,
94                                   kernel_size=5, stride=1, **self.hebb_params,
95                                   padding=0, t_invert=1)
96         self.pool1 = nn.MaxPool2d(kernel_size=2, stride=2, padding
97                                   =0)
98         self.activ1 = Triangle(power=0.7)
99
100        # Layer 2
101        self.bn2 = nn.BatchNorm2d(96, affine=False)
102        self.conv2 = HebbianDepthConv2d(in_channels=96, out_channels
103                                      =96, kernel_size=3, stride=1, **self.hebb_params,
104                                      t_invert=0.65, padding=0)
105        self.bn_point2 = nn.BatchNorm2d(96, affine=False)
106        self.conv_point2 = HebbianConv2d(in_channels=96,
107                                       out_channels=384, kernel_size=1, stride=1,
108                                       **self.hebb_params,
109                                       t_invert=0.65, padding
110                                       =0)
111
112        self.activ2 = Triangle(power=1.4)
113
114        # Layer 3
115        self.bn3 = nn.BatchNorm2d(384, affine=False)
116        self.conv3 = HebbianDepthConv2d(in_channels=384,
117                                      out_channels=384, kernel_size=3, stride=1,

```

```

101                                     **self.hebb_params, t_invert
102                                         =0.25, padding=0)
103     self.bn_point3 = nn.BatchNorm2d(384, affine=False)
104     self.conv_point3 = HebbianConv2d(in_channels=384,
105                                         out_channels=1536, kernel_size=1, stride=1,
106                                         **self.hebb_params,
107                                         t_invert=0.25, padding
108                                         =0)
109     self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
110                                         =0)
111     self.activ3 = Triangle(power=1.)
112
113     # Output layers
114     self.flatten = nn.Flatten()
115     self.fc1 = nn.Linear(38400, 10)
116     self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
117                                         38400)
118     self.dropout = nn.Dropout(0.5)
119
120     def _build_lagani_network(self):
121         # Layer 1
122         self.bn1 = nn.BatchNorm2d(3, affine=False)
123         self.conv1 = HebbianConv2d(in_channels=3, out_channels=96,
124                                         kernel_size=5, stride=1, **self.hebb_params,
125                                         padding=0, t_invert=1)
126         self.pool1 = nn.MaxPool2d(kernel_size=2, stride=2, padding
127                                         =0)
128         self.activ1 = Triangle(power=1.)
129
130         # Layer 2
131         self.bn2 = nn.BatchNorm2d(96, affine=False)
132         self.conv2 = HebbianDepthConv2d(in_channels=96, out_channels
133                                         =96, kernel_size=3, stride=1, **self.hebb_params,
134                                         t_invert=0.65, padding=0)
135         self.bn_point2 = nn.BatchNorm2d(96, affine=False)
136         self.conv_point2 = HebbianConv2d(in_channels=96,
137                                         out_channels=128, kernel_size=1, stride=1, **self.
138                                         hebb_params,
139                                         t_invert=0.65, padding=0)
140         self.pool2 = nn.MaxPool2d(kernel_size=2, stride=2, padding
141                                         =0)
142         self.activ2 = Triangle(power=1.)
143
144         # Layer 3
145         self.bn3 = nn.BatchNorm2d(128, affine=False)
146         self.conv3 = HebbianDepthConv2d(in_channels=128,
147                                         out_channels=128, kernel_size=3, stride=1, **self.
148                                         hebb_params,
149                                         t_invert=0.25, padding=0)
150         self.bn_point3 = nn.BatchNorm2d(128, affine=False)
151         self.conv_point3 = HebbianConv2d(in_channels=128,
152                                         out_channels=192, kernel_size=1, stride=1, **self.
153                                         hebb_params,
154                                         t_invert=0.25, padding=0)
155         self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
156                                         =0)
157         self.activ3 = Triangle(power=1.)

```

```

142     # Layer 4
143     self.bn4 = nn.BatchNorm2d(192, affine=False)
144     self.conv4 = HebbianDepthConv2d(in_channels=192,
145         out_channels=192, kernel_size=3, stride=1, **self.
146             hebb_params,
147                 t_invert=0.25, padding=0)
148     self.bn_point4 = nn.BatchNorm2d(192, affine=False)
149     self.conv_point4 = HebbianConv2d(in_channels=192,
150         out_channels=256, kernel_size=1, stride=1, **self.
151             hebb_params,
152                 t_invert=0.25, padding=0)
153     self.activ4 = Triangle(power=1.)
154
155     # Output layers
156     self.flatten = nn.Flatten()
157     self.fc1 = nn.Linear(1728, 10)
158     self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
159         1728)
160     self.dropout = nn.Dropout(0.5)
161
162     def forward_features(self, x):
163         x = self.pool1(self.activ1(self.conv1(self.bn1(x))))
164         return x
165
166     def features_extract(self, x):
167         x = self.forward_features(x)
168         if self.version == "lagani":
169             x = self.pool2(self.activ2(self.conv_point2(self.
170                 bn_point2(self.conv2(self.bn2(x))))))
171             x = self.pool3(self.activ3(self.conv_point3(self.
172                 bn_point3(self.conv3(self.bn3(x))))))
173             x = self.activ4(self.conv_point4(self.bn_point4(self.
174                 conv4(self.bn4(x)))))

175             elif self.version == "hardhebb":
176                 x = self.activ2(self.conv_point2(self.bn_point2(self.
177                     conv2(self.bn2(x)))))
178                 x = self.pool3(self.activ3(self.conv_point3(self.
179                     bn_point3(self.conv3(self.bn3(x))))))

180             elif self.version == "softhebb":
181                 x = self.pool2(self.activ2(self.conv_point2(self.
182                     bn_point2(self.conv2(self.bn2(x))))))
183                 x = self.pool3(self.activ3(self.conv_point3(self.
184                     bn_point3(self.conv3(self.bn3(x))))))

185             return x
186
187     def forward(self, x):
188         x = self.features_extract(x)
189         x = self.flatten(x)
190         x = self.fc1(self.dropout(x))
191         return x
192
193     def plot_grid(self, tensor, path, num_rows=5, num_cols=5,
194         layer_name=""):
195         # Ensure we're working with the first 25 filters (or less if
196             there are fewer)
197         excitatory = tensor[:20]
198         inhibitory = tensor[-5:]
199         # Symmetric normalization for excitatory weights

```

```

186     max_abs_exc = torch.max(torch.abs(excitatory))
187     norm_exc = excitatory / (max_abs_exc + 1e-8)
188     # Symmetric normalization for inhibitory weights
189     max_abs_inh = torch.max(torch.abs(inhibitory))
190     norm_inh = inhibitory / (max_abs_inh + 1e-8)
191     tensor = torch.cat((norm_exc, norm_inh))
192     # Move to CPU and convert to numpy
193     tensor = tensor.cpu().detach().numpy()
194
195     if tensor.shape[2] == 1 and tensor.shape[3] == 1: # 1x1
196         convolution case
197         out_channels, in_channels = tensor.shape[:2]
198         fig = plt.figure(figsize=(14, 10))
199         # Create a gridspec for the layout
200         gs = fig.add_gridspec(2, 2, width_ratios=[20, 1],
201                               height_ratios=[1, 3],
202                               left=0.1, right=0.9, bottom=0.1,
203                               top=0.9, wspace=0.05, hspace
204                               =0.2)
205         ax1 = fig.add_subplot(gs[0, 0])
206         ax2 = fig.add_subplot(gs[1, 0])
207         cbar_ax = fig.add_subplot(gs[:, 1])
208         # Bar plot for average weights per filter
209         avg_weights = tensor.mean(axis=(1, 2, 3))
210         norm = plt.Normalize(vmin=avg_weights.min(), vmax=
211                               avg_weights.max())
212         im1 = ax1.bar(range(out_channels), avg_weights, color=
213                       plt.cm.RdYlGn(norm(avg_weights)))
214         ax1.set_xlabel('Filter Index')
215         ax1.set_ylabel('Average Weight')
216         ax1.set_title(f'Average Weights for 1x1 Kernels in {
217                         layer_name}')
218         ax1.axhline(y=0, color='black', linestyle='--',
219                     linewidth=0.5)
220         # Heatmap for detailed weight distribution
221         im2 = ax2.imshow(tensor.reshape(out_channels,
222                                     in_channels), cmap='RdYlGn', aspect='auto', norm=norm
223                                     )
224         ax2.set_xlabel('Input Channel')
225         ax2.set_ylabel('Output Channel (Filter)')
226         ax2.set_title('Detailed Weight Distribution')
227         # Add colorbar to the right of both subplots
228         fig.colorbar(im2, cax=cbar_ax, label='Normalized Weight
Value')
229
230     else:
231         fig, axes = plt.subplots(num_rows, num_cols, figsize
232                                 =(15, 10))
233         fig.suptitle(f'First 25 Filters of {layer_name}')
234         for i, ax in enumerate(axes.flat):
235             if i < tensor.shape[0]:
236                 filter_img = tensor[i]
237                 # Handle different filter shapes
238                 if filter_img.shape[0] == 3: # RGB filter (3, H
239                     , W)
240                     filter_img = np.transpose(filter_img, (1, 2,
241                                               0))

```

```

229             filter_img = (filter_img - filter_img.min())
230             / (filter_img.max() - filter_img.min() + 1e-8)
231         elif filter_img.shape[0] == 1: # Grayscale
232             filter = (1, H, W)
233             filter_img = filter_img.squeeze()
234         else: # Multi-channel filter (C, H, W), take
235             mean across channels
236             filter_img = np.mean(filter_img, axis=0)
237             ax.imshow(filter_img, cmap='viridis' if
238             filter_img.ndim == 2 else None)
239             ax.set_title(f'Filter {i + 1}')
240             ax.axis('off')
241
242     if path:
243         fig.savefig(path, bbox_inches='tight')
244         wandb.log({f'{layer_name} filters': wandb.Image(fig)})
245         plt.close(fig)
246
247     def visualize_filters(self, layer_name='conv1', save_path=None):
248         weights = getattr(self, layer_name).weight.data
249         self.plot_grid(weights, save_path, layer_name=layer_name)

```

A.2.8 Model_residual.py

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4
5 from hebb import HebbianConv2d
6 from hebb_depthwise import HebbianDepthConv2d
7
8 import matplotlib.pyplot as plt
9 import numpy as np
10 import wandb
11 import seaborn as sns
12
13 """
14 File handles the creation of Hebbian Residual models
15 Similar to mode_hebb.py, refer to this file for further explanations
16 """
17
18 torch.manual_seed(0)
19 default_hebb_params = {'mode': HebbianConv2d.MODE_SOFTWTA, 'w_nrm': True,
20                         'k': 50, 'act': nn.Identity(), 'alpha': 1.}
21
22 class Triangle(nn.Module):
23     def __init__(self, power: float = 1, inplace: bool = True):
24         super(Triangle, self).__init__()
25         self.inplace = inplace
26         self.power = power
27
28     def forward(self, input: torch.Tensor) -> torch.Tensor:
29         input = input - torch.mean(input.data, axis=1, keepdims=True)
30         return F.relu(input, inplace=self.inplace) ** self.power

```

```

32 class LongSkipConnection(nn.Module):
33     def __init__(self, in_channels, out_channels, spatial_change=1,
34                  hebb_params=None):
35         super(LongSkipConnection, self).__init__()
36         if hebb_params is None:
37             hebb_params = default_hebb_params
38
39         self.conv = HebbianConv2d(in_channels, out_channels,
40                                kernel_size=1, stride=1, **hebb_params)
41         self.bn = nn.BatchNorm2d(out_channels, affine=False)
42         self.spatial_change = spatial_change
43
44     def forward(self, x):
45         x = self.conv(x)
46         x = self.bn(x)
47         if self.spatial_change != 1:
48             x = F.adaptive_avg_pool2d(x, output_size=x.size()[2:] //
49                                       self.spatial_change)
50
51     return x
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
    class HebbianResidualBlock(nn.Module):
        def __init__(self, in_channels, out_channels, kernel_size,
                     stride=1, hebb_params=None, t_invert=1., expansion_factor=6,
                     act=1.):
            super(HebbianResidualBlock, self).__init__()
            if hebb_params is None:
                hebb_params = default_hebb_params
            # Calculate padding to maintain spatial dimensions
            # Doubt regarding when additional padding is required
            padding = (kernel_size - 1) // 2
            hidden_dim = in_channels * expansion_factor
            self.bn1 = nn.BatchNorm2d(in_channels, affine=False)
            self.conv1 = HebbianConv2d(in_channels, hidden_dim,
                                      kernel_size=1, stride=1, **hebb_params, t_invert=t_invert
                                      ,
                                      padding=0)
            self.bn2 = nn.BatchNorm2d(hidden_dim, affine=False)
            self.conv2 = HebbianDepthConv2d(hidden_dim, hidden_dim,
                                           kernel_size, stride=1, **hebb_params,
                                           t_invert=t_invert, padding=
                                           padding)
            self.bn3 = nn.BatchNorm2d(hidden_dim, affine=False)
            self.conv3 = HebbianConv2d(hidden_dim, out_channels,
                                      kernel_size=1, stride=1, **hebb_params, t_invert=t_invert
                                      ,
                                      padding=0)
            self.activ = Triangle(power=act)
            self.shortcut = nn.Sequential()
            if stride != 1 or in_channels != out_channels:
                self.shortcut = nn.Sequential(
                    nn.BatchNorm2d(in_channels, affine=False),
                    HebbianConv2d(in_channels, out_channels, kernel_size=1,
                                  stride=1, **hebb_params, padding=0)
                )
        def forward(self, x):
            residual = x

```

```

78         out = self.activ(self.conv1(self.bn1(x)))
79         out = self.activ(self.conv2(self.bn2(out)))
80         out = self.conv3(self.bn3(out))
81         out += self.shortcut(residual)
82         return self.activ(out)
83
84 class Net_Depthwise_Residual(nn.Module):
85     def __init__(self, hebb_params=None):
86         super(Net_Depthwise_Residual, self).__init__()
87
88         if hebb_params is None:
89             hebb_params = default_hebb_params
90
91         self.bn1 = nn.BatchNorm2d(3, affine=False)
92         self.conv1 = HebbianConv2d(in_channels=3, out_channels=96,
93             kernel_size=5, stride=1, **hebb_params, padding=2)
94         self.pool1 = nn.MaxPool2d(kernel_size=4, stride=2, padding
95             =1)
96         self.activ1 = Triangle(power=0.7)
97
98         self.res1 = HebbianResidualBlock(96, 384, kernel_size=3,
99             expansion_factor=4, hebb_params=hebb_params, t_invert
100             =0.65, act=1.4)
101        self.pool2 = nn.MaxPool2d(kernel_size=4, stride=2, padding
102             =1)
103
104        self.res2 = HebbianResidualBlock(384, 1536, kernel_size=3,
105            expansion_factor=4, hebb_params=hebb_params, t_invert
106            =0.25)
107        self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
108            =0)
109
110        self.flatten = nn.Flatten()
111        self.fc1 = nn.Linear(24576, 10) # Adjust this based on your
112            input size
113        self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
114            24576)
115        self.dropout = nn.Dropout(0.5)
116
117    def forward_features(self, x):
118        x = self.pool1(self.activ1(self.conv1(self.bn1(x))))
119        return x
120
121    def features_extract(self, x):
122        x = self.forward_features(x)
123        x = self.pool2(self.res1(x))
124        x = self.pool3(self.res2(x))
125        return x
126
127    def forward(self, x):
128        x = self.features_extract(x)
129        x = self.flatten(x)
130        x = self.fc1(self.dropout(x))
131        return x
132
133    # Plot neurons/filter of a target layer
134    def plot_grid(self, tensor, path, num_rows=5, num_cols=5,
135        layer_name=""):

```

```

125     # Ensure we're working with the first 25 filters (or less if
126     # there are fewer)
127     excitatory = tensor[:20]
128     inhibitory = tensor[-5:]
129     # Symmetric normalization for excitatory weights
130     max_abs_exc = torch.max(torch.abs(excitatory))
131     norm_exc = excitatory / (max_abs_exc + 1e-8)
132     # Symmetric normalization for inhibitory weights
133     max_abs_inh = torch.max(torch.abs(inhibitory))
134     norm_inh = inhibitory / (max_abs_inh + 1e-8)
135     tensor = torch.cat((norm_exc, norm_inh))
136     # Normalize the tensor
137     # Move to CPU and convert to numpy
138     tensor = tensor.cpu().detach().numpy()
139
140     if tensor.shape[2] == 1 and tensor.shape[3] == 1: # 1x1
141         convolution case
142         out_channels, in_channels = tensor.shape[:2]
143         fig = plt.figure(figsize=(14, 10))
144         # Create a gridspec for the layout
145         gs = fig.add_gridspec(2, 2, width_ratios=[20, 1],
146                               height_ratios=[1, 3],
147                               left=0.1, right=0.9, bottom=0.1,
148                               top=0.9, wspace=0.05, hspace
149                               =0.2)
150         ax1 = fig.add_subplot(gs[0, 0])
151         ax2 = fig.add_subplot(gs[1, 0])
152         cbar_ax = fig.add_subplot(gs[:, 1])
153         # Bar plot for average weights per filter
154         avg_weights = tensor.mean(axis=(1, 2, 3))
155         norm = plt.Normalize(vmin=avg_weights.min(), vmax=
156                               avg_weights.max())
157         im1 = ax1.bar(range(out_channels), avg_weights, color=
158                       plt.cm.RdYlGn(norm(avg_weights)))
159         ax1.set_xlabel('Filter Index')
160         ax1.set_ylabel('Average Weight')
161         ax1.set_title(f'Average Weights for 1x1 Kernels in {
162                         layer_name}')
163         ax1.axhline(y=0, color='black', linestyle='--',
164                     linewidth=0.5)
165         # Heatmap for detailed weight distribution
166         im2 = ax2.imshow(tensor.reshape(out_channels,
167                                   in_channels), cmap='RdYlGn', aspect='auto', norm=norm
168 )
169         ax2.set_xlabel('Input Channel')
170         ax2.set_ylabel('Output Channel (Filter)')
171         ax2.set_title('Detailed Weight Distribution')
172         # Add colorbar to the right of both subplots
173         fig.colorbar(im2, cax=cbar_ax, label='Normalized Weight
174                      Value')
175
176     else:
177         fig, axes = plt.subplots(num_rows, num_cols, figsize
178                               =(15, 10))
179         fig.suptitle(f'First 25 Filters of {layer_name}')
180         for i, ax in enumerate(axes.flat):
181             if i < tensor.shape[0]:
182                 filter_img = tensor[i]

```

```

170     # Handle different filter shapes
171     if filter_img.shape[0] == 3: # RGB filter (3, H
172         , W)
173         filter_img = np.transpose(filter_img, (1, 2,
174         0))
175         filter_img = (filter_img - filter_img.min())
176         / (filter_img.max() - filter_img.min() +
177         1e-8)
178     elif filter_img.shape[0] == 1: # Grayscale
179         filter (1, H, W)
180         filter_img = filter_img.squeeze()
181     else: # Multi-channel filter (C, H, W), take
182         mean across channels
183         filter_img = np.mean(filter_img, axis=0)
184     ax.imshow(filter_img, cmap='viridis' if
185         filter_img.ndim == 2 else None)
186     ax.set_title(f'Filter {i + 1}')
187     ax.axis('off')
188
189     if path:
190         fig.savefig(path, bbox_inches='tight')
191         wandb.log({f'{layer_name} filters': wandb.Image(fig)})
192         plt.close(fig)
193
194     def visualize_filters(self, layer_name='conv1', save_path=
195         None):
196         weights = getattr(self, layer_name).weight.data
197         self.plot_grid(weights, save_path, layer_name=layer_name
198             )

```

A.2.9 Model_BackProp.py

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 from sklearn.decomposition import PCA
5
6 import matplotlib.pyplot as plt
7 import numpy as np
8 from sklearn.manifold import TSNE
9 from sklearn.preprocessing import StandardScaler
10 from tqdm import tqdm
11 from itertools import islice
12 import umap
13 import wandb
14 import math
15 import seaborn as sns
16
17 """
18 Backpropagation model similar to model_hebb.py
19 TODO: make code more similar to model_hebb.py, with different
20     architectures controlled by an argument
21 """
22 torch.manual_seed(0)
23
24 class Triangle(nn.Module):
25     def __init__(self, power: float = 1, inplace: bool = True):

```

```

26     super(Triangle, self).__init__()
27     self.inplace = inplace
28     self.power = power
29
30     def forward(self, input: torch.Tensor) -> torch.Tensor:
31         input = input - torch.mean(input.data, axis=1, keepdims=True)
32         return F.relu(input, inplace=self.inplace) ** self.power
33
34     def custom_depthwise_conv_init(conv_layer):
35         in_channels = conv_layer.in_channels
36         kernel_size = conv_layer.kernel_size[0] # Assuming square
37         kernels
38         weight_range = 25 / math.sqrt(in_channels * kernel_size *
39             kernel_size)
40         # Depthwise separable weights
41         conv_layer.weight.data = weight_range * torch.randn(in_channels,
42             1, *conv_layer.kernel_size)
43         if conv_layer.bias is not None:
44             conv_layer.bias.data.zero_()
45
46     def custom_pointwise_conv_init(conv_layer):
47         in_channels = conv_layer.in_channels
48         out_channels = conv_layer.out_channels
49         kernel_size = conv_layer.kernel_size[0] # Should be 1 for
50         pointwise
51         groups = conv_layer.groups
52         weight_range = 25 / math.sqrt(in_channels * kernel_size *
53             kernel_size)
54         conv_layer.weight.data = weight_range * torch.randn(out_channels,
55             in_channels // groups, *conv_layer.kernel_size)
56         if conv_layer.bias is not None:
57             conv_layer.bias.data.zero_()
58
59     class Net_Backpropagation_depth(nn.Module):
60         def __init__(self):
61             super(Net_Backpropagation_depth, self).__init__()
62
63             # A single Depthwise convolutional layer
64             self.bn1 = nn.BatchNorm2d(3, affine=False)
65             self.conv1 = nn.Conv2d(in_channels=3, out_channels=96,
66                 kernel_size=(5,5), padding=2, bias=False)
67             self.pool1 = nn.MaxPool2d(kernel_size=4, stride=2, padding
68                 =1)
69             self.activ1 = nn.ReLU()
70
71             self.bn2 = nn.BatchNorm2d(96, affine=False)
72             self.conv2 = nn.Conv2d(in_channels=96, out_channels=96,
73                 kernel_size=(3,3), padding=1, bias=False, groups=96)
74             self.bn_point2 = nn.BatchNorm2d(96, affine=False)
75             self.conv_point2 = nn.Conv2d(in_channels=96, out_channels
76                 =384, kernel_size=(1,1), padding=0, bias=False)
77             self.pool2 = nn.MaxPool2d(kernel_size=4, stride=2, padding
78                 =1)
79             self.activ2 = nn.ReLU()
80
81             self.bn3 = nn.BatchNorm2d(384, affine=False)

```

```

71     self.conv3 = nn.Conv2d(in_channels=384, out_channels=384,
72         kernel_size=(3,3), padding=1, bias=False, groups=384)
73     self.bn_point3 = nn.BatchNorm2d(384, affine=False)
74     self.conv_point3 = nn.Conv2d(in_channels=384, out_channels
75         =1536, kernel_size=(1,1), padding=0, bias=False)
76     self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
77         =0)
78     self.activ3 = nn.ReLU()
79
80     self.flatten = nn.Flatten()
81     # Final fully-connected layer classifier
82     self.fc1 = nn.Linear(24576, 10)
83     self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
84         24576)
85     self.dropout = nn.Dropout(0.5)
86
87     # Apply custom initialization to depthwise convolutional
88     # layers
89     custom_pointwise_conv_init(self.conv1)
90     custom_depthwise_conv_init(self.conv2)
91     custom_depthwise_conv_init(self.conv3)
92
93     # Apply custom initialization to pointwise convolutional
94     # layers
95     custom_pointwise_conv_init(self.conv_point2)
96     custom_pointwise_conv_init(self.conv_point3)
97
98     def forward_features(self, x):
99         x = self.pool1(self.activ1(self.conv1(self.bn1(x))))
100        return x
101
102    def features_extract(self, x):
103        x = self.forward_features(x)
104        x = self.pool2(self.activ2(self.conv_point2(self.bn_point2(
105            self.conv2(self.bn2(x))))))
106        x = self.pool3(self.activ3(self.conv_point3(self.bn_point3(
107            self.conv3(self.bn3(x))))))
108        return x
109
110    def forward(self, x):
111        x = self.features_extract(x)
112        x = self.flatten(x)
113        x = self.fc1(self.dropout(x))
114        return x
115
116    def plot_grid(self, tensor, path, num_rows=5, num_cols=5,
117        layer_name=""):
118        # Ensure we're working with the first 25 filters (or less if
119        # there are fewer)
120        tensor = tensor[:25]
121        # Normalize the tensor
122        tensor = (tensor - tensor.min()) / (tensor.max() - tensor.
123            min() + 1e-8)
124        # Move to CPU and convert to numpy
125        tensor = tensor.cpu().detach().numpy()
126
127        if tensor.shape[2] == 1 and tensor.shape[3] == 1: # 1x1
128            convolution case

```

```

117     out_channels, in_channels = tensor.shape[:2]
118     fig = plt.figure(figsize=(14, 10))
119
120     # Create a gridspec for the layout
121     gs = fig.add_gridspec(2, 2, width_ratios=[20, 1],
122                           height_ratios=[1, 3],
123                           left=0.1, right=0.9, bottom=0.1,
124                           top=0.9, wspace=0.05, hspace
125                           =0.2)
126
127     ax1 = fig.add_subplot(gs[0, 0])
128     ax2 = fig.add_subplot(gs[1, 0])
129     cbar_ax = fig.add_subplot(gs[:, 1])
130
131     # Bar plot for average weights per filter
132     avg_weights = tensor.mean(axis=(1, 2, 3))
133     norm = plt.Normalize(vmin=avg_weights.min(), vmax=
134                           avg_weights.max())
135     im1 = ax1.bar(range(out_channels), avg_weights, color=
136                   plt.cm.RdYlGn(norm(avg_weights)))
137     ax1.set_xlabel('Filter Index')
138     ax1.set_ylabel('Average Weight')
139     ax1.set_title(f'Average Weights for 1x1 Kernels in {layer_name}')
140     ax1.axhline(y=0, color='black', linestyle='--',
141                  linewidth=0.5)
142
143     # Heatmap for detailed weight distribution
144     im2 = ax2.imshow(tensor.reshape(out_channels,
145                           in_channels), cmap='RdYlGn', aspect='auto', norm=norm
146                           )
147     ax2.set_xlabel('Input Channel')
148     ax2.set_ylabel('Output Channel (Filter)')
149     ax2.set_title('Detailed Weight Distribution')
150
151     # Add colorbar to the right of both subplots
152     fig.colorbar(im2, cax=cbar_ax, label='Normalized Weight
153                 Value')
154
155 else:
156     fig, axes = plt.subplots(num_rows, num_cols, figsize
157                             =(15, 10))
158     fig.suptitle(f'First 25 Filters of {layer_name}')
159     for i, ax in enumerate(axes.flat):
160         if i < tensor.shape[0]:
161             filter_img = tensor[i]
162             # Handle different filter shapes
163             if filter_img.shape[0] == 3: # RGB filter (3, H
164               , W)
165                 filter_img = np.transpose(filter_img, (1, 2,
166                                           0))
167             elif filter_img.shape[0] == 1: # Grayscale
168                 filter_img = filter_img.squeeze()
169             else: # Multi-channel filter (C, H, W), take
170                   mean across channels
171                 filter_img = np.mean(filter_img, axis=0)

```

```

159             ax.imshow(filter_img, cmap='viridis' if
160                     filter_img.ndim == 2 else None)
161             ax.set_title(f'Filter {i + 1}')
162             ax.axis('off')
163
164     if path:
165         fig.savefig(path, bbox_inches='tight')
166     wandb.log({f'{layer_name} filters': wandb.Image(fig)})
167     plt.close(fig)
168
169     def visualize_filters(self, layer_name='conv1', save_path=None):
170         weights = getattr(self, layer_name).weight.data
171         self.plot_grid(weights, save_path, layer_name=layer_name)
172
173 class Net_Backpropagation(nn.Module):
174     def __init__(self):
175         super(Net_Backpropagation, self).__init__()
176
177         # A single Depthwise convolutional layer
178         self.bn1 = nn.BatchNorm2d(3, affine=False)
179         self.conv1 = nn.Conv2d(in_channels=3, out_channels=96,
180                             kernel_size=(5,5), padding=2, bias=False)
181         self.pool1 = nn.MaxPool2d(kernel_size=4, stride=2, padding
182                                 =1)
183         self.activ1 = nn.ReLU()
184
185         self.bn2 = nn.BatchNorm2d(96, affine=False)
186         self.conv2 = nn.Conv2d(in_channels=96, out_channels=384,
187                             kernel_size=(3,3), padding=1, bias=False)
188         self.pool2 = nn.MaxPool2d(kernel_size=4, stride=2, padding
189                                 =1)
190         self.activ2 = nn.ReLU()
191
192         self.bn3 = nn.BatchNorm2d(384, affine=False)
193         self.conv3 = nn.Conv2d(in_channels=384, out_channels=1536,
194                             kernel_size=(3,3), padding=1, bias=False)
195         self.pool3 = nn.AvgPool2d(kernel_size=2, stride=2, padding
196                                 =0)
197         self.activ3 = nn.ReLU()
198
199         self.flatten = nn.Flatten()
200         # Final fully-connected layer classifier
201         self.fc1 = nn.Linear(24576, 10)
202         self.fc1.weight.data = 0.11048543456039805 * torch.rand(10,
203                                                               24576)
204         self.dropout = nn.Dropout(0.5)
205
206         # Apply custom initialization to depthwise convolutional
207         # layers
208         custom_pointwise_conv_init(self.conv1)
209         custom_pointwise_conv_init(self.conv2)
210         custom_pointwise_conv_init(self.conv3)
211
212     def forward_features(self, x):
213         x = self.pool1(self.activ1(self.conv1(self.bn1(x))))
214         return x

```

```

208     def features_extract(self, x):
209         x = self.forward_features(x)
210         x = self.pool2(self.activ2(self.conv2(self.bn2(x))))
211         x = self.pool3(self.activ3(self.conv3(self.bn3(x))))
212         return x
213
214     def forward(self, x):
215         x = self.features_extract(x)
216         x = self.flatten(x)
217         x = self.fc1(self.dropout(x))
218         return x
219
220     def plot_grid(self, tensor, path, num_rows=5, num_cols=5,
221                  layer_name=""):
222         # Ensure we're working with the first 25 filters (or less if
223         # there are fewer)
224         tensor = tensor[:25]
225         # Normalize the tensor
226         tensor = (tensor - tensor.min()) / (tensor.max() - tensor.
227             min() + 1e-8)
228         # Move to CPU and convert to numpy
229         tensor = tensor.cpu().detach().numpy()
230
231         if tensor.shape[2] == 1 and tensor.shape[3] == 1: # 1x1
232             convolution case
233             out_channels, in_channels = tensor.shape[:2]
234             fig = plt.figure(figsize=(14, 10))
235
236             # Create a gridspec for the layout
237             gs = fig.add_gridspec(2, 2, width_ratios=[20, 1],
238                                   height_ratios=[1, 3],
239                                   left=0.1, right=0.9, bottom=0.1,
240                                   top=0.9, wspace=0.05, hspace
241                                   =0.2)
242
243             ax1 = fig.add_subplot(gs[0, 0])
244             ax2 = fig.add_subplot(gs[1, 0])
245             cbar_ax = fig.add_subplot(gs[:, 1])
246
247             # Bar plot for average weights per filter
248             avg_weights = tensor.mean(axis=(1, 2, 3))
249             norm = plt.Normalize(vmin=avg_weights.min(), vmax=
250                 avg_weights.max())
251             im1 = ax1.bar(range(out_channels), avg_weights, color=
252                 plt.cm.RdYlGn(norm(avg_weights)))
253             ax1.set_xlabel('Filter Index')
254             ax1.set_ylabel('Average Weight')
255             ax1.set_title(f'Average Weights for 1x1 Kernels in {
256                 layer_name}')
257             ax1.axhline(y=0, color='black', linestyle='--',
258                         linewidth=0.5)
259
260             # Heatmap for detailed weight distribution
261             im2 = ax2.imshow(tensor.reshape(out_channels,
262                 in_channels), cmap='RdYlGn', aspect='auto', norm=norm
263             )
264             ax2.set_xlabel('Input Channel')
265             ax2.set_ylabel('Output Channel (Filter)')

```

```

253     ax2.set_title('Detailed Weight Distribution')
254
255     # Add colorbar to the right of both subplots
256     fig.colorbar(im2, cax=cbar_ax, label='Normalized Weight
257                 Value')
258
259     else:
260         fig, axes = plt.subplots(num_rows, num_cols, figsize
261                                 =(15, 10))
262         fig.suptitle(f'First 25 Filters of {layer_name}')
263         for i, ax in enumerate(axes.flat):
264             if i < tensor.shape[0]:
265                 filter_img = tensor[i]
266                 # Handle different filter shapes
267                 if filter_img.shape[0] == 3: # RGB filter (3, H
268                     , W)
269                     filter_img = np.transpose(filter_img, (1, 2,
270                                         0))
271                 elif filter_img.shape[0] == 1: # Grayscale
272                     filter (1, H, W)
273                     filter_img = filter_img.squeeze()
274                 else: # Multi-channel filter (C, H, W), take
275                     mean across channels
276                     filter_img = np.mean(filter_img, axis=0)
277                     ax.imshow(filter_img, cmap='viridis' if
278                               filter_img.ndim == 2 else None)
279                     ax.set_title(f'Filter {i + 1}')
280                     ax.axis('off')
281
282         if path:
283             fig.savefig(path, bbox_inches='tight')
284             wandb.log({f'{layer_name} filters': wandb.Image(fig)})
285             plt.close(fig)
286
287     def visualize_filters(self, layer_name='conv1', save_path=None):
288         weights = getattr(self, layer_name).weight.data
289         self.plot_grid(weights, save_path, layer_name=layer_name)

```

A.2.10 Experiment_hebbian.py

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torch.optim.lr_scheduler as sched
5
6 import data
7 from model_hebb import Net_Hebbian
8 import matplotlib.pyplot as plt
9 import warnings
10
11 from logger import Logger
12 from torchmetrics import Accuracy, Precision, Recall, F1Score,
13     ConfusionMatrix
14 import seaborn as sns
15 import wandb
16 from visualizer import plot_ltp_ltd, print_weight_statistics,
17     visualize_data_clusters
18 from receptive_fields import visualize_filters

```

```

17 """
18 """
19 Code structure inspired and modified from https://github.com/
    NeuromorphicComputing/SoftHebb
20 Similar training and testing loop, with some identical function for
    the custom learning rate schedule
21 """
22
23 torch.manual_seed(0)
24
25 # Calculate evaluation metrics
26 def calculate_metrics(preds, labels, num_classes):
27     if num_classes == 2:
28         accuracy = Accuracy(task='binary', num_classes=num_classes) .
29             to(device)
30         precision = Precision(task='binary', average='weighted',
31             num_classes=num_classes).to(device)
32         recall = Recall(task='binary', average='weighted',
33             num_classes=num_classes).to(device)
34         f1 = F1Score(task='binary', average='weighted', num_classes=
35             num_classes).to(device)
36         confusion_matrix = ConfusionMatrix(task='binary',
37             num_classes=num_classes).to(device)
38     else:
39         accuracy = Accuracy(task='multiclass', num_classes=
40             num_classes).to(device)
41         precision = Precision(task='multiclass', average='macro',
42             num_classes=num_classes).to(device)
43         recall = Recall(task='multiclass', average='macro',
44             num_classes=num_classes).to(device)
45         f1 = F1Score(task='multiclass', average='macro', num_classes
46             =num_classes).to(device)
47         confusion_matrix = ConfusionMatrix(task='multiclass',
48             num_classes=num_classes).to(device)
49
50     acc = accuracy(preds, labels)
51     prec = precision(preds, labels)
52     rec = recall(preds, labels)
53     f1_score = f1(preds, labels)
54     conf_matrix = confusion_matrix(preds, labels)
55
56     return acc, prec, rec, f1_score, conf_matrix
57
58 class WeightNormDependentLR(optim.lr_scheduler._LRScheduler):
59 """
60     Custom Learning Rate Scheduler for unsupervised training of
61         SoftHebb Convolutional blocks.
62     Difference between current neuron norm and theoretical converged
63         norm (=1) scales the initial lr.
64 """
65
66     def __init__(self, optimizer, power_lr, last_epoch=-1, verbose=
67         False):
68         self.optimizer = optimizer
69         self.initial_lr_groups = [group['lr'] for group in self.
70             optimizer.param_groups] # store initial lrs
71         self.power_lr = power_lr
72         super().__init__(optimizer, last_epoch, verbose)

```

```

59
60     def get_lr(self):
61         if not self._get_lr_called_within_step:
62             warnings.warn("To get the last learning rate computed by
63                           the scheduler, "
64                           "please use 'get_last_lr()' .", UserWarning
65                           )
66
67         new_lr = []
68         for i, group in enumerate(self.optimizer.param_groups):
69             for param in group['params']:
70                 # difference between current neuron norm and
71                 # theoretical converged norm (=1) scales the
72                 # initial lr
73                 # initial_lr * |neuron_norm - 1| ** 0.5
74                 norm_diff = torch.abs(torch.linalg.norm(param.view(
75                     param.shape[0], -1), dim=1, ord=2) - 1) + 1e-10
76                 new_lr.append(self.initial_lr_groups[i] * (norm_diff
77                         ** self.power_lr)[:, None, None, None])
78
79         return new_lr
80
81
82     class TensorLRSGD(optim.SGD):
83         @torch.no_grad()
84         def step(self, closure=None):
85             """Performs a single optimization step, using a non-scalar (
86                 tensor) learning rate.
87             Arguments:
88                 closure (callable, optional): A closure that reevaluates
89                     the model
90                     and returns the loss.
91             """
92             loss = None
93             if closure is not None:
94                 with torch.enable_grad():
95                     loss = closure()
96             for group in self.param_groups:
97                 weight_decay = group['weight_decay']
98                 momentum = group['momentum']
99                 dampening = group['dampening']
100                nesterov = group['nesterov']
101                for p in group['params']:
102                    if p.grad is None:
103                        continue
104                    d_p = p.grad
105                    if weight_decay != 0:
106                        d_p = d_p.add(p, alpha=weight_decay)
107                    if momentum != 0:
108                        param_state = self.state[p]
109                        if 'momentum_buffer' not in param_state:
110                            buf = param_state['momentum_buffer'] = torch
111                                .clone(d_p).detach()
112                        else:
113                            buf = param_state['momentum_buffer']
114                            buf.mul_(momentum).add_(d_p, alpha=1 -
115                                dampening)
116                        if nesterov:
117                            d_p = d_p.add(buf, alpha=momentum)
118                        else:
119                            d_p = buf

```

```

107             p.add_(-group['lr'] * d_p)
108     return loss
109
110 # Main training and testing loop
111 # Experiment is configured to replicate SoftHebb-Surr/HardWTA/Cos-
112 # Instar configuration results
113 if __name__ == "__main__":
114     hebb_param = {'mode': 'hard', 'w_nrm': False, 'act': nn.Identity
115                 (), 'k': 1, 'alpha': 1.}
116     device = torch.device('cuda:0')
117     model = Net_Hebbian(hebb_params=hebb_param, version="softthebb")
118     model.to(device)
119
120     wandb_logger = Logger(
121         f"SoftHebb-Surr/HardWTA/Cos-Instar", project='Final-ReRuns -'
122         'HebbianCNN', model=model)
123     logger = wandb_logger.get_logger()
124     num_parameters = sum(p.numel() for p in model.parameters() if p.
125                           requires_grad)
126     print(f"Parameter Count Total: {num_parameters}")
127
128     # Custom learning rate and scheduler for only SoftHebb
129     # implementation
130     # unsup_optimizer = TensorLRSGD([
131     #     {"params": model.conv1.parameters(), "lr": 0.08, },
132     #     {"params": model.conv2.parameters(), "lr": 0.005, },
133     #     {"params": model.conv3.parameters(), "lr": 0.01, },
134     # ], lr=0)
135     # unsup_lr_scheduler = WeightNormDependentLR(unsup_optimizer,
136     #                                              power_lr=0.5)
137
138     # Learning rate for all other hebbian experiments
139     hebb_params = [
140         {'params': model.conv1.parameters(), 'lr': 0.1},
141         {'params': model.conv2.parameters(), 'lr': 0.1},
142         {'params': model.conv3.parameters(), 'lr': 0.1}
143     ]
144     unsup_optimizer = optim.SGD(hebb_params, lr=0) # The lr here
145     # will be overridden by the individual lrs
146
147     sup_optimizer = optim.Adam(model.fc1.parameters(), lr=0.001)
148     criterion = nn.CrossEntropyLoss()
149
150     trn_set, tst_set, zca = data.get_data(dataset='cifar10', root='
151         datasets', batch_size=64,
152                                         whiten_lvl=1e-3)
153     print(f'Processing Training batches: {len(trn_set)})')
154
155     print("Initial Weight statistics")
156     print_weight_statistics(model.conv1, 'conv1')
157     print_weight_statistics(model.conv2, 'conv2')
158     print_weight_statistics(model.conv3, 'conv3')
159
160     print("Visualizing Initial Filters")
161     model.visualize_filters('conv1')
162
163     running_loss = 0.0

```

```

157 # Hebbian feature extraction
158 for epoch in range(1):
159     print(f"Training Hebbian epoch {epoch}")
160     for i, data in enumerate(trn_set, 0):
161         inputs, _ = data
162         inputs = inputs.to(device)
163         # zero the parameter gradients
164         unsup_optimizer.zero_grad()
165         with torch.no_grad():
166             outputs = model(inputs)
167             # Visualize changes before updating
168             if i % 200 == 0: # Every 100 datapoint
169                 print(f'Saving details after batch {i}')
170                 plot_ltp_ltd(model.conv1, 'conv1', num_filters=10,
171                               detailed_mode=True)
172                 plot_ltp_ltd(model.conv2, 'conv2', num_filters=10,
173                               detailed_mode=True)
174                 plot_ltp_ltd(model.conv3, 'conv3', num_filters=10,
175                               detailed_mode=True)
176
177                 model.visualize_filters('conv1')
178                 model.visualize_filters('conv2')
179                 model.visualize_filters('conv3')
180             for layer in [model.conv1, model.conv2, model.conv3]:
181                 if hasattr(layer, 'local_update'):
182                     layer.local_update()
183             unsup_optimizer.step()
184             # Scheduler only for SoftHebb
185             # unsup_lr_scheduler.step()
186             print("Visualizing Filters")
187             model.visualize_filters('conv1', f'results>{"demo"}/'
188                                     demo_conv1_filters_epoch_{1}.png')
189             model.visualize_filters('conv2', f'results>{"demo"}/'
190                                     demo_conv2_filters_epoch_{1}.png')
191             model.visualize_filters('conv3', f'results>{"demo"}/'
192                                     demo_conv3_filters_epoch_{1}.png')
193
194             # set requires_grad false and eval mode for all modules but
195             # classifier
196             unsup_optimizer.zero_grad()
197             model.conv1.requires_grad = False
198             model.conv2.requires_grad = False
199             model.conv3.requires_grad = False
200             # model.conv4.requires_grad = False
201             model.conv1.eval()
202             model.conv2.eval()
203             model.conv3.eval()
204             # model.conv4.eval()
205             model.bn1.eval()
206             model.bn2.eval()
207             model.bn3.eval()
208             # model.bn4.eval()
209             print("Visualizing Test Class separation")
210             visualize_data_clusters(tst_set, model=model, method='umap', dim
211                                     =2)
212             # Supervised training of classifier
213             # Train classifier with backpropagation
214             print("Training Classifier")

```

```

207     for epoch in range(10):
208         model.fc1.train()
209         model.dropout.train()
210         running_loss = 0.0
211         correct = 0
212         total = 0
213         train_preds = []
214         train_labels = []
215         for i, data in enumerate(trn_set, 0):
216             inputs, labels = data
217             inputs = inputs.to(device)
218             labels = labels.to(device)
219             # zero the parameter gradients
220             sup_optimizer.zero_grad()
221             outputs = model(inputs)
222             loss = criterion(outputs, labels)
223             loss.backward()
224             sup_optimizer.step()
225             # compute training statistics
226             running_loss += loss.item()
227             total += labels.size(0)
228             _, predicted = torch.max(outputs.data, 1)
229             correct += (predicted == labels).sum().item()
230             # For wandb logs
231             preds = torch.argmax(outputs, dim=1)
232             train_preds.append(preds)
233             train_labels.append(labels)
234
235             print(f'Accuracy of the network on the train images: {100 * correct / total} %')
236             print(f'{epoch + 1}] loss: {running_loss / total:.3f}')
237
238             train_preds = torch.cat(train_preds, dim=0)
239             train_labels = torch.cat(train_labels, dim=0)
240             acc, prec, rec, f1_score, conf_matrix = calculate_metrics(
241                 train_preds, train_labels, 10)
242             logger.log({'train_accuracy': acc, 'train_precision': prec,
243                         'train_recall': rec, 'train_f1_score': f1_score})
244             f, ax = plt.subplots(figsize=(15, 10))
245             sns.heatmap(conf_matrix.clone().detach().cpu().numpy(),
246                         annot=True, ax=ax)
247             logger.log({"train_confusion_matrix": wandb.Image(f)})
248             plt.close(f)
249
250             # Evaluation on test set
251             model.eval()
252             running_loss = 0.
253             correct = 0
254             total = 0
255             test_preds = []
256             test_labels = []
257             # since we're not training, we don't need to calculate the
258             # gradients for our outputs/model
259             with torch.no_grad():
260                 for data in tst_set:
261                     images, labels = data
262                     images = images.to(device)
263                     labels = labels.to(device)

```

```

260         # calculate outputs by running images through the
261         # network
262         outputs = model(images)
263         # the class with the highest energy is what we
264         # choose as prediction
265         _, predicted = torch.max(outputs.data, 1)
266         total += labels.size(0)
267         correct += (predicted == labels).sum().item()
268         loss = criterion(outputs, labels)
269         running_loss += loss.item()
270         # For wandb logs
271         preds = torch.argmax(outputs, dim=1)
272         test_preds.append(preds)
273         test_labels.append(labels)
274
275         print(f'Accuracy of the network on the test images: {100 * correct / total}%')
276         print(f'test loss: {running_loss / total:.3f}')
277
278         test_preds = torch.cat(test_preds, dim=0)
279         test_labels = torch.cat(test_labels, dim=0)
280         acc, prec, rec, f1_score, conf_matrix = calculate_metrics(
281             test_preds, test_labels, 10)
282         logger.log({'test_accuracy': acc, 'test_precision': prec,
283                     'test_recall': rec, 'test_f1_score': f1_score})
284         f, ax = plt.subplots(figsize=(15, 10))
285         sns.heatmap(conf_matrix.clone().detach().cpu().numpy(),
286                     annot=True, ax=ax)
287         logger.log({"test_confusion_matrix": wandb.Image(f)})
288         plt.close(f)
289
290     # Visualise receptive fields at end, as this code alters the
291     # model architecture (padding removal)
292     print("Visualizing Receptive fields")
293     visualize_filters(model, model.conv1, num_filters=25)
294     visualize_filters(model, model.conv2, num_filters=25)
295     visualize_filters(model, model.conv3, num_filters=25)

```

A.2.11 Experiment_hebbian_depthwise.py

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torch.optim.lr_scheduler as sched
5
6 import data
7 from model_depthwise import Net_Depthwise
8 import numpy as np
9 import matplotlib.pyplot as plt
10 import warnings
11
12 from logger import Logger
13 from torchmetrics import Accuracy, Precision, Recall, F1Score,
14     ConfusionMatrix
15 import seaborn as sns
16 import wandb
17 from visualizer import plot_ltp_ltd, plot_ltp_ltd_ex_in,
18     print_weight_statistics, visualize_data_clusters

```

```

17 import pandas as pd
18 from receptive_fields import visualize_filters
19
20 """
21 Identical to experiment_hebbian.py, but adapted for depthwise
22   convolutions
23 """
24 torch.manual_seed(0)
25
26 def calculate_metrics(preds, labels, num_classes):
27     if num_classes == 2:
28         accuracy = Accuracy(task='binary', num_classes=num_classes).to(device)
29         precision = Precision(task='binary', average='weighted',
30                               num_classes=num_classes).to(device)
31         recall = Recall(task='binary', average='weighted',
32                          num_classes=num_classes).to(device)
33         f1 = F1Score(task='binary', average='weighted', num_classes=
34                      num_classes).to(device)
35         confusion_matrix = ConfusionMatrix(task='binary',
36                                             num_classes=num_classes).to(device)
37     else:
38         accuracy = Accuracy(task='multiclass', num_classes=
39                               num_classes).to(device)
40         precision = Precision(task='multiclass', average='macro',
41                               num_classes=num_classes).to(device)
42         recall = Recall(task='multiclass', average='macro',
43                          num_classes=num_classes).to(device)
44         f1 = F1Score(task='multiclass', average='macro', num_classes
45                      =num_classes).to(device)
46         confusion_matrix = ConfusionMatrix(task='multiclass',
47                                             num_classes=num_classes).to(device)
48
49     acc = accuracy(preds, labels)
50     prec = precision(preds, labels)
51     rec = recall(preds, labels)
52     f1_score = f1(preds, labels)
53     conf_matrix = confusion_matrix(preds, labels)
54
55     return acc, prec, rec, f1_score, conf_matrix
56
57 class WeightNormDependentLR(optim.lr_scheduler._LRScheduler):
58     """
59     Custom Learning Rate Scheduler for unsupervised training of
60     SoftHebb Convolutional blocks.
61     Difference between current neuron norm and theoretical converged
62     norm (=1) scales the initial lr.
63     """
64
65     def __init__(self, optimizer, power_lr, last_epoch=-1, verbose=False):
66         self.optimizer = optimizer
67         self.initial_lr_groups = [group['lr'] for group in self.
68                                  optimizer.param_groups] # store initial lrs
69         self.power_lr = power_lr
70         super().__init__(optimizer, last_epoch, verbose)
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159

```

```

60     def get_lr(self):
61         if not self._get_lr_called_within_step:
62             warnings.warn("To get the last learning rate computed by
63                             the scheduler, "
64                             "please use 'get_last_lr()' .", UserWarning
65                             )
66         new_lr = []
67         for i, group in enumerate(self.optimizer.param_groups):
68             for param in group['params']:
69                 # difference between current neuron norm and
70                 # theoretical converged norm (=1) scales the
71                 # initial lr
72                 # initial_lr * |neuron_norm - 1| ** 0.5
73                 norm_diff = torch.abs(torch.linalg.norm(param.view(
74                     param.shape[0], -1), dim=1, ord=2) - 1) + 1e-10
75                 new_lr.append(self.initial_lr_groups[i] * (norm_diff
76                     ** self.power_lr)[:, None, None, None])
77         return new_lr
78
79
80
81     class TensorLRSGD(optim.SGD):
82         @torch.no_grad()
83         def step(self, closure=None):
84             """Performs a single optimization step, using a non-scalar (
85                 tensor) learning rate.
86             Arguments:
87                 closure (callable, optional): A closure that reevaluates
88                     the model
89                     and returns the loss.
90             """
91             loss = None
92             if closure is not None:
93                 with torch.enable_grad():
94                     loss = closure()
95             for group in self.param_groups:
96                 weight_decay = group['weight_decay']
97                 momentum = group['momentum']
98                 dampening = group['dampening']
99                 nesterov = group['nesterov']
100                for p in group['params']:
101                    if p.grad is None:
102                        continue
103                    d_p = p.grad
104                    if weight_decay != 0:
105                        d_p = d_p.add(p, alpha=weight_decay)
106                    if momentum != 0:
107                        param_state = self.state[p]
108                        if 'momentum_buffer' not in param_state:
109                            buf = param_state['momentum_buffer'] = torch
110                                .clone(d_p).detach()
111                        else:
112                            buf = param_state['momentum_buffer']
113                            buf.mul_(momentum).add_(d_p, alpha=1 -
114                                dampening)
115                        if nesterov:
116                            d_p = d_p.add(buf, alpha=momentum)
117                        else:
118                            d_p = buf

```

```

108         p.add_(-group['lr'] * d_p)
109     return loss
110
111 if __name__ == "__main__":
112
113     hebb_param = {'mode': 'hard', 'w_nrm': False, 'act': nn.Identity
114                 (), 'k': 1, 'alpha': 1.}
115     device = torch.device('cuda:0')
116     model = Net_Depthwise(hebb_params=hebb_param, version="softthebb"
117                           )
118     model.to(device)
119
120     wandb_logger = Logger(
121         f"Depthwise_HardHebb-Surr/HardWTA/Cos-Instar", project='Final
122             -ReRuns-HebbianCNN', model=model)
123     logger = wandb_logger.get_logger()
124     num_parameters = sum(p.numel() for p in model.parameters() if p.
125                           requires_grad)
126     print(f"Parameter Count Total: {num_parameters}")
127
128     # unsup_optimizer = TensorLRSGD([
129     #     {"params": model.conv1.parameters(), "lr": 0.08, },
130     #     {"params": model.conv2.parameters(), "lr": 0.005, },
131     #     {"params": model.conv_point2.parameters(), "lr": 0.005, },
132     #     {"params": model.conv3.parameters(), "lr": 0.01, },
133     #     {"params": model.conv_point3.parameters(), "lr": 0.01, }
134     # ], lr=0)
135     # unsup_lr_scheduler = WeightNormDependentLR(unsup_optimizer,
136     #                                              power_lr=0.5)
137
138     hebb_params = [
139         {'params': model.conv1.parameters(), 'lr': 0.1},
140         {'params': model.conv2.parameters(), 'lr': 0.1},
141         {'params': model.conv_point2.parameters(), 'lr': 0.1},
142         {'params': model.conv3.parameters(), 'lr': 0.1},
143         {'params': model.conv_point3.parameters(), 'lr': 0.1}
144     ]
145     unsup_optimizer = optim.SGD(hebb_params, lr=0) # The lr here
146     will be overridden by the individual lrs
147
148     sup_optimizer = optim.Adam(model.fc1.parameters(), lr=0.001)
149     criterion = nn.CrossEntropyLoss()
150
151     trn_set, tst_set, zca = data.get_data(dataset='cifar10', root='
152         datasets', batch_size=64,
153                                         whiten_lvl=1e-3)
154     print(f'Processing Training batches: {len(trn_set)}')
155
156     print("Initial Weight statistics")
157     print_weight_statistics(model.conv1, 'conv1')
158     print_weight_statistics(model.conv2, 'conv2')
159     print_weight_statistics(model.conv3, 'conv3')
160     print_weight_statistics(model.conv_point2, 'conv_point2')
161
162     running_loss = 0.0
163     for epoch in range(1):
164         print(f"Training Hebbian epoch {epoch}")
165         for i, data in enumerate(trn_set, 0):

```

```

159     inputs, _ = data
160     inputs = inputs.to(device)
161     # zero the parameter gradients
162     # unsup_optimizer.zero_grad()
163     with torch.no_grad():
164         outputs = model(inputs)
165     # Visualize changes before updating
166     if i % 200 == 0: # Every 100 datapoint
167         print(f'Saving details after batch {i}')
168         plot_ltp_ltd(model.conv1, 'conv1', num_filters=10,
169                       detailed_mode=True)
170         plot_ltp_ltd(model.conv2, 'conv2', num_filters=10,
171                       detailed_mode=True)
172         # plot_ltp_ltd(model.conv_point2, 'conv_point2',
173         #               num_filters=10, detailed_mode=True)
174         model.visualize_filters('conv1')
175         model.visualize_filters('conv2')
176         model.visualize_filters('conv_point2')
177
178         model.visualize_filters('conv3')
179         model.visualize_filters('conv_point3')
180     for layer in [model.conv1, model.conv2, model.conv3,
181                   model.conv_point2, model.conv_point3]:
182     # for layer in [model.conv1, model.conv2, model.conv3]:
183         if hasattr(layer, 'local_update'):
184             layer.local_update()
185         # unsup_optimizer.step()
186         # unsup_lr_scheduler.step()
187     print("Visualizing Filters")
188     model.visualize_filters('conv1', f'results>{"demo"}/'
189                             demo_conv1_filters_epoch_{1}.png')
190     model.visualize_filters('conv2', f'results>{"demo"}/'
191                             demo_conv2_filters_epoch_{1}.png')
192     model.visualize_filters('conv3', f'results>{"demo"}/'
193                             demo_conv3_filters_epoch_{1}.png')
194     # model.visualize_filters('conv_point2', f'results>{"demo"}/'
195     #   demo_conv_point2_filters_epoch_{1}.png')
196
197     # Supervised training of classifier
198     # set requires_grad false and eval mode for all modules but
199     # classifier
200     unsup_optimizer.zero_grad()
201     model.conv1.requires_grad = False
202     model.conv2.requires_grad = False
203     model.conv3.requires_grad = False
204     # model.conv4.requires_grad = False
205     model.conv1.eval()
206     model.conv2.eval()
207     model.conv3.eval()
208     # model.conv4.eval()
209     model.bn1.eval()
210     model.bn2.eval()
211     model.bn3.eval()
212     # model.bn4.eval()
213     model.conv_point2.requires_grad = False
214     model.conv_point3.requires_grad = False
215     # model.conv_point4.requires_grad = False

```

```

208     model.conv_point2.eval()
209     model.conv_point3.eval()
210     # model.conv_point4.eval()
211     model.bn_point2.eval()
212     model.bn_point3.eval()
213     # model.bn_point4.eval()
214     print("Visualizing Test Class separation")
215     visualize_data_clusters(tst_set, model=model, method='umap', dim
216                               =2)
216     # Train classifier with backpropagation
217     print("Training Classifier")
218     for epoch in range(5):
219         model.fc1.train()
220         model.dropout.train()
221         running_loss = 0.0
222         correct = 0
223         total = 0
224         train_preds = []
225         train_labels = []
226         for i, data in enumerate(trn_set, 0):
227             inputs, labels = data
228             inputs = inputs.to(device)
229             labels = labels.to(device)
230             # zero the parameter gradients
231             sup_optimizer.zero_grad()
232             outputs = model(inputs)
233             loss = criterion(outputs, labels)
234             loss.backward()
235             sup_optimizer.step()
236             # compute training statistics
237             running_loss += loss.item()
238             total += labels.size(0)
239             _, predicted = torch.max(outputs.data, 1)
240             correct += (predicted == labels).sum().item()
241             # For wandb logs
242             preds = torch.argmax(outputs, dim=1)
243             train_preds.append(preds)
244             train_labels.append(labels)
245
246             print(f'Accuracy of the network on the train images: {100 *
247                   correct // total} %')
248             print(f'{[epoch + 1]} loss: {running_loss / total:.3f}')
249
250             train_preds = torch.cat(train_preds, dim=0)
251             train_labels = torch.cat(train_labels, dim=0)
252             acc, prec, rec, f1_score, conf_matrix = calculate_metrics(
253                 train_preds, train_labels, 10)
254             logger.log({'train_accuracy': acc, 'train_precision': prec,
255                         'train_recall': rec, 'train_f1_score': f1_score})
256             f, ax = plt.subplots(figsize=(15, 10))
257             sns.heatmap(conf_matrix.clone().detach().cpu().numpy(),
258                         annot=True, ax=ax)
259             logger.log({'train_confusion_matrix': wandb.Image(f)})
260             plt.close(f)
261
262             # Evaluation on test set
263             model.eval()
264             running_loss = 0.

```

```

261     correct = 0
262     total = 0
263     test_preds = []
264     test_labels = []
265     # since we're not training, we don't need to calculate the
266     # gradients for our outputs
267     with torch.no_grad():
268         for data in tst_set:
269             images, labels = data
270             images = images.to(device)
271             labels = labels.to(device)
272             # calculate outputs by running images through the
273             # network
274             outputs = model(images)
275             # the class with the highest energy is what we
276             # choose as prediction
277             _, predicted = torch.max(outputs.data, 1)
278             total += labels.size(0)
279             correct += (predicted == labels).sum().item()
280             loss = criterion(outputs, labels)
281             running_loss += loss.item()
282             # For wandb logs
283             preds = torch.argmax(outputs, dim=1)
284             test_preds.append(preds)
285             test_labels.append(labels)
286
287             print(f'Accuracy of the network on the test images: {100 *'
288                   ' correct / total} %')
289             print(f'test loss: {running_loss / total:.3f}')
290
291             test_preds = torch.cat(test_preds, dim=0)
292             test_labels = torch.cat(test_labels, dim=0)
293             acc, prec, rec, f1_score, conf_matrix = calculate_metrics(
294                 test_preds, test_labels, 10)
295             logger.log({'test_accuracy': acc, 'test_precision': prec, '
296                         'test_recall': rec, 'test_f1_score': f1_score})
297             f, ax = plt.subplots(figsize=(15, 10))
298             sns.heatmap(conf_matrix.clone().detach().cpu().numpy(),
299                         annot=True, ax=ax)
300             logger.log({"test_confusion_matrix": wandb.Image(f)})
301             plt.close(f)
302
303             print("Visualizing Receptive fields")
304             visualize_filters(model, model.conv1, num_filters=25)
305             visualize_filters(model, model.conv2, num_filters=25)
306             visualize_filters(model, model.conv3, num_filters=25)
307
308             visualize_filters(model, model.conv_point2, num_filters=25)
309             visualize_filters(model, model.conv_point3, num_filters=25)

```

A.2.12 Experiment_hebbian_residual.py

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torch.optim.lr_scheduler as sched
5
6 import data

```

```

7  from model_residual import Net_Depthwise_Residual
8  import numpy as np
9  import matplotlib.pyplot as plt
10 import warnings
11
12 from logger import Logger
13 from torchmetrics import Accuracy, Precision, Recall, F1Score,
14     ConfusionMatrix
15 import seaborn as sns
16 import wandb
17 from visualizer import plot_ltp_ltd, print_weight_statistics,
18     visualize_data_clusters
19 import pandas as pd
20 from receptive_fields_residual import visualize_filters
21
22 """
23 Identical to experiment_hebbian.py, but adapted for depthwise
24 convolutions
25 """
26
27 torch.manual_seed(0)
28
29 def calculate_metrics(preds, labels, num_classes):
30     if num_classes == 2:
31         accuracy = Accuracy(task='binary', num_classes=num_classes).
32             to(device)
33         precision = Precision(task='binary', average='weighted',
34             num_classes=num_classes).to(device)
35         recall = Recall(task='binary', average='weighted',
36             num_classes=num_classes).to(device)
37         f1 = F1Score(task='binary', average='weighted', num_classes=
38             num_classes).to(device)
39         confusion_matrix = ConfusionMatrix(task='binary',
40             num_classes=num_classes).to(device)
41     else:
42         accuracy = Accuracy(task='multiclass', num_classes=
43             num_classes).to(device)
44         precision = Precision(task='multiclass', average='macro',
45             num_classes=num_classes).to(device)
46         recall = Recall(task='multiclass', average='macro',
47             num_classes=num_classes).to(device)
48         f1 = F1Score(task='multiclass', average='macro', num_classes
49             =num_classes).to(device)
50         confusion_matrix = ConfusionMatrix(task='multiclass',
51             num_classes=num_classes).to(device)
52
53     acc = accuracy(preds, labels)
54     prec = precision(preds, labels)
55     rec = recall(preds, labels)
56     f1_score = f1(preds, labels)
57     conf_matrix = confusion_matrix(preds, labels)
58
59     return acc, prec, rec, f1_score, conf_matrix
60
61
62 class WeightNormDependentLR(optim.lr_scheduler._LRScheduler):
63     """

```

```

51     Custom Learning Rate Scheduler for unsupervised training of
52         SoftHebb Convolutional blocks.
53     Difference between current neuron norm and theoretical converged
54         norm (=1) scales the initial lr.
55     """
56
57
58     def __init__(self, optimizer, power_lr, last_epoch=-1, verbose=
59         False):
60         self.optimizer = optimizer
61         self.initial_lr_groups = [group['lr'] for group in self.
62             optimizer.param_groups] # store initial lrs
63         self.power_lr = power_lr
64         super().__init__(optimizer, last_epoch, verbose)
65
66     def get_lr(self):
67         if not self._get_lr_called_within_step:
68             warnings.warn("To get the last learning rate computed by
69                 the scheduler, "
70                         "please use 'get_last_lr()' .", UserWarning
71                         )
72         new_lr = []
73         for i, group in enumerate(self.optimizer.param_groups):
74             for param in group['params']:
75                 # difference between current neuron norm and
76                 # theoretical converged norm (=1) scales the
77                 # initial lr
78                 # initial_lr * |neuron_norm - 1| ** 0.5
79                 norm_diff = torch.abs(torch.linalg.norm(param.view(
80                     param.shape[0], -1), dim=1, ord=2) - 1) + 1e-10
81                 new_lr.append(self.initial_lr_groups[i] * (norm_diff
82                     ** self.power_lr)[:, None, None, None])
83         return new_lr
84
85
86     class TensorLRSGD(optim.SGD):
87         @torch.no_grad()
88         def step(self, closure=None):
89             """Performs a single optimization step, using a non-scalar (
90                 tensor) learning rate.
91             Arguments:
92                 closure (callable, optional): A closure that reevaluates
93                     the model
94                     and returns the loss.
95             """
96             loss = None
97             if closure is not None:
98                 with torch.enable_grad():
99                     loss = closure()
100             for group in self.param_groups:
101                 weight_decay = group['weight_decay']
102                 momentum = group['momentum']
103                 dampening = group['dampening']
104                 nesterov = group['nesterov']
105                 for p in group['params']:
106                     if p.grad is None:
107                         continue
108                     d_p = p.grad
109                     if weight_decay != 0:

```

```

97             d_p = d_p.add(p, alpha=weight_decay)
98         if momentum != 0:
99             param_state = self.state[p]
100            if 'momentum_buffer' not in param_state:
101                buf = param_state['momentum_buffer'] = torch
102                    .clone(d_p).detach()
103            else:
104                buf = param_state['momentum_buffer']
105                buf.mul_(momentum).add_(d_p, alpha=1 -
106                    dampening)
107            if nesterov:
108                d_p = d_p.add(buf, alpha=momentum)
109            else:
110                d_p = buf
111            p.add_(-group['lr'] * d_p)
112
113    return loss
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146

```

```

147 # Unsupervised training with SoftHebb
148 running_loss = 0.0
149 for epoch in range(1):
150     print(f"Training Hebbian epoch {epoch}")
151     for i, data in enumerate(trn_set, 0):
152         inputs, _ = data
153         inputs = inputs.to(device)
154         # zero the parameter gradients
155         # unsup_optimizer.zero_grad()
156         # forward + update computation
157         with torch.no_grad():
158             outputs = model(inputs)
159         # Visualize changes before updating
160         if i % 200 == 0: # Every 100 datapoint
161             print(f'Saving details after batch {i}')
162             plot_ltp_ltd(model.conv1, 'conv1', num_filters=10,
163                           detailed_mode=True)
163             plot_ltp_ltd(model.res1.conv1, 'res1.conv1',
164                           num_filters=10, detailed_mode=True)
164             plot_ltp_ltd(model.res1.conv2, 'res1.conv2',
165                           num_filters=10, detailed_mode=True)
165             plot_ltp_ltd(model.res1.conv3, 'res1.conv3',
166                           num_filters=10, detailed_mode=True)
166             plot_ltp_ltd(model.res2.conv1, 'res2.conv1',
167                           num_filters=10, detailed_mode=True)
167             plot_ltp_ltd(model.res2.conv2, 'res2.conv2',
168                           num_filters=10, detailed_mode=True)
168             plot_ltp_ltd(model.res2.conv3, 'res2.conv3',
169                           num_filters=10, detailed_mode=True)
169             model.visualize_filters('conv1')
170             model.visualize_filters('res1.conv2')
171             model.visualize_filters('res2.conv2')
172
173         layers = hebbian_layers = [model.conv1, model.res1, model.
174                                     res2]
174         for layer in layers:
175             if hasattr(layer, 'local_update'):
176                 layer.local_update()
177             else:
178                 for sublayer in layer.modules():
179                     if hasattr(sublayer, 'local_update'):
180                         sublayer.local_update()
181             # optimize
182             unsup_optimizer.step()
183             # unsup_lr_scheduler.step()
184
185         print("Visualizing Filters")
186         # model.visualize_filters('conv1', f'results>{"demo"}/
187         #     demo_conv1_filters_epoch_{1}.png')
187         # model.visualize_filters('res1.conv1', f'results>{"demo"}/
188         #     demo_res1_conv1_filters_epoch_{1}.png')
188         model.visualize_filters('res1.conv2', f'results>{"demo"}/
189         #     demo_res1_conv2_filters_epoch_{1}.png')
189         # model.visualize_filters('res1.conv3', f'results>{"demo"}/
190         #     demo_res1_conv3_filters_epoch_{1}.png')
190         # model.visualize_filters('res2.conv1', f'results>{"demo"}/
191         #     demo_res2_conv1_filters_epoch_{1}.png')

```

```

191     model.visualize_filters('res2.conv2', f'results>{"demo"}/
192         demo_res2_conv2_filters_epoch_{1}.png')
193 # model.visualize_filters('res2.conv3', f'results>{"demo"}/
194 #         demo_res2_conv3_filters_epoch_{1}.png')
195 # Supervised training of classifier
196 # set requires_grad false and eval mode for all modules but
197 # classifier
198 unsup_optimizer.zero_grad()
199 model.conv1.requires_grad = False
200 model.res1.requires_grad = False
201 model.res2.requires_grad = False
202 model.conv1.eval()
203 model.bn1.eval()
204 model.res1.eval()
205 model.res2.eval()

206 print("Visualizing Class separation")
207 visualize_data_clusters(tst_set, model=model, method='umap', dim
208 =2)
209 print("Training Classifier")
210 for epoch in range(5):
211     model.fc1.train()
212     model.dropout.train()
213     running_loss = 0.0
214     correct = 0
215     total = 0
216     train_preds = []
217     train_labels = []
218     for i, data in enumerate(trn_set, 0):
219         inputs, labels = data
220         inputs = inputs.to(device)
221         labels = labels.to(device)
222         # zero the parameter gradients
223         sup_optimizer.zero_grad()
224         # forward + backward + optimize
225         outputs = model(inputs)
226         loss = criterion(outputs, labels)
227         loss.backward()
228         sup_optimizer.step()
229         # compute training statistics
230         running_loss += loss.item()
231         total += labels.size(0)
232         _, predicted = torch.max(outputs.data, 1)
233         correct += (predicted == labels).sum().item()
234         # For wandb logs
235         preds = torch.argmax(outputs, dim=1)
236         train_preds.append(preds)
237         train_labels.append(labels)

238     print(f'Accuracy of the network on the train images: {100 *
239         correct // total} %')
240     print(f'{[epoch + 1]} loss: {running_loss / total:.3f}'')

241     train_preds = torch.cat(train_preds, dim=0)
242     train_labels = torch.cat(train_labels, dim=0)
243     acc, prec, rec, f1_score, conf_matrix = calculate_metrics(
244         train_preds, train_labels, 10)

```

```

242     logger.log({'train_accuracy': acc, 'train_precision': prec,
243                 'train_recall': rec, 'train_f1_score': f1_score})
244     f, ax = plt.subplots(figsize=(15, 10))
245     sns.heatmap(conf_matrix.clone().detach().cpu().numpy(),
246                  annot=True, ax=ax)
246     logger.log({"train_confusion_matrix": wandb.Image(f)})
247     plt.close(f)
248
249     # Evaluation on test set
250     model.eval()
251     running_loss = 0.
252     correct = 0
253     total = 0
254     # since we're not training, we don't need to calculate the
255     # gradients for our outputs
256     test_preds = []
257     test_labels = []
258     with torch.no_grad():
259         for data in tst_set:
260             images, labels = data
261             images = images.to(device)
262             labels = labels.to(device)
263             # calculate outputs by running images through the
264             # network
265             outputs = model(images)
266             # the class with the highest energy is what we
267             # choose as prediction
268             _, predicted = torch.max(outputs.data, 1)
269             total += labels.size(0)
270             correct += (predicted == labels).sum().item()
271             loss = criterion(outputs, labels)
272             running_loss += loss.item()
273             # For wandb logs
274             preds = torch.argmax(outputs, dim=1)
275             test_preds.append(preds)
276             test_labels.append(labels)
277
278             print(f'Accuracy of the network on the test images: {100 * correct / total}%')
279             print(f'test loss: {running_loss / total:.3f}')
280
281             test_preds = torch.cat(test_preds, dim=0)
282             test_labels = torch.cat(test_labels, dim=0)
283             acc, prec, rec, f1_score, conf_matrix = calculate_metrics(
284                 test_preds, test_labels, 10)
285             logger.log({'test_accuracy': acc, 'test_precision': prec,
286                         'test_recall': rec, 'test_f1_score': f1_score})
287             f, ax = plt.subplots(figsize=(15, 10))
288             sns.heatmap(conf_matrix.clone().detach().cpu().numpy(),
289                         annot=True, ax=ax)
290             logger.log({"test_confusion_matrix": wandb.Image(f)})
291             plt.close(f)
292
293             visualize_filters(model, model.conv1)
294             visualize_filters(model, model.res1.conv2)
295             visualize_filters(model, model.res2.conv2)

```

A.2.13 Experiment_bp.py

```
 1 from sklearn.manifold import TSNE
 2 from sklearn.preprocessing import StandardScaler
 3 import torch
 4 import torch.nn as nn
 5 import torch.optim as optim
 6 import torch.optim.lr_scheduler as sched
 7
 8 import data
 9 from Model_BackProp import Net_Backpropagation,
   Net_Backpropagation_depth
10 import numpy as np
11 import matplotlib.pyplot as plt
12 import umap
13 import warnings
14
15 from logger import Logger
16 from torchmetrics import Accuracy, Precision, Recall, F1Score,
   ConfusionMatrix
17 import seaborn as sns
18 import wandb
19 import pandas as pd
20
21 from receptive_fields import visualize_filters
22
23 torch.manual_seed(0)
24
25 def print_weight_statistics(layer, layer_name):
26     """
27         Prints statistics of the weights from a given layer.
28
29     Args:
30         layer (nn.Module): The layer to analyze
31         layer_name (str): Name of the layer for printing purposes
32     """
33     weights = layer.weight.data
34
35     print(f"Statistics for {layer_name}:")
36     print(f"Mean: {weights.mean().item():.4f}")
37     print(f"Max: {weights.max().item():.4f}")
38     print(f"Min: {weights.min().item():.4f}")
39     print(f"Standard Deviation: {weights.std().item():.4f}")
40     print(f"Median: {weights.median().item():.4f}")
41     print(f"25th Percentile: {weights.quantile(0.25).item():.4f}")
42     print(f"75th Percentile: {weights.quantile(0.75).item():.4f}")
43     print(f"Number of positive weights: {(weights > 0).sum().item()}")
44     print(f"Number of negative weights: {(weights < 0).sum().item()}")
45     print(f"Total number of weights: {weights.numel()}")
46     print()
47
48 def plot_ltp_ltd(layer, layer_name, num_filters=10, detailed_mode=
49     False):
50     weights = layer.weight.data
51     delta_w = layer.weight.grad
```

```

53     if not detailed_mode:
54         fig, ax = plt.subplots(figsize=(12, 6))
55
56         for i in range(min(num_filters, weights.shape[0])):
57             weight_change = delta_w[i].sum().item()
58             color = 'green' if weight_change > 0 else 'red'
59             ax.bar(i, weight_change, color=color)
60
61             ax.axhline(y=0, color='black', linestyle='--', linewidth=0.5)
62             ax.set_xlabel('Filter Index')
63             ax.set_ylabel('Total Weight Change')
64             ax.set_title(f'Overall LTP/LTD for first {num_filters}
65                           filters in {layer_name}')
66
67             # Log the plot to wandb
68             wandb.log({f"{layer_name}_Overall_LTP_LTD": wandb.Image(fig)
69                         })
70             plt.close(fig)
71
72     else:
73         # Plot 1: Overall weight change
74         fig, ax1 = plt.subplots(figsize=(12, 6))
75         for i in range(min(num_filters, weights.shape[0])):
76             weight_change = delta_w[i].sum().item()
77             color = 'green' if weight_change > 0 else 'red'
78             ax1.bar(i, weight_change, color=color)
79
80             ax1.axhline(y=0, color='black', linestyle='--', linewidth
81                         =0.5)
82             ax1.set_xlabel('Filter Index')
83             ax1.set_ylabel('Total Weight Change')
84             ax1.set_title(f'Overall LTP/LTD for first {num_filters}
85                           filters')
86
87             # Log the plot to wandb
88             wandb.log({f"{layer_name}_Overall_LTP_LTD": wandb.Image(fig)
89                         })
90             plt.close(fig)
91
92         # Plot 2: Detailed weight changes within each filter
93         fig, ax2 = plt.subplots(figsize=(12, 6))
94         data = []
95         filter_indices = []
96         for i in range(min(num_filters, weights.shape[0])):
97             changes = delta_w[i].view(-1).tolist()
98             data.extend(changes)
99             filter_indices.extend([i] * len(changes))
100
101             sns.violinplot(x=filter_indices, y=data, ax=ax2)
102             ax2.axhline(y=0, color='black', linestyle='--', linewidth
103                         =0.5)
104             ax2.set_xlabel('Filter Index')
105             ax2.set_ylabel('Weight Change')
106             ax2.set_title(f'Detailed Weight Changes within first {
107                           num_filters} filters')
108
109             # Log the plot to wandb

```

```

103     wandb.log({f"{layer_name}_Detailed_Weight_Changes": wandb.
104         Image(fig)})
105     plt.close(fig)
106
106 # Plot 3: Detailed statistics for each filter and overall
107 # layer statistics
107 fig, (ax3, ax4) = plt.subplots(2, 1, figsize=(12, 12),
108     gridspec_kw={'height_ratios': [3, 1]})  

108 stats = []
109 for i in range(min(num_filters, weights.shape[0])):
110     filter_weights = weights[i].view(-1) #could be delta_w
111     stats.append({
112         'Mean': filter_weights.mean().item(),
113         'Median': filter_weights.median().item(),
114         'Std Dev': filter_weights.std().item(),
115         '% Positive': (filter_weights > 0).float().mean().
116             item() * 100,
116         '% Negative': (filter_weights < 0).float().mean().
117             item() * 100
117     })
118 # Per-filter statistics
119 stat_df = pd.DataFrame(stats)
120 sns.heatmap(stat_df.T, annot=True, cmap='coolwarm', center
121     =0, ax=ax3)
121 ax3.set_xlabel('Filter Index')
122 ax3.set_title('Detailed Statistics for Each Filter')
123 # Overall layer statistics
124 all_weights = weights.view(-1) # This includes ALL weights
125     in the layer
125 overall_stats = pd.DataFrame({
126     'Layer Overall': {
127         'Mean': all_weights.mean().item(),
128         'Median': all_weights.median().item(),
129         'Std Dev': all_weights.std().item(),
130         '% Positive': (all_weights > 0).float().mean().item
131             () * 100,
131         '% Negative': (all_weights < 0).float().mean().item
132             () * 100
132     }
133 })
134 sns.heatmap(overall_stats, annot=True, cmap='coolwarm',
135     center=0, ax=ax4)
135 ax4.set_title('Overall Layer Statistics')
136 plt.tight_layout()
137 # Log the plot to wandb
138 wandb.log({f"{layer_name}_Weight_Statistics": wandb.Image(
139     fig)})
139 plt.close(fig)
140
141 # Plot 4: LTP/LTD per Weight (mean across channels)
142 num_filters_to_show = min(25, weights.shape[0])
143 if weights.shape[2] == 1 and weights.shape[3] == 1: # Check
144     if kernels are 1x1
145         fig, ax = plt.subplots(figsize=(20, 5))
146         filter_changes = [delta_w[i].mean().item() for i in
147             range(num_filters_to_show)]
148         norm = plt.Normalize(vmin=min(filter_changes), vmax=max(
149             filter_changes))

```

```

147     colors = plt.cm.RdYlGn(norm(filter_changes))
148     ax.bar(range(num_filters_to_show), filter_changes, color
149         =colors)
150     ax.set_xlabel('Filter Index')
151     ax.set_ylabel('Weight Change')
152     ax.set_title('LTP/LTD for 1x1 Kernels')
153     ax.axhline(y=0, color='black', linestyle='--', linewidth
154         =0.5)
155     sm = plt.cm.ScalarMappable(cmap="RdYlGn", norm=norm)
156     sm.set_array([])
157     cbar = plt.colorbar(sm, ax=ax, orientation='horizontal',
158         aspect=30)
159     cbar.set_label('Weight Change')
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181 def calculate_metrics(preds, labels, num_classes):
182     if num_classes == 2:
183         accuracy = Accuracy(task='binary', num_classes=num_classes).
184             to(device)
185         precision = Precision(task='binary', average='weighted',
186             num_classes=num_classes).to(device)
187         recall = Recall(task='binary', average='weighted',
188             num_classes=num_classes).to(device)
189         f1 = F1Score(task='binary', average='weighted', num_classes=
190             num_classes).to(device)
191         confusion_matrix = ConfusionMatrix(task='binary',
192             num_classes=num_classes).to(device)
193
194     else:
195         accuracy = Accuracy(task='multiclass', num_classes=
196             num_classes).to(device)

```

```

190     precision = Precision(task='multiclass', average='macro',
191                           num_classes=num_classes).to(device)
192     recall = Recall(task='multiclass', average='macro',
193                      num_classes=num_classes).to(device)
194     f1 = F1Score(task='multiclass', average='macro', num_classes
195                  =num_classes).to(device)
196     confusion_matrix = ConfusionMatrix(task='multiclass',
197                                         num_classes=num_classes).to(device)
198
199     acc = accuracy(preds, labels)
200     prec = precision(preds, labels)
201     rec = recall(preds, labels)
202     f1_score = f1(preds, labels)
203     conf_matrix = confusion_matrix(preds, labels)
204
205     return acc, prec, rec, f1_score, conf_matrix
206
207
208 def visualize_data_clusters(dataloader, model=None, method='tsne',
209                             dim=2, perplexity=30, n_neighbors=15, min_dist=0.1,
210                                         n_components=2, random_state=42):
211     device = torch.device("cuda" if torch.cuda.is_available() else "
212                           cpu")
213     features_list = []
214     labels_list = []
215     if model is not None:
216         model.eval()
217     with torch.no_grad():
218         for data, labels in dataloader:
219             data = data.to(device)
220             if model is not None:
221                 if hasattr(model, 'features_extract'):
222                     features = model.features_extract(data)
223                 else:
224                     features = model.conv1(data)
225             else:
226                 features = data
227             features = features.view(features.size(0), -1).cpu()._
228                         numpy()
229             features_list.append(features)
230             labels_list.append(labels.numpy())
231     features = np.vstack(features_list)
232     labels = np.concatenate(labels_list)
233     # Normalize features
234     scaler = StandardScaler()
235     features_normalized = scaler.fit_transform(features)
236     # Apply dimensionality reduction
237     if method == 'tsne':
238         reducer = TSNE(n_components=dim, perplexity=perplexity,
239                        n_iter=1000, random_state=random_state)
240     elif method == 'umap':
241         reducer = umap.UMAP(n_neighbors=n_neighbors, min_dist=
242                            min_dist, n_components=dim, random_state=random_state)
243     else:
244         raise ValueError("Method must be either 'tsne' or 'umap'")
245     projected_data = reducer.fit_transform(features_normalized)
246     # Plotting
247     if dim == 2:
248         fig = plt.figure(figsize=(12, 10))

```

```

239     scatter = plt.scatter(projected_data[:, 0], projected_data
240                         [:, 1], c=labels, alpha=0.5, cmap='tab10')
241     plt.colorbar(scatter, label='Class Labels')
242     plt.title(f'CIFAR-10 Data Clusters using {method.upper()} (2
243               D)')
244     plt.xlabel(f'{method.upper()} Component 1')
245     plt.ylabel(f'{method.upper()} Component 2')
246 elif dim == 3:
247     fig = plt.figure(figsize=(12, 10))
248     ax = fig.add_subplot(111, projection='3d')
249     scatter = ax.scatter(projected_data[:, 0], projected_data[:, 1],
250                          projected_data[:, 2], c=labels, alpha=0.5,
251                          cmap='tab10')
252     fig.colorbar(scatter, label='Class Labels')
253     ax.set_title(f'CIFAR-10 Data Clusters using {method.upper()} (3D)')
254     ax.set_xlabel(f'{method.upper()} Component 1')
255     ax.set_ylabel(f'{method.upper()} Component 2')
256     ax.set_zlabel(f'{method.upper()} Component 3')
257 else:
258     raise ValueError("dim must be either 2 or 3")
259
260 # Main file to train backpropagation network
261 # TODO: Same visualisation code as Hebbian is not fully implemented
262
263 if __name__ == "__main__":
264
265     device = torch.device('cuda:0')
266     model = Net_Backpropagation()
267     model.to(device)
268
269     wandb_logger = Logger(
270         f"SoftHebb-BackPropagation-Init",
271         project='Final-HebbianCNN', model=model)
272     logger = wandb_logger.get_logger()
273     num_parameters = sum(p.numel() for p in model.parameters() if p.
274                           requires_grad)
275     print(f"Parameter Count Total: {num_parameters}")
276
277     sup_optimizer = optim.Adam(model.parameters(), lr=0.01)
278     criterion = nn.CrossEntropyLoss()
279
280     trn_set, tst_set, zca = data.get_data(dataset='cifar10', root='
281                                         datasets', batch_size=64,
282                                         whiten_lvl=None)
283
284     print(f'Processing Training batches: {len(trn_set)}')
285     for epoch in range(5):
286         print(f"Training BP epoch {epoch}")
287         model.train()
288         running_loss = 0.0
289         correct = 0
290         total = 0
291         train_preds = []
292         train_labels = []

```

```

291     for i, data in enumerate(trn_set, 0):
292         inputs, labels = data
293         inputs = inputs.to(device)
294         labels = labels.to(device)
295         # zero the parameter gradients
296         sup_optimizer.zero_grad()
297         # forward + backward + optimize
298         outputs = model(inputs)
299         loss = criterion(outputs, labels)
300         loss.backward()
301
302         if i % 200 == 0: # Every 100 datapoint
303             print(f'Saving details after batch {i}')
304             plot_ltp_ltd(model.conv1, 'conv1', num_filters=10,
305                           detailed_mode=True)
306             plot_ltp_ltd(model.conv2, 'conv2', num_filters=10,
307                           detailed_mode=True)
308             # plot_ltp_ltd(model.conv_point2, 'conv_point2',
309                           # num_filters=10, detailed_mode=True)
310             model.visualize_filters('conv1')
311             model.visualize_filters('conv2')
312             # model.visualize_filters('conv_point2')
313
314             sup_optimizer.step()
315             # compute training statistics
316             running_loss += loss.item()
317             total += labels.size(0)
318             _, predicted = torch.max(outputs.data, 1)
319             correct += (predicted == labels).sum().item()
320             # For wandb logs
321             preds = torch.argmax(outputs, dim=1)
322             train_preds.append(preds)
323             train_labels.append(labels)
324
325             # print("Visualizing Receptive fields")
326             # visualize_filters(model, model.conv1, num_filters=25)
327             # visualize_filters(model, model.conv2, num_filters=25)
328             # visualize_filters(model, model.conv3, num_filters=25)
329
330             print(f'Accuracy of the network on the train images: {100 * correct / total} %')
331             print(f'{epoch + 1}] loss: {running_loss / total:.3f}')
332
333             train_preds = torch.cat(train_preds, dim=0)
334             train_labels = torch.cat(train_labels, dim=0)
335             acc, prec, rec, f1_score, conf_matrix = calculate_metrics(
336                 train_preds, train_labels, 10)
337             logger.log({'train_accuracy': acc, 'train_precision': prec,
338                         'train_recall': rec, 'train_f1_score': f1_score})
339             f, ax = plt.subplots(figsize=(15, 10))
340             sns.heatmap(conf_matrix.clone().detach().cpu().numpy(),
341                         annot=True, ax=ax)
342             logger.log({'train_confusion_matrix': wandb.Image(f)})
343             plt.close(f)
344
345             # Evaluation on test set
346             model.eval()
347             running_loss = 0.

```

```

342     correct = 0
343     total = 0
344     # since we're not training, we don't need to calculate the
345     # gradients for our outputs
346     test_preds = []
347     test_labels = []
348     with torch.no_grad():
349         for data in tst_set:
350             images, labels = data
351             images = images.to(device)
352             labels = labels.to(device)
353             # calculate outputs by running images through the
354             # network
355             outputs = model(images)
356             # the class with the highest energy is what we
357             # choose as prediction
358             _, predicted = torch.max(outputs.data, 1)
359             total += labels.size(0)
360             correct += (predicted == labels).sum().item()
361             loss = criterion(outputs, labels)
362             running_loss += loss.item()
363             # For wandb logs
364             preds = torch.argmax(outputs, dim=1)
365             test_preds.append(preds)
366             test_labels.append(labels)

367             print(f'Accuracy of the network on the test images: {100 * 
368                 correct / total} %')
369             print(f'test loss: {running_loss / total:.3f}')

370             test_preds = torch.cat(test_preds, dim=0)
371             test_labels = torch.cat(test_labels, dim=0)
372             acc, prec, rec, f1_score, conf_matrix = calculate_metrics(
373                 test_preds, test_labels, 10)
374             logger.log({'test_accuracy': acc, 'test_precision': prec, ' 
375                 test_recall': rec, 'test_f1_score': f1_score})
376             f, ax = plt.subplots(figsize=(15, 10))
377             sns.heatmap(conf_matrix.clone().detach().cpu().numpy(), 
378                 annot=True, ax=ax)
379             logger.log({"test_confusion_matrix": wandb.Image(f)})
380             plt.close(f)

381             print("Visualizing Filters")
382             model.visualize_filters('conv1', f'results>{"demo"}/ 
383                 demo_conv1_filters_epoch_{1}.png')
384             model.visualize_filters('conv2', f'results>{"demo"}/ 
385                 demo_conv2_filters_epoch_{1}.png')
386             model.visualize_filters('conv3', f'results>{"demo"}/ 
387                 demo_conv3_filters_epoch_{1}.png')
388             # model.visualize_filters('conv_point2', f'results>{"demo"}/ 
389                 demo_conv_point2_filters_epoch_{1}.png')
390             print("Weight statistics")
391             print_weight_statistics(model.conv1, 'conv1')
392             print_weight_statistics(model.conv2, 'conv2')
393             print_weight_statistics(model.conv3, 'conv3')
394             # print_weight_statistics(model.conv_point2, 'conv3')

395             print("Visualizing Class separation")

```

```

389     visualize_data_clusters(tst_set, model=model, method='umap', dim
390                             =2)
391
392     print("Visualizing Receptive fields")
393     visualize_filters(model, model.conv1, num_filters=25)
394     visualize_filters(model, model.conv2, num_filters=25)
395     visualize_filters(model, model.conv3, num_filters=25)

```

A.2.14 Visualizer.py

```

1 import torch
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import wandb
6 import pandas as pd
7 import umap
8 from sklearn.manifold import TSNE
9 from sklearn.preprocessing import StandardScaler
10
11 torch.manual_seed(0)
12
13 def print_weight_statistics(layer, layer_name):
14     """
15         Prints statistics of the weights from a given layer.
16     Args:
17     """
18     weights = layer.weight.data
19
20     print(f"Statistics for {layer_name}:")
21     print(f"Mean: {weights.mean().item():.4f}")
22     print(f"Max: {weights.max().item():.4f}")
23     print(f"Min: {weights.min().item():.4f}")
24     print(f"Standard Deviation: {weights.std().item():.4f}")
25     print(f"Median: {weights.median().item():.4f}")
26     print(f"25th Percentile: {weights.quantile(0.25).item():.4f}")
27     print(f"75th Percentile: {weights.quantile(0.75).item():.4f}")
28     print(f"Number of positive weights: {(weights >= 0).sum().item()}")
29     print(f"Number of negative weights: {(weights < 0).sum().item()}")
30     print(f"Total number of weights: {weights.numel()}")
31     print()
32
33
34 # Visualise different information regarding weights and changes in
35 # weights
36 def plot_ltp_ltd_ex_in(layer, layer_name, num_filters=10,
37                         detailed_mode=False):
38     weights = {
39         'wee': layer.weight_ee.data,
40         'wei': layer.weight_ei.data,
41         'wie': layer.weight_ie.data
42     }
43     delta_w = {
44         'wee': layer.delta_w_ee.data,
45         'wei': layer.delta_w_ei.data,
46         'wie': layer.delta_w_ie.data

```

```

45 }
46
47 if not detailed_mode:
48     fig, axes = plt.subplots(1, 3, figsize=(18, 6))
49     for idx, (weight_type, weight) in enumerate(weights.items()):
50         :
51         ax = axes[idx]
52         delta_w = delta_w[weight_type]
53         for i in range(min(num_filters, weight.shape[0])):
54             weight_change = delta_w[i].sum().item()
55             color = 'green' if weight_change > 0 else 'red'
56             ax.bar(i, weight_change, color=color)
57             ax.axhline(y=0, color='black', linestyle='--', linewidth=0.5)
58             ax.set_xlabel('Filter Index')
59             ax.set_ylabel('Total Weight Change')
60             ax.set_title(f'Overall LTP/LTD for {weight_type} in {layer_name}')
61             plt.tight_layout()
62             wandb.log({f'{layer_name}_Overall_LTP_LTD': wandb.Image(fig)})
63         plt.close(fig)
64
65     else:
66         for weight_type in weights.keys():
67             weight = weights[weight_type]
68             delta_w = delta_w[weight_type]
69
70             # Plot 1: Overall weight change
71             fig, ax1 = plt.subplots(figsize=(12, 6))
72             for i in range(min(num_filters, weight.shape[0])):
73                 weight_change = delta_w[i].sum().item()
74                 color = 'green' if weight_change > 0 else 'red'
75                 ax1.bar(i, weight_change, color=color)
76                 ax1.axhline(y=0, color='black', linestyle='--', linewidth=0.5)
77                 ax1.set_xlabel('Filter Index')
78                 ax1.set_ylabel('Total Weight Change')
79                 ax1.set_title(f'Overall LTP/LTD for {weight_type} in {layer_name}')
80                 wandb.log({f'{layer_name}_{weight_type}_Overall_LTP_LTD':
81                             wandb.Image(fig)})
82             plt.close(fig)
83
84             # Plot 2: Detailed weight changes within each filter
85             fig, ax2 = plt.subplots(figsize=(12, 6))
86             data = []
87             filter_indices = []
88             for i in range(min(num_filters, weight.shape[0])):
89                 changes = delta_w[i].view(-1).tolist()
90                 data.extend(changes)
91                 filter_indices.extend([i] * len(changes))
92             sns.violinplot(x=filter_indices, y=data, ax=ax2)
93             ax2.axhline(y=0, color='black', linestyle='--', linewidth=0.5)
94             ax2.set_xlabel('Filter Index')
95             ax2.set_ylabel('Weight Change')

```

```

94     ax2.set_title(f'Detailed Weight Changes for {weight_type}
95         } in {layer_name}')
96     wandb.log({f'{layer_name}_{weight_type}'
97         '_Detailed_Weight_Changes': wandb.Image(fig)})
98     plt.close(fig)
99
100    # Plot 3: Detailed statistics for each filter and
101        overall layer statistics
102    fig, (ax3, ax4) = plt.subplots(2, 1, figsize=(12, 12),
103        gridspec_kw={'height_ratios': [3, 1]})
104    stats = []
105    for i in range(min(num_filters, weight.shape[0])):
106        filter_weights = weight[i].view(-1)
107        stats.append({
108            'Mean': filter_weights.mean().item(),
109            'Median': filter_weights.median().item(),
110            'Std Dev': filter_weights.std().item(),
111            '% Positive': (filter_weights >= 0).float().mean()
112                .item() * 100,
113            '% Negative': (filter_weights < 0).float().mean()
114                .item() * 100
115        })
116    stat_df = pd.DataFrame(stats)
117    sns.heatmap(stat_df.T, annot=True, cmap='coolwarm',
118        center=0, ax=ax3)
119    ax3.set_xlabel('Filter Index')
120    ax3.set_title(f'Detailed Statistics for Each Filter in {
121        weight_type}')
122    all_weights = weight.view(-1)
123    overall_stats = pd.DataFrame({
124        'Layer Overall': {
125            'Mean': all_weights.mean().item(),
126            'Median': all_weights.median().item(),
127            'Std Dev': all_weights.std().item(),
128            '% Positive': (all_weights > 0).float().mean().
129                item() * 100,
130            '% Negative': (all_weights < 0).float().mean().
131                item() * 100
132        }
133    })
134    sns.heatmap(overall_stats, annot=True, cmap='coolwarm',
135        center=0, ax=ax4)
136    ax4.set_title(f'Overall Layer Statistics for {
137        weight_type}')
138    plt.tight_layout()
139    wandb.log({f'{layer_name}_{weight_type}'
140        '_Weight_Statistics': wandb.Image(fig)})
141    plt.close(fig)
142
143    # Plot 4: LTP/LTD per Weight (mean across channels)
144    num_filters_to_show = min(25, weight.shape[0])
145    if weight.shape[2] == 1 and weight.shape[3] == 1:  #
146        Check if kernels are 1x1
147        fig, ax = plt.subplots(figsize=(20, 5))
148        filter_changes = [delta[i].mean().item() for i in
149            range(num_filters_to_show)]
150        norm = plt.Normalize(vmin=min(filter_changes), vmax=
151            max(filter_changes))

```

```

136         colors = plt.cm.RdYlGn(norm(filter_changes))
137         ax.bar(range(num_filters_to_show), filter_changes,
138                color=colors)
139         ax.set_xlabel('Filter Index')
140         ax.set_ylabel('Weight Change')
141         ax.set_title(f'LTP/LTD for 1x1 Kernels in {'
142                     weight_type}')
143         ax.axhline(y=0, color='black', linestyle='--',
144                     linewidth=0.5)
145         sm = plt.cm.ScalarMappable(cmap="RdYlGn", norm=norm)
146         sm.set_array([])
147         cbar = plt.colorbar(sm, ax=ax, orientation='
148                             horizontal', aspect=30)
149         cbar.set_label('Weight Change')
150     else:
151         rows = int(np.ceil(np.sqrt(num_filters_to_show)))
152         cols = int(np.ceil(num_filters_to_show / rows))
153         fig, axes = plt.subplots(rows, cols, figsize=(20,
154                                         20))
155         axes = axes.flatten()
156         for i in range(num_filters_to_show):
157             filter_changes = delta[i].mean(dim=0).cpu().
158                         numpy() # Mean across channels
159             im = axes[i].imshow(filter_changes, cmap='RdYlGn
160                               ', interpolation='nearest')
161             axes[i].set_title(f'Filter {i}')
162             axes[i].axis('off')
163             plt.colorbar(im, ax=axes[i], fraction=0.046, pad
164                         =0.04)
165             for j in range(i + 1, rows * cols):
166                 axes[j].axis('off')
167             fig.suptitle(f'LTP/LTD per Weight (Mean across
168                         channels) for {weight_type}', fontsize=16)
169             wandb.log({f"{layer_name}_{weight_type}"
170                         '_LTP_LTD_per_Weight': wandb.Image(fig)})
171             plt.close(fig)
172
173     # Plot 5: Detailed statistics for each filter and
174     # overall layer statistics: delta_w
175     fig, (ax3, ax4) = plt.subplots(2, 1, figsize=(12, 12),
176                                 gridspec_kw={'height_ratios': [3, 1]})
177     stats = []
178     for i in range(min(num_filters, weight.shape[0])):
179         filter_weights = delta[i].view(-1)
180         stats.append({
181             'Mean': filter_weights.mean().item(),
182             'Median': filter_weights.median().item(),
183             'Std Dev': filter_weights.std().item(),
184             '% Positive': (filter_weights > 0).float().mean
185                           ().item() * 100,
186             '% Negative': (filter_weights < 0).float().mean
187                           ().item() * 100
188         })
189     stat_df = pd.DataFrame(stats)
190     sns.heatmap(stat_df.T, annot=True, cmap='coolwarm',
191                  center=0, ax=ax3)
192     ax3.set_xlabel('Filter Index')

```

```

178     ax3.set_title(f'Detailed Statistics for Each Filter in {weight_type}')
179     all_weights = delta.view(-1)
180     overall_stats = pd.DataFrame({
181         'Layer Overall': {
182             'Mean': all_weights.mean().item(),
183             'Median': all_weights.median().item(),
184             'Std Dev': all_weights.std().item(),
185             '% Positive': (all_weights > 0).float().mean().item() * 100,
186             '% Negative': (all_weights < 0).float().mean().item() * 100
187         }
188     })
189     sns.heatmap(overall_stats, annot=True, cmap='coolwarm',
190                 center=0, ax=ax4)
190     ax4.set_title(f'Overall Layer Statistics for Delta_w in {weight_type}')
191     plt.tight_layout()
192     wandb.log({f'{layer_name}_{weight_type}_Delta_w_Statistics': wandb.Image(fig)})
193     plt.close(fig)
194
195     # Plot 6: Weight distribution plots
196     fig, ax6 = plt.subplots(figsize=(12, 6))
197     sns.histplot(weight.view(-1).cpu().numpy(), bins=50, kde=True, ax=ax6)
198     ax6.set_xlabel('Weight Value')
199     ax6.set_ylabel('Frequency')
200     ax6.set_title(f'Weight Distribution in {layer_name} for {weight_type}')
201     wandb.log({f'{layer_name}_{weight_type}_Weight_Distribution': wandb.Image(fig)})
202     plt.close(fig)
203
204 def plot_ltp_ltd(layer, layer_name, num_filters=10, detailed_mode=False):
205     weights = layer.weight.data
206     delta_w = layer.delta_w.data
207     # print(f"Layer {layer_name}")
208     # print(delta_w.mean())
209     # print(delta_w.max())
210
211     if not detailed_mode:
212         fig, ax = plt.subplots(figsize=(12, 6))
213
214         for i in range(min(num_filters, weights.shape[0])):
215             weight_change = delta_w[i].sum().item()
216             color = 'green' if weight_change > 0 else 'red'
217             ax.bar(i, weight_change, color=color)
218
219             ax.axhline(y=0, color='black', linestyle='--', linewidth=0.5)
220             ax.set_xlabel('Filter Index')
221             ax.set_ylabel('Total Weight Change')
222             ax.set_title(f'Overall LTP/LTD for first {num_filters} filters in {layer_name}')
223
224     # Log the plot to wandb

```

```

225     wandb.log({f"{layer_name}_Overall_LTP_LTD": wandb.Image(fig)
226         })
227     plt.close(fig)
228
229     # Plot 1: Overall weight change
230     fig, ax1 = plt.subplots(figsize=(12, 6))
231     for i in range(min(num_filters, weights.shape[0])):
232         weight_change = delta_w[i].sum().item()
233         color = 'green' if weight_change > 0 else 'red'
234         ax1.bar(i, weight_change, color=color)
235
236         ax1.axhline(y=0, color='black', linestyle='--', linewidth
237             =0.5)
238         ax1.set_xlabel('Filter Index')
239         ax1.set_ylabel('Total Weight Change')
240         ax1.set_title(f'Overall LTP/LTD for first {num_filters}
241             filters')
242
243     # Log the plot to wandb
244     wandb.log({f"{layer_name}_Overall_LTP_LTD": wandb.Image(fig)
245         })
246     plt.close(fig)
247
248     # Plot 2: Detailed weight changes within each filter
249     fig, ax2 = plt.subplots(figsize=(12, 6))
250     data = []
251     filter_indices = []
252     for i in range(min(num_filters, weights.shape[0])):
253         changes = delta_w[i].view(-1).tolist()
254         data.extend(changes)
255         filter_indices.extend([i] * len(changes))
256
257         sns.violinplot(x=filter_indices, y=data, ax=ax2)
258         ax2.axhline(y=0, color='black', linestyle='--', linewidth
259             =0.5)
260         ax2.set_xlabel('Filter Index')
261         ax2.set_ylabel('Weight Change')
262         ax2.set_title(f'Detailed Weight Changes within first {
263             num_filters} filters')
264
265     # Log the plot to wandb
266     wandb.log({f"{layer_name}_Detailed_Weight_Changes": wandb.
267         Image(fig)})
268     plt.close(fig)
269
270     # Plot 3: Detailed statistics for each filter and overall
271         layer statistics
272     fig, (ax3, ax4) = plt.subplots(2, 1, figsize=(12, 12),
273         gridspec_kw={'height_ratios': [3, 1]})
274     stats = []
275     for i in range(min(num_filters, weights.shape[0])):
276         filter_weights = weights[i].view(-1) #could be delta_w
277         stats.append({
278             'Mean': filter_weights.mean().item(),
279             'Median': filter_weights.median().item(),
280             'Std Dev': filter_weights.std().item(),

```

```

273             '% Positive': (filter_weights > 0).float().mean().
274                 item() * 100,
275             '% Negative': (filter_weights < 0).float().mean().
276                 item() * 100
277         })
278     # Per-filter statistics
279     stat_df = pd.DataFrame(stats)
280     sns.heatmap(stat_df.T, annot=True, cmap='coolwarm', center
281                 =0, ax=ax3)
282     ax3.set_xlabel('Filter Index')
283     ax3.set_title('Detailed Statistics for Each Filter')
284     # Overall layer statistics
285     all_weights = weights.view(-1) # This includes ALL weights
286         in the layer
287     overall_stats = pd.DataFrame({
288         'Layer Overall': {
289             'Mean': all_weights.mean().item(),
290             'Median': all_weights.median().item(),
291             'Std Dev': all_weights.std().item(),
292             '% Positive': (all_weights > 0).float().mean().item()
293                 () * 100,
294             '% Negative': (all_weights < 0).float().mean().item()
295                 () * 100
296         }
297     })
298     sns.heatmap(overall_stats, annot=True, cmap='coolwarm',
299                 center=0, ax=ax4)
300     ax4.set_title('Overall Layer Statistics')
301     plt.tight_layout()
302     # Log the plot to wandb
303     wandb.log({f"{layer_name}_Weight_Statistics": wandb.Image(
304         fig)})
305     plt.close(fig)

# Plot 4: LTP/LTD per Weight (mean across channels)
306     num_filters_to_show = min(25, weights.shape[0])
307     if weights.shape[2] == 1 and weights.shape[3] == 1: # Check
308         if kernels are 1x1
309             fig, ax = plt.subplots(figsize=(20, 5))
310             filter_changes = [delta_w[i].mean().item() for i in
311                 range(num_filters_to_show)]
312             norm = plt.Normalize(vmin=min(filter_changes), vmax=max(
313                 filter_changes))
314             colors = plt.cm.RdYlGn(norm(filter_changes))
315             ax.bar(range(num_filters_to_show), filter_changes, color
316                 =colors)
317             ax.set_xlabel('Filter Index')
318             ax.set_ylabel('Weight Change')
319             ax.set_title('LTP/LTD for 1x1 Kernels')
320             ax.axhline(y=0, color='black', linestyle='--', linewidth
321                 =0.5)
322             sm = plt.cm.ScalarMappable(cmap="RdYlGn", norm=norm)
323             sm.set_array([])
324             cbar = plt.colorbar(sm, ax=ax, orientation='horizontal',
325                 aspect=30)
326             cbar.set_label('Weight Change')
327
328     else:

```

```

317     rows = int(np.ceil(np.sqrt(num_filters_to_show)))
318     cols = int(np.ceil(num_filters_to_show / rows))
319     fig, axes = plt.subplots(rows, cols, figsize=(20, 20))
320     axes = axes.flatten()
321
322     for i in range(num_filters_to_show):
323         filter_changes = delta_w[i].mean(dim=0).cpu().numpy
324             () # Mean across channels
325         im = axes[i].imshow(filter_changes, cmap='RdYlGn',
326             interpolation='nearest')
327         axes[i].set_title(f'Filter {i}')
328         axes[i].axis('off')
329         plt.colorbar(im, ax=axes[i], fraction=0.046, pad
330             =0.04)
331
332     for j in range(i + 1, rows * cols):
333         axes[j].axis('off')
334
335     fig.suptitle('LTP/LTD per Weight (Mean across channels)',
336             , fontsize=16)
337 # Log the plot to wandb
338 wandb.log({f"{layer_name}_LTP_LTD_per_Weight": wandb.Image(
339             fig)})
340 plt.close(fig)
341
342 # Plot 5: Detailed statistics for each filter and overall
343     layer statistics: delta_w
344 fig, (ax3, ax4) = plt.subplots(2, 1, figsize=(12, 12),
345     gridspec_kw={'height_ratios': [3, 1]})
346 stats = []
347 for i in range(min(num_filters, weights.shape[0])):
348     filter_weights = delta_w[i].view(-1) # could be delta_w
349     stats.append({
350         'Mean': filter_weights.mean().item(),
351         'Median': filter_weights.median().item(),
352         'Std Dev': filter_weights.std().item(),
353         '% Positive': (filter_weights > 0).float().mean().
354             item() * 100,
355         '% Negative': (filter_weights < 0).float().mean().
356             item() * 100
357     })
358 # Per-filter statistics
359 stat_df = pd.DataFrame(stats)
360 sns.heatmap(stat_df.T, annot=True, cmap='coolwarm', center
361             =0, ax=ax3)
362 ax3.set_xlabel('Filter Index')
363 ax3.set_title('Detailed Statistics for Each Filter')
364 # Overall layer statistics
365 all_weights = delta_w.view(-1) # This includes ALL weights
366     in the layer
367 overall_stats = pd.DataFrame({
368     'Layer Overall': {
369         'Mean': all_weights.mean().item(),
370         'Median': all_weights.median().item(),
371         'Std Dev': all_weights.std().item(),
372         '% Positive': (all_weights > 0).float().mean().item
373             () * 100,

```

```

362             '% Negative': (all_weights < 0).float().mean().item
363             () * 100
364         }
365     sns.heatmap(overall_stats, annot=True, cmap='coolwarm',
366                 center=0, ax=ax4)
367     ax4.set_title('Overall Layer Statistics for Delta_w')
368     plt.tight_layout()
369     # Log the plot to wandb
370     wandb.log({f"{layer_name}_Delta_w_Statistics": wandb.Image(
371         fig)})})
372     plt.close(fig)
373
374     # Plot 6: Weight distribution plots to understand
375     # connectivity of weights better
376     fig, ax6 = plt.subplots(figsize=(12, 6))
377     sns.histplot(weights.view(-1).cpu().numpy(), bins=50, kde=
378                   True, ax=ax6)
379     ax6.set_xlabel('Weight Value')
380     ax6.set_ylabel('Frequency')
381     ax6.set_title(f'Weight Distribution in {layer_name}')
382     # Log the plot to wandb
383     wandb.log({f"{layer_name}_Weight_Distribution": wandb.Image(
384         fig)})})
385     plt.close(fig)
386
387 # Visualise class separation
388 def visualize_data_clusters(dataloader, model=None, method='tsne',
389                             dim=2, perplexity=30, n_neighbors=15, min_dist=0.1,
390                             n_components=2, random_state=42):
391     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
392     features_list = []
393     labels_list = []
394     if model is not None:
395         model.eval()
396     with torch.no_grad():
397         for data, labels in dataloader:
398             data = data.to(device)
399             if model is not None:
400                 if hasattr(model, 'features_extract'):
401                     features = model.features_extract(data)
402                 else:
403                     features = model.conv1(data)
404             else:
405                 features = data
406             features = features.view(features.size(0), -1).cpu().
407                         numpy()
408             features_list.append(features)
409             labels_list.append(labels.numpy())
410     features = np.vstack(features_list)
411     labels = np.concatenate(labels_list)
412     # Normalize features
413     scaler = StandardScaler()
414     features_normalized = scaler.fit_transform(features)
415     # Apply dimensionality reduction
416     if method == 'tsne':

```

```

411     reducer = TSNE(n_components=dim, perplexity=perplexity,
412                      n_iter=1000, random_state=random_state)
413     elif method == 'umap':
414         reducer = umap.UMAP(n_neighbors=n_neighbors, min_dist=
415                             min_dist, n_components=dim, random_state=random_state)
416     else:
417         raise ValueError("Method must be either 'tsne' or 'umap'")
418     projected_data = reducer.fit_transform(features_normalized)
419     # Plotting
420     if dim == 2:
421         fig = plt.figure(figsize=(12, 10))
422         scatter = plt.scatter(projected_data[:, 0], projected_data
423                              [:, 1], c=labels, alpha=0.5, cmap='tab10')
424         plt.colorbar(scatter, label='Class Labels')
425         plt.title(f'CIFAR-10 Data Clusters using {method.upper()} (2
426                   D)')
427         plt.xlabel(f'{method.upper()} Component 1')
428         plt.ylabel(f'{method.upper()} Component 2')
429     elif dim == 3:
430         fig = plt.figure(figsize=(12, 10))
431         ax = fig.add_subplot(111, projection='3d')
432         scatter = ax.scatter(projected_data[:, 0], projected_data[:, 1],
433                               projected_data[:, 2], c=labels, alpha=0.5,
434                               cmap='tab10')
435         fig.colorbar(scatter, label='Class Labels')
436         ax.set_title(f'CIFAR-10 Data Clusters using {method.upper()} (3D)')
437         ax.set_xlabel(f'{method.upper()} Component 1')
438         ax.set_ylabel(f'{method.upper()} Component 2')
439         ax.set_zlabel(f'{method.upper()} Component 3')
440     else:
441         raise ValueError("dim must be either 2 or 3")
442
443     wandb.log({"Class separation": wandb.Image(fig)})
444     plt.close(fig)

```

A.2.15 receptive_field.py

```

1 import matplotlib.pyplot as plt
2 import torch
3 from torch import nn, optim
4 import torch.nn.functional as F
5 from hebb import HebbianConv2d
6 from hebb_depthwise import HebbianDepthConv2d
7
8 # Manually change if using Dale network
9 # from hebb_abs import HebbianConv2d
10 # from hebb_abs_depthwise import HebbianDepthConv2d
11 import wandb
12
13 # Code to visualise receptive fields
14
15 def get_partial_model(model, target_layer):
16     layers = []
17     for layer in model.children():
18         layers.append(layer)
19         if layer == target_layer:
20             break

```

```

21     return torch.nn.Sequential(*layers)
22
23
24 def calculate_receptive_field(model, target_layer):
25     current_rf = 1
26     current_stride = 1
27
28     for layer in model.children():
29         if isinstance(layer, (nn.Conv2d, HebbianConv2d,
30                             HebbianDepthConv2d)):
31             kernel_size = layer.kernel_size[0] if isinstance(layer.
32                 kernel_size, tuple) else layer.kernel_size
33             stride = layer.stride[0] if isinstance(layer.stride,
34                 tuple) else layer.stride
35             current_rf += (kernel_size - 1) * current_stride
36             current_stride *= stride
37         elif isinstance(layer, (nn.MaxPool2d, nn.AvgPool2d)):
38             kernel_size = layer.kernel_size if isinstance(layer.
39                 kernel_size, int) else layer.kernel_size[0]
40             stride = layer.stride if isinstance(layer.stride, int)
41                 else layer.stride[0]
42             current_rf += (kernel_size - 1) * current_stride
43             current_stride *= stride
44
45             if layer == target_layer:
46                 break
47
48     return current_rf
49
50
51 def get_layer_output(model, x, target_layer):
52     """
53     Forward pass through the model, stopping at the target custom
54     Hebbian layer.
55     """
56
57     for layer in model.children():
58         x = layer(x) # Pass through each layer
59         if layer == target_layer:
60             return x # Return the output of the target Hebbian
61             layer
62     raise ValueError(f"Target layer {target_layer} not found in the
63                     model.")
64
65 def remove_padding(model, target_layer):
66     """
67     Remove padding from layers before the target layer.
68     """
69     for layer in model.children():
70         if isinstance(layer, nn.Conv2d):
71             layer.padding = (0, 0)
72         elif isinstance(layer, (HebbianConv2d, HebbianDepthConv2d)):
73             # For custom SoftHebbConv2d layers, we need to modify
74             # the padding directly
75             layer.padding = 0
76         if isinstance(layer, (nn.MaxPool2d, nn.AvgPool2d)):
77             layer.padding = (0, 0)
78         if layer == target_layer:
79             break
80
81     return model
82
83
84 def gaussian_blur(x, kernel_size=5, sigma=1.0):

```

```

70     channels = x.shape[1]
71     kernel = torch.tensor([
72         [1., 4., 6., 4., 1.],
73         [4., 16., 24., 16., 4.],
74         [6., 24., 36., 24., 6.],
75         [4., 16., 24., 16., 4.],
76         [1., 4., 6., 4., 1.]
77     ], device=x.device).unsqueeze(0).unsqueeze(0) / 256.0
78     kernel = kernel.repeat(channels, 1, 1, 1)
79     padding = kernel_size // 2
80     return F.conv2d(x, kernel, padding=padding, groups=channels)
81
82 class SingleMax:
83     def __init__(self, max_val: float, eps: float):
84         self.max_val = max_val
85         self.eps = eps
86
87     def __call__(self, outputs):
88         if self.max_val is None:
89             return outputs > 0
90         else:
91             return (self.max_val - outputs) < self.eps
92
93 class L2ProjGradientDescent:
94     def __init__(self, steps, random_start=True, rel_stepsize=0.1):
95         self.steps = steps
96         self.random_start = random_start
97         self.rel_stepsize = rel_stepsize
98
99     def get_random_start(self, x0, epsilon):
100        batch_size, c, h, w = x0.shape
101        r = torch.randn(batch_size, c * h * w, device=x0.device)
102        r = r / r.norm(dim=1, keepdim=True)
103        r = r.view_as(x0)
104        return x0 + 0.00001 * epsilon * r
105
106     def normalize_gradient(self, grad):
107         return grad / (grad.view(grad.shape[0], -1).norm(dim=1).view
108                         (-1, 1, 1, 1) + 1e-8)
109
110     def project(self, x, x0, epsilon):
111         delta = x - x0
112         delta = epsilon * delta / delta.view(delta.shape[0], -1).
113             norm(dim=1).view(-1, 1, 1, 1).clamp(min=1e-12)
114         return x0 + delta
115
116     def run(self, model, x0, target_layer, filter_idx, epsilon,
117            criterion):
118         x = x0.clone()
119         if self.random_start:
120             x = self.get_random_start(x0, epsilon)
121         for _ in range(self.steps):
122             x.requires_grad_(True)
123             x_smooth = gaussian_blur(x, sigma=1.0) # Apply Gaussian
124                 smoothing
125             activation = get_layer_output(model, x, target_layer)
126             loss = -activation[0, filter_idx].sum()
127             if criterion(loss.item()):

```

```

124         break
125     grad = torch.autograd.grad(loss, x)[0]
126     grad = self.normalize_gradient(grad)
127     with torch.no_grad():
128         x = x - self.rel_stepsize * epsilon * grad
129         x = self.project(x, x0, epsilon)
130         x.clamp_(0, 1)
131     return x
132
133 def visualize_filters(model, layer, num_filters=25, input_shape=(1,
134                         3, 32, 32), step_size=0.001, iterations=500,
135                         random_start=True, nb_start=1, l2_norm=True,
136                         max_val=1e10, eps=1e-05, epsilon=None):
137     if epsilon is None:
138         epsilon = torch.tensor([255.0], device='cuda') / 255.0
139     else:
140         epsilon = torch.tensor(epsilon, device='cuda')
141     model.eval()
142     model = remove_padding(model, layer)
143     receptive_field_size = calculate_receptive_field(model, layer)
144     print(f'Receptive field size: {receptive_field_size}x{receptive_field_size}')
145     input_shape = (1, 3, receptive_field_size, receptive_field_size)
146     if hasattr(layer, 'out_channels'):
147         out_channels = layer.out_channels
148     else:
149         raise ValueError("The target layer does not have an 'out_channels' attribute.")
150     num_filters = min(num_filters, out_channels)
151     filter_images = []
152     criterion = SingleMax(max_val, eps)
153     pgd = L2ProjGradientDescent(steps=iterations, random_start=random_start,
154                                 rel_stepsize=step_size)
155
156     for filter_idx in range(num_filters):
157         best_image = None
158         best_activation = float('-inf')
159         for start in range(nb_start):
160             x0 = torch.zeros(input_shape, device='cuda',
161                             requires_grad=True)
162             optimized_image = pgd.run(model, x0, layer, filter_idx,
163                                       epsilon, criterion)
164             activation = -get_layer_output(model, optimized_image,
165                                            layer)[0, filter_idx].sum().item()
166             if activation > best_activation:
167                 print(f'New best Receptive Field found for filter {filter_idx} in Reboot {start}')
168                 best_activation = activation
169                 best_image = optimized_image
170             optimized_image = best_image.cpu().squeeze(0).permute(1, 2, 0)
171             optimized_image = (optimized_image - optimized_image.min()) /
172                               (optimized_image.max() - optimized_image.min())
173             filter_images.append(optimized_image)
174
175     # Plot the filter visualizations in a grid
176     grid_size = int(num_filters ** 0.5) + (1 if num_filters ** 0.5 %
177                                             1 > 0 else 0)
178     fig, axes = plt.subplots(grid_size, grid_size, figsize=(20, 20))

```

```
170     axes = axes.flatten()
171     for i in range(num_filters):
172         axes[i].imshow(filter_images[i].numpy())
173         axes[i].set_title(f'Filter {i + 1}')
174         axes[i].axis('off')
175     # Turn off unused subplots
176     for j in range(num_filters, len(axes)):
177         axes[j].axis('off')
178     plt.tight_layout()
179     wandb.log({"Receptive Fields": wandb.Image(fig)})
180     plt.close(fig)
```

A.2.16 receptive_field_residual.py

```

1 import matplotlib.pyplot as plt
2 import torch
3 from torch import nn
4 import torch.nn.functional as F
5 from hebb import HebbianConv2d
6 from hebb_depthwise import HebbianDepthConv2d
7 from model_residual import HebbianResidualBlock
8 import wandb
9
10 # Code to visualise receptive fields. Modified version of
11 # receptive_fields.py to work with residual blocks
12
13 def remove_padding_except_conv2(model):
14     for module in model.modules():
15         if isinstance(module, HebbianResidualBlock):
16             module.conv1.padding = 0
17             # Keep padding for conv2
18             module.conv3.padding = 0
19             if isinstance(module.shortcut, nn.Sequential):
20                 for layer in module.shortcut:
21                     if isinstance(layer, (nn.Conv2d, HebbianConv2d)):
22                         :
23                         layer.padding = 0
24             elif isinstance(module, (nn.Conv2d, HebbianConv2d,
25             HebbianDepthConv2d)):
26                 module.padding = 0
27             elif isinstance(module, (nn.MaxPool2d, nn.AvgPool2d)):
28                 module.padding = 0
29
30 def get_partial_model(model, target_layer):
31     layers = []
32     for name, module in model.named_modules():
33         if isinstance(module, (nn.Conv2d, HebbianConv2d,
34             HebbianDepthConv2d, HebbianResidualBlock, nn.MaxPool2d,
35             nn.BatchNorm2d)):
36             if isinstance(module, HebbianResidualBlock):
37                 layers.append((f"{name}.bn1", module.bn1))
38                 layers.append((f"{name}.conv1", module.conv1))
39                 layers.append((f"{name}.bn2", module.bn2))
40                 layers.append((f"{name}.conv2", module.conv2))
41                 layers.append((f"{name}.bn3", module.bn3))
42                 layers.append((f"{name}.conv3", module.conv3))
43                 if isinstance(module.shortcut, nn.Sequential):
44
45

```

```

40             for i, shortcut_layer in enumerate(module.
41                 shortcut):
42                 layers.append(f"{name}.shortcut.{i}",
43                     shortcut_layer)
44             else:
45                 layers.append((name, module))
46         if module == target_layer:
47             break
48     return layers
49
50
51
52
53     def calculate_receptive_field(model, target_layer):
54         current_rf = 1
55         current_stride = 1
56
57         print(f"Calculating receptive field for {target_layer}")
58         print(f"Initial RF: {current_rf}, Initial stride: {current_stride}")
59
60         for name, module in model.named_modules():
61             if isinstance(module, (nn.Conv2d, HebbianConv2d,
62                 HebbianDepthConv2d)):
63                 kernel_size = module.kernel_size[0] if isinstance(module.
64                     .kernel_size, tuple) else module.kernel_size
65                 stride = module.stride[0] if isinstance(module.stride,
66                     tuple) else module.stride
67                 prev_rf = current_rf
68                 current_rf += (kernel_size - 1) * current_stride
69                 prev_stride = current_stride
70                 current_stride *= stride
71                 print(f"Layer: {name}, Type: Conv, Kernel: {kernel_size},
72                     Stride: {stride}")
73                 print(f"    RF: {prev_rf} -> {current_rf}, Stride: {prev_stride} -> {current_stride}")
74             elif isinstance(module, nn.MaxPool2d):
75                 kernel_size = module.kernel_size if isinstance(module.
76                     .kernel_size, int) else module.kernel_size[0]
77                 stride = module.stride if isinstance(module.stride, int)
78                     else module.stride[0]
79                 prev_rf = current_rf
80                 current_rf += (kernel_size - 1) * current_stride
81                 prev_stride = current_stride
82                 current_stride *= stride
83                 print(f"Layer: {name}, Type: MaxPool, Kernel: {
84                     kernel_size}, Stride: {stride}")
85                 print(f"    RF: {prev_rf} -> {current_rf}, Stride: {prev_stride} -> {current_stride}")
86             else:
87                 print(f"Layer: {name}, Type: {type(module).__name__} (no
88                     effect on RF)")
89
90         if module == target_layer:
91             print(f"Reached target layer. Final receptive field: {
92                 current_rf}x{current_rf}")
93             return current_rf
94
95     print(f"Target layer not found. Final receptive field: {
96                 current_rf}x{current_rf}")

```

```

83     return current_rf
84
85
86 def get_layer_output(model, x, target_layer, debug=False):
87     if debug:
88         print(f"Initial input shape: {x.shape}")
89
90     for name, module in model.named_children():
91         if debug:
92             print(f"Processing module: {name}")
93
94         if isinstance(module, HebbianResidualBlock):
95             if debug:
96                 print(f"Entering {name}, input shape: {x.shape}")
97             residual = x
98             out = module.activ(module.conv1(module.bn1(x)))
99             if debug:
100                 print(f"  After conv1: {out.shape}")
101             out = module.activ(module.conv2(module.bn2(out)))
102             if debug:
103                 print(f"  After conv2: {out.shape}")
104             out = module.conv3(module.bn3(out))
105             if debug:
106                 print(f"  After conv3: {out.shape}")
107             shortcut_output = module.shortcut(residual)
108             if debug:
109                 print(f"  Shortcut output: {shortcut_output.shape}")
110
111             # Ensure out and shortcut_output have the same size
112             if out.size() != shortcut_output.size():
113                 min_size = min(out.size(2), shortcut_output.size(2))
114                 out = out[:, :, :min_size, :min_size]
115                 shortcut_output = shortcut_output[:, :, :min_size, :
116                                 min_size]
117
118             x = module.activ(out + shortcut_output)
119             if debug:
120                 print(f"Exiting {name}, output shape: {x.shape}")
121             elif isinstance(module, (nn.Conv2d, HebbianConv2d,
122                                     HebbianDepthConv2d, nn.MaxPool2d, nn.AvgPool2d)):
123                 if debug:
124                     print(f"Entering {name}, input shape: {x.shape}")
125                 x = module(x)
126                 if debug:
127                     print(f"Exiting {name}, output shape: {x.shape}")
128             elif isinstance(module, nn.BatchNorm2d):
129                 if debug:
130                     print(
131                         f"Entering {name} (BatchNorm2d), input shape: {x.
132                                         .shape}, expected features: {module.
133                                         num_features}")
134                 x = module(x)
135                 if debug:
136                     print(f"Exiting {name} (BatchNorm2d), output shape:
137                         {x.shape}")
138
139             if module == target_layer:
140                 if debug:

```

```

136         print(f'Reached target layer {name}, output shape: {x.shape}')
137     return x
138
139     # If the target is inside a HebbianResidualBlock, we need to
140     # look inside it
141     if isinstance(module, HebbianResidualBlock):
142         for sub_name, sub_module in module.named_children():
143             if sub_module == target_layer:
144                 if debug:
145                     print(f'Reached target layer {name}.{sub_name}, output shape: {x.shape}')
146             return x
147
148     raise ValueError(f'Target layer {target_layer} not found in the
149     model.')
150
151 def gaussian_blur(x, kernel_size=5, sigma=1.0):
152     channels = x.shape[1]
153     kernel = torch.tensor([
154         [1., 4., 6., 4., 1.],
155         [4., 16., 24., 16., 4.],
156         [6., 24., 36., 24., 6.],
157         [4., 16., 24., 16., 4.],
158         [1., 4., 6., 4., 1.]
159     ], device=x.device).unsqueeze(0).unsqueeze(0) / 256.0
160     kernel = kernel.repeat(channels, 1, 1, 1)
161     padding = kernel_size // 2
162     return F.conv2d(x, kernel, padding=padding, groups=channels)
163
164 class SingleMax:
165     def __init__(self, max_val: float, eps: float):
166         self.max_val = max_val
167         self.eps = eps
168
169     def __call__(self, outputs):
170         if self.max_val is None:
171             return outputs > 0
172         else:
173             return (self.max_val - outputs) < self.eps
174
175 class L2ProjGradientDescent:
176     def __init__(self, steps, random_start=True, rel_stepsize=0.1):
177         self.steps = steps
178         self.random_start = random_start
179         self.rel_stepsize = rel_stepsize
180
181     def get_random_start(self, x0, epsilon):
182         batch_size, c, h, w = x0.shape
183         r = torch.randn(batch_size, c * h * w, device=x0.device)
184         r = r / r.norm(dim=1, keepdim=True)
185         r = r.view_as(x0)
186         return x0 + 0.00001 * epsilon * r
187
188     def normalize_gradient(self, grad):
189         return grad / (grad.view(grad.shape[0], -1).norm(dim=1).view
190                         (-1, 1, 1, 1) + 1e-8)

```

```

189
190     def project(self, x, x0, epsilon):
191         delta = x - x0
192         delta = epsilon * delta / delta.view(delta.shape[0], -1).
193             norm(dim=1).view(-1, 1, 1, 1).clamp(min=1e-12)
194         return x0 + delta
195
196
197     def run(self, model, x0, target_layer, filter_idx, epsilon,
198           criterion, debug=False):
199         x = x0.clone()
200         if self.random_start:
201             x = self.get_random_start(x0, epsilon)
202         for _ in range(self.steps):
203             x.requires_grad_(True)
204             x_smooth = gaussian_blur(x, sigma=1.0) # Apply Gaussian
205                 smoothing
206             activation = get_layer_output(model, x, target_layer,
207                 debug=debug)
208             loss = -activation[0, filter_idx].sum()
209             if criterion(loss.item()):
210                 break
211             grad = torch.autograd.grad(loss, x)[0]
212             grad = self.normalize_gradient(grad)
213             with torch.no_grad():
214                 x = x - self.rel_stepsize * epsilon * grad
215                 x = self.project(x, x0, epsilon)
216                 x.clamp_(0, 1)
217         return x
218
219
220     def visualize_filters(model, target_layer, num_filters=25,
221                           input_shape=(1, 3, 32, 32), step_size=0.001, iterations=500,
222                           random_start=True, nb_start=1, l2_norm=True,
223                           max_val=1e10, eps=1e-05, epsilons=None,
224                           debug=False):
225         if epsilons is None:
226             epsilons = torch.tensor([255.0], device='cuda') / 255.0
227         else:
228             epsilons = torch.tensor(epsilons, device='cuda')
229         model.eval()
230         remove_padding_except_conv2(model)
231         receptive_field_size = calculate_receptive_field(model,
232             target_layer)
233         print(f'Receptive field size: {receptive_field_size}x{{
234             receptive_field_size}}')
235         input_shape = (1, 3, receptive_field_size, receptive_field_size)
236
237         if isinstance(target_layer, HebbianResidualBlock):
238             out_channels = target_layer.conv3.out_channels
239         elif hasattr(target_layer, 'out_channels'):
240             out_channels = target_layer.out_channels
241         else:
242             raise ValueError("The target layer does not have an 'out_channels' attribute.")
243
244         num_filters = min(num_filters, out_channels)
245         filter_images = []
246         criterion = SingleMax(max_val, eps)

```

```

237 pgd = L2ProjGradientDescent(steps=iterations, random_start=
238     random_start, rel_stepsize=step_size)
239
240     for filter_idx in range(num_filters):
241         best_image = None
242         best_activation = float('-inf')
243         for start in range(nb_start):
244             x0 = torch.zeros(input_shape, device='cuda',
245                 requires_grad=True)
246             optimized_image = pgd.run(model, x0, target_layer,
247                 filter_idx, epsilon, criterion, debug=debug)
248             activation = -get_layer_output(model, optimized_image,
249                 target_layer, debug=debug)[0, filter_idx].sum().item()
250             if activation > best_activation:
251                 print(f"New best Receptive Field found for filter {filter_idx} in Reboot {start}")
252                 best_activation = activation
253                 best_image = optimized_image
254             optimized_image = best_image.cpu().squeeze(0).permute(1, 2,
255                 0)
256             optimized_image = (optimized_image - optimized_image.min())
257                 / (optimized_image.max() - optimized_image.min())
258             filter_images.append(optimized_image)
259
260     # Plot the filter visualizations in a grid
261     grid_size = int(num_filters ** 0.5) + (1 if num_filters ** 0.5 %
262         1 > 0 else 0)
263     fig, axes = plt.subplots(grid_size, grid_size, figsize=(20, 20))
264     axes = axes.flatten()
265     for i in range(num_filters):
266         axes[i].imshow(filter_images[i].numpy())
267         axes[i].set_title(f'Filter {i + 1}')
268         axes[i].axis('off')
269     # Turn off unused subplots
270     for j in range(num_filters, len(axes)):
271         axes[j].axis('off')
272     plt.tight_layout()
273     wandb.log({"Receptive Fields": wandb.Image(fig)})
274     plt.close(fig)

```

A.2.17 params.py

```

1 import torch
2
3 AVAILABLE_DEVICES = ['cpu']
4 if torch.cuda.is_available(): AVAILABLE_DEVICES += ['cuda:{}'.format
4     (d) for d in range(torch.cuda.device_count())]
5 DEFAULT_DEVICE = 'cuda:0' #'cpu'
6 NUM_WORKERS = 4
7 DIST_BINS = 20

```

A.3 Architectures

Table 5: SoftHebb CNN Architecture (Non-Depthwise)

Layer	Type	Output Shape	Kernel	Stride	Padding	Activation
Input	-	(3, 32, 32)	-	-	-	-
1	BatchNorm2d	(3, 32, 32)	-	-	-	-
	HebbianConv2d	(96, 32, 32)	5x5	1	2	-
	Triangle	(96, 32, 32)	-	-	-	power=0.7
	MaxPool2d	(96, 16, 16)	4x4	2	1	-
2	BatchNorm2d	(96, 16, 16)	-	-	-	-
	HebbianConv2d	(384, 16, 16)	3x3	1	1	-
	Triangle	(384, 16, 16)	-	-	-	power=1.4
	MaxPool2d	(384, 8, 8)	4x4	2	1	-
3	BatchNorm2d	(384, 8, 8)	-	-	-	-
	HebbianConv2d	(1536, 8, 8)	3x3	1	1	-
	Triangle	(1536, 8, 8)	-	-	-	power=1.0
	AvgPool2d	(1536, 4, 4)	2x2	2	0	-
Output	Linear	(10)	-	-	-	-

Table 6: HardHebb CNN Architecture (Non-Depthwise)

Layer	Type	Output Shape	Kernel	Stride	Padding	Activation
Input	-	(3, 32, 32)	-	-	-	-
1	BatchNorm2d	(3, 32, 32)	-	-	-	-
	HebbianConv2d	(96, 28, 28)	5x5	1	0	Cosine
	Triangle	(96, 28, 28)	-	-	-	power=0.7
	MaxPool2d	(96, 14, 14)	2x2	2	0	-
2	BatchNorm2d	(96, 14, 14)	-	-	-	-
	HebbianConv2d	(384, 12, 12)	3x3	1	0	Cosine
	Triangle	(384, 12, 12)	-	-	-	power=1.4
3	BatchNorm2d	(384, 12, 12)	-	-	-	-
	HebbianConv2d	(1536, 10, 10)	3x3	1	0	Cosine
	Triangle	(1536, 10, 10)	-	-	-	power=1.0
	AvgPool2d	(1536, 5, 5)	2x2	2	0	-
Output	Linear	(10)	-	-	-	-

* Cosine = Cosine Similarity between weights and input

Table 7: Lagani Short CNN Architecture (Non-Depthwise)

Layer	Type	Output Shape	Kernel	Stride	Padding	Activation
Input	-	(3, 32, 32)	-	-	-	-
1	BatchNorm2d	(3, 32, 32)	-	-	-	-
	HebbianConv2d	(96, 28, 28)	5x5	1	0	Cosine
	Triangle	(96, 28, 28)	-	-	-	power=1.0
	MaxPool2d	(96, 14, 14)	2x2	2	0	-
2	BatchNorm2d	(96, 14, 14)	-	-	-	-
	HebbianConv2d	(128, 12, 12)	3x3	1	0	Cosine
	Triangle	(128, 12, 12)	-	-	-	power=1.0
3	BatchNorm2d	(128, 12, 12)	-	-	-	-
	HebbianConv2d	(192, 10, 10)	3x3	1	0	Cosine
	Triangle	(192, 10, 10)	-	-	-	power=1.0
	AvgPool2d	(192, 5, 5)	2x2	2	0	-
Output	Linear	(10)	-	-	-	-

Table 8: Lagani Long CNN Architecture (Non-Depthwise)

Layer	Type	Output Shape	Kernel	Stride	Padding	Activation
Input	-	(3, 32, 32)	-	-	-	-
1	BatchNorm2d	(3, 32, 32)	-	-	-	-
	HebbianConv2d	(96, 28, 28)	5x5	1	0	Cosine
	Triangle	(96, 28, 28)	-	-	-	power=1.0
	MaxPool2d	(96, 14, 14)	2x2	2	0	-
2	BatchNorm2d	(96, 14, 14)	-	-	-	-
	HebbianConv2d	(128, 12, 12)	3x3	1	0	Cosine
	Triangle	(128, 12, 12)	-	-	-	power=1.0
3	BatchNorm2d	(128, 12, 12)	-	-	-	-
	HebbianConv2d	(192, 10, 10)	3x3	1	0	Cosine
	Triangle	(192, 10, 10)	-	-	-	power=1.0
	AvgPool2d	(192, 5, 5)	2x2	2	0	-
4	BatchNorm2d	(192, 5, 5)	-	-	-	-
	HebbianConv2d	(256, 3, 3)	3x3	1	0	Cosine
	Triangle	(256, 3, 3)	-	-	-	power=1.0
Output	Linear	(10)	-	-	-	-

Table 9: SoftHebb CNN Architecture (Depthwise)

Layer	Type	Output Shape	Kernel	Stride	Padding	Activation
Input	-	(3, 32, 32)	-	-	-	-
1	BatchNorm2d	(3, 32, 32)	-	-	-	-
	HebbianConv2d	(96, 32, 32)	5x5	1	2	-
	Triangle	(96, 32, 32)	-	-	-	power=0.7
	MaxPool2d	(96, 16, 16)	4x4	2	1	-
2	BatchNorm2d	(96, 16, 16)	-	-	-	-
	HebbianDepthConv2d	(96, 16, 16)	3x3	1	1	-
	BatchNorm2d	(96, 16, 16)	-	-	-	-
	HebbianConv2d	(384, 16, 16)	1x1	1	0	-
	Triangle	(384, 16, 16)	-	-	-	power=1.4
	MaxPool2d	(384, 8, 8)	4x4	2	1	-
3	BatchNorm2d	(384, 8, 8)	-	-	-	-
	HebbianDepthConv2d	(384, 8, 8)	3x3	1	1	-
	BatchNorm2d	(384, 8, 8)	-	-	-	-
	HebbianConv2d	(1536, 8, 8)	1x1	1	0	-
	Triangle	(1536, 8, 8)	-	-	-	power=1.0
	AvgPool2d	(1536, 4, 4)	2x2	2	0	-
Output	Linear	(10)	-	-	-	-

Table 10: HardHebb CNN Architecture (Depthwise)

Layer	Type	Output Shape	Kernel	Stride	Padding	Activation
Input	-	(3, 32, 32)	-	-	-	-
1	BatchNorm2d	(3, 32, 32)	-	-	-	-
	HebbianConv2d	(96, 28, 28)	5x5	1	0	Cosine power=0.7
	Triangle	(96, 28, 28)	-	-	-	-
	MaxPool2d	(96, 14, 14)	2x2	2	0	-
2	BatchNorm2d	(96, 14, 14)	-	-	-	-
	HebbianDepthConv2d	(96, 12, 12)	3x3	1	0	Cosine
	BatchNorm2d	(96, 12, 12)	-	-	-	-
	HebbianConv2d	(384, 12, 12)	1x1	1	0	Cosine
	Triangle	(384, 12, 12)	-	-	-	power=1.4
3	BatchNorm2d	(384, 12, 12)	-	-	-	-
	HebbianDepthConv2d	(384, 10, 10)	3x3	1	0	Cosine
	BatchNorm2d	(384, 10, 10)	-	-	-	-
	HebbianConv2d	(1536, 10, 10)	1x1	1	0	Cosine
	Triangle	(1536, 10, 10)	-	-	-	power=1.0
	AvgPool2d	(1536, 5, 5)	2x2	2	0	-
Output	Linear	(10)	-	-	-	-

Table 11: Net_Depthwise_Residual CNN Architecture

Layer	Type	Output Shape	Kernel	Stride	Padding	Activation
Input	-	(3, 32, 32)	-	-	-	-
1	BatchNorm2d	(3, 32, 32)	-	-	-	-
	HebbianConv2d	(96, 32, 32)	5x5	1	2	Cosine power=0.7
	Triangle	(96, 32, 32)	-	-	-	-
	MaxPool2d	(96, 16, 16)	4x4	2	1	-
2	HebbianResidualBlock	(384, 16, 16)	-	-	-	power=1.4
	MaxPool2d	(384, 8, 8)	4x4	2	1	-
3	HebbianResidualBlock	(1536, 8, 8)	-	-	-	power=1.0
	AvgPool2d	(1536, 4, 4)	2x2	2	0	-
Output	Linear	(10)	-	-	-	-
HebbianResidualBlock Internal Structure						
Main Path	BatchNorm2d	(in_ch, H, W)	-	-	-	-
	HebbianConv2d	(hidden_dim, H, W)	1x1	1	0	Cosine
	Triangle	(hidden_dim, H, W)	-	-	-	power=act
	BatchNorm2d	(hidden_dim, H, W)	-	-	-	-
	HebbianDepthConv2d	(hidden_dim, H, W)	3x3	1	1	Cosine
	Triangle	(hidden_dim, H, W)	-	-	-	power=act
	BatchNorm2d	(hidden_dim, H, W)	-	-	-	-
	HebbianConv2d	(out_ch, H, W)	1x1	1	0	Cosine
Shortcut	BatchNorm2d*	(in_ch, H, W)	-	-	-	-
	HebbianConv2d*	(out_ch, H, W)	1x1	1	0	Cosine
	Add	(out_ch, H, W)	-	-	-	-
	Triangle	(out_ch, H, W)	-	-	-	power=act

* Only if in_channels \neq out_channels