

Dynamically Modular and Sparse General Continual Learning

Arnav Varma¹, Elahe Arani^{1,2,†} and Bahram Zonooz^{1,2,†}

¹Advanced Research Lab, NavInfo Europe, Eindhoven, The Netherlands

²Department of Mathematics and Computer Science, Eindhoven University of Technology, The Netherlands

Keywords: Dynamic Neural Networks, Policy Gradients, Lifelong Learning.

Abstract: Real-world applications often require learning continuously from a stream of data under ever-changing conditions. When trying to learn from such non-stationary data, deep neural networks (DNNs) undergo catastrophic forgetting of previously learned information. Among the common approaches to avoid catastrophic forgetting, rehearsal-based methods have proven effective. However, they are still prone to forgetting due to task-interference as all parameters respond to all tasks. To counter this, we take inspiration from sparse coding in the brain and introduce dynamic modularity and sparsity (*Dynamos*) for rehearsal-based general continual learning. In this setup, the DNN learns to respond to stimuli by activating relevant subsets of neurons. We demonstrate the effectiveness of *Dynamos* on multiple datasets under challenging continual learning evaluation protocols. Finally, we show that our method learns representations that are modular and specialized, while maintaining reusability by activating subsets of neurons with overlaps corresponding to the similarity of stimuli. The code is available at <https://github.com/NeurAI-Lab/DynamicContinualLearning>.

1 INTRODUCTION

Deep neural networks (DNNs) have achieved human-level performance in several applications (Greenwald et al., 2021; Taigman et al., 2014). These networks are trained on the multiple tasks within an application with the data being received under an independent and identically distributed (i.i.d.) assumption. This assumption is satisfied by shuffling the data from all tasks and balancing and normalizing the samples from each task in the application (Hadsell et al., 2020). Consequently, DNNs can achieve human-level performance on all tasks in these applications by modeling the joint distribution of the data as a stationary process. Humans, on the other hand, can model the world from inherently non-stationary and sequential observations (French, 1999). Learning continually from the more realistic sequential and non-stationary data is crucial for many applications such as lifelong learning robots (Thrun and Mitchell, 1995) and self-driving cars (Nose et al., 2019). However, vanilla gradient-based training for such continual learning setups with a continuous stream of tasks and data leads to task interference in the DNN's parameters, and consequently, catastrophic forgetting on old tasks (McCloskey and Cohen, 1989; Kirkpatrick et al., 2017). Therefore, there is a need for methods to

alleviate catastrophic forgetting in continual learning.

Previous works have aimed to address these challenges in continual learning. These can be broadly classified into three categories. First, regularization-based methods (Kirkpatrick et al., 2017; Schwarz et al., 2018; Zenke et al., 2017) that penalize changes to the parameters of DNNs to reduce task interference. Second, parameter isolation methods (Adel et al., 2020) that assign distinct subsets of parameters to different tasks. Finally, rehearsal-based methods (Chaudhry et al., 2019) that co-train on current and stored previous samples. Among these, regularization-based and parameter isolation-based methods often require additional information (such as task-identity at test time and task-boundaries during training), or unconstrained growth of networks. These requirements fail to meet general continual learning (GCL) desiderata (Delange et al., 2021; Farquhar and Gal, 2018), making these methods unsuitable for GCL.

Although rehearsal-based methods improve over other categories and meet GCL desiderata, they still suffer from catastrophic forgetting through task interference in the DNN parameters, as all parameters respond to all examples and tasks. This could be resolved by inculcating task or example specific parameter isolation in the rehearsal-based methods. However, it is worth noting that unlike parameter isolation

[†]Contributed equally

methods, modularity and sparsity in the brain is not static. There is evidence that the brain responds to stimuli in a dynamic and sparse manner, with different modules or subsets of neurons responding "dynamically" to different stimuli (Graham and Field, 2006). The advantages of a dynamic and sparse response to stimuli have been explored in deep learning in stationary settings through mechanisms such as gating of modules (Veit and Belongie, 2018), early-exiting (Li et al., 2017; Hu et al., 2020), and dynamic routing (Wang et al., 2018), along with training losses that incentivize sparsity of neural activations (Wu et al., 2018). These studies observed that DNNs trained to predict dynamically also learn to respond differently to different inputs. Furthermore, the learned DNNs demonstrate clustering of parameters in terms of tasks such as similarity, difficulty, and resolution of inputs (Wang et al., 2018; Veit and Belongie, 2018), indicating dynamic modularity. Hence, we hypothesize that combining rehearsal-based methods with dynamic sparsity and modularity could help further mitigate catastrophic forgetting in a more biologically plausible fashion while adhering to GCL desiderata.

To this end, we propose Dynamic Modularity and Sparsity (*Dynamos*), a general continual learning algorithm that combines rehearsal-based methods with dynamic modularity and sparsity. Concretely, we seek to achieve three objectives: dynamic and sparse response to inputs with specialized modules, competent performance, and reducing catastrophic forgetting. To achieve dynamic and sparse responses to inputs, we define multiple agents in our DNN, each responsible for dynamically zeroing out filter activations of a convolutional layer based on the input to that layer. The agents are rewarded for choosing actions that remove activations (sparse responses) if the network predictions are accurate, but are penalized heavily for choosing actions that lead to inaccurate predictions. Agents also rely on prototype losses to learn specialized features. To reduce forgetting and achieve competent performance, we maintain a constant-size memory buffer in which we store previously seen examples. The network is retrained on previous examples alongside current examples to both maintain performance on current and previous tasks, as well as to enforce consistency between current and previous responses to stimuli. *Dynamos* demonstrates competent performance on multiple continual learning datasets under multiple evaluation protocols, including general continual learning. Additionally, our method demonstrates similar and overlapping responses for similar inputs and disparate responses for dissimilar inputs. Finally, we demonstrate that our

method can simulate the trial-to-trial variability observed in humans (Faisal et al., 2008; Werner and Mountcastle, 1963).

2 RELATED WORK

Research in deep learning has approached the dynamic compositionality and sparsity observed in the human brain through dynamic neural networks, where different subsets of neurons or different sub-networks are activated for different stimuli (Bengio et al., 2015; Bolukbasi et al., 2017). This can be achieved through early exiting (Hu et al., 2020), dynamic routing through mixtures of experts or multiple branches (Collier et al., 2020; Wang et al., 2022), and through gating of modules (Wang et al., 2018). Early-exiting might force the DNN to learn specific features in its earlier layers and consequently hurt performance (Wu et al., 2018) as the earlier layers of DNNs are known to learn general purpose features (Yosinski et al., 2014). Dynamic routing, on the other hand, would require the growth of new experts in response to new tasks that risk unconstrained growth, or the initialization of a larger DNN with branches corresponding to the expected number of tasks (Chen et al., 2020). Dynamic networks with gating mechanisms, meanwhile, have been shown to achieve competent performance in i.i.d. training with standard DNNs embedded with small gating networks (Veit and Belongie, 2018; Wu et al., 2018; Wang et al., 2018). These gating networks emit a discrete keep/drop decision for each module, depending on the input to the module or the DNN. As this operation is non-differentiable, a Gumbel Softmax approximation (Veit and Belongie, 2018; Wang et al., 2018), or an agent trained with policy gradients (Wu et al., 2018; Sutton and Barto, 2018) is commonly used in each module to enable backpropagation. However, unlike the latter, the Gumbel-Softmax approximation induces an asymmetry between the forward pass activations at inference and training (Wang et al., 2018). Furthermore, these methods are not applicable to continual learning.

Recent works have attempted to build dynamic networks for continual learning setups (Chen et al., 2020; Abati et al., 2020), where data arrive in a more realistic sequential manner. InstAParam (Chen et al., 2020), Random Path Selection (RPS) (Rajasegaran et al., 2019), and MoE (Collier et al., 2020) start with multiple parallel blocks at each layer, finding input-specific or task-specific paths within this large network. Nevertheless, this requires knowledge of the number of tasks to be learned ahead of training.

More importantly, initializing a large network might be unnecessary as indicated by the competent performance of dynamic networks with gating mechanisms in i.i.d training. In contrast to this, MNTDP (Veniat et al., 2021), LMC (Ostapenko et al., 2021), and CCGN (Abati et al., 2020) start with a standard architecture and grow units to respond to new data or tasks. Of these, MNTDP and LMC develop task-specific networks where all inputs from the same task elicit the same response and therefore do not show a truly dynamic response to stimuli. CCGN, however, composes convolutional filters dynamically to respond to stimuli, using a task-specific vector for every convolutional filter, and task boundaries to freeze frequently active filters. However, this leads to unrestrained growth and fails in the absence of task-boundaries, which makes it unsuitable for general continual learning.

Therefore, we propose a general continual learning method with dynamic modularity and sparsity (Dynamos) induced through reinforcement learning agents trained with policy gradients.

3 METHODOLOGY

Humans learn continually from inherently non-stationary and sequential observations of the world without catastrophic forgetting, even without supervision about tasks to be performed or the arrival of new tasks, maintaining a bounded memory throughout. This involves, among other things, making multi-scale associations between current and previous observations (Goyal and Bengio, 2020) and responding sparsely and dynamically to stimuli (Graham and Field, 2006). The former concerns consolidation of previous experiences and ensuring that learned experiences evoke a similar response. The latter concern dynamically composing a subset of the specialized neural modules available to respond to stimuli, reusing only the relevant previously learned information. This also avoids erasure of information irrelevant to current stimuli but relevant to previous experiences. We now formulate an approach for dynamic sparse and modular general continual learning that mimics these procedures with DNNs.

3.1 Dynamic, Modular, and Sparse Response to Stimuli

To achieve a dynamic, modular, and sparse response to inputs, we use a DNN F with a policy to compose a subset of the available modules in each layer to respond to the input to that layer. More specifically, we

use a CNN which is incentivized to drop some channels in its activations adaptively using policy gradients (Sutton and Barto, 2018; Williams, 1992).

Let us consider the l^{th} convolutional layer with c_l output channels $\forall l \in \{1, 2, \dots, L\}$, where L is the total number of convolutional layers in the network. The input to the convolutional layer is processed using an agent module with actions $a_l \in \{0, 1\}^{c_l}$ as output, where each action represents the decision to drop (action = 0) or keep (action = 1) the corresponding channel of the output of the convolutional layer. The agent module uses a self-attention network to obtain a channel-wise attention vector v_l of dimension c_l , which is converted into "action probabilities" using a probability layer. The policy for choosing actions is then sampled from a c_l -dimensional Bernoulli distribution;

$$p_l = \sigma(v_l)$$

$$\pi_l(a_l) = \prod_{i=1}^{c_l} p_{l,i}^{a_{l,i}} (1 - p_{l,i})^{(1-a_{l,i})}, \quad (1)$$

where $p_l \in (0, 1)^{c_l}$ is the output of the probability layer σ , and π_l is the policy function. The final output of the convolutional layer is the channel-wise product of the actions with the output of the convolution. This policy formulation is used at each convolutional layer in the CNN, leading to L agents in total. The overall structure of an agent for a convolutional layer is shown in Figure 1.

These agents are rewarded for dropping channels while making accurate predictions through a reward function. For an input to the DNN X applied to classification with label Y :

$$Z, V = F(X), V = [v_1 | v_2, \dots | v_L]$$

$$\hat{Y} = \arg \max Z, \quad (2)$$

where Z refers to the logits. Now, the ratio of activations or channels that were retained in the layer l is determined by $\frac{1}{c_l} \sum_{i=1}^{c_l} a_{l,i}$. So, for a target activation retention rate per layer or "keep ratio" kr , the reward function is as follows:

$$R_l(X, Y) = \begin{cases} -(kr - \frac{1}{c_l} \sum_{i=1}^{c_l} a_{l,i})^2, & \text{if } \hat{Y} = Y \\ -\lambda(kr - \frac{1}{c_l} \sum_{i=1}^{c_l} a_{l,i})^2, & \text{otherwise.} \end{cases} \quad (3)$$

Therefore, when the DNN's predictions are correct, each agent is rewarded for dropping enough activations to match the "keep ratio" from its corresponding convolutional layer. However, when the prediction is incorrect, each agent is penalized for the same, scaled by a constant penalty factor λ . The global nature of the reward function, achieved through dependence on the correctness of the prediction, also enforces coordination between agents. Following REINFORCE (Williams, 1992), the loss from all agents

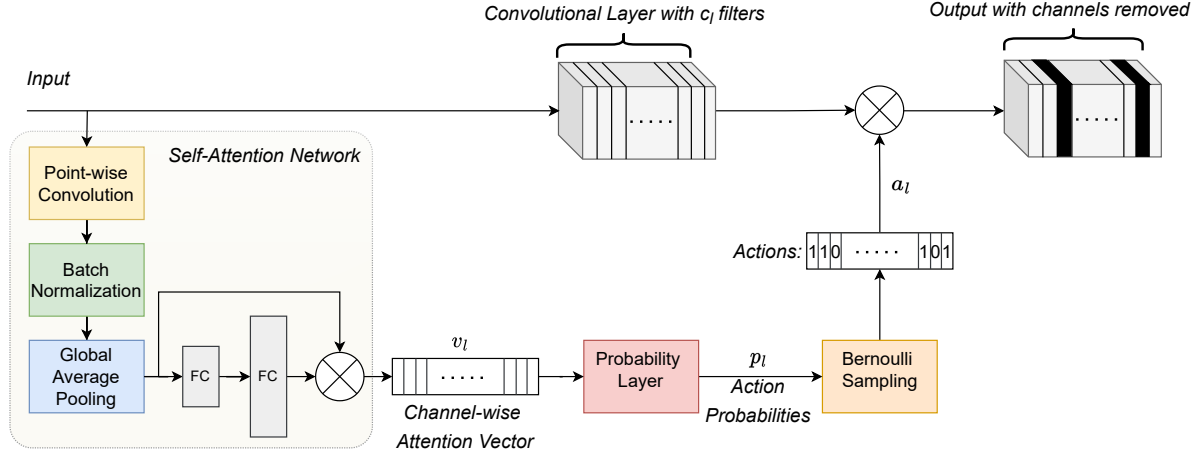


Figure 1: An overview of *Dynamos*' dynamic and sparse response mechanism at the l^{th} convolutional layer. Blacked activations are removed. The agent (bottom path) self-attention network uses a pointwise convolution to match output channels and global average pooling to get a channel-length flattened vector. This is sent through an MLP with one hidden layer and Sigmoid activation, and multiplied with the original channel-length representation to get the channel-wise self-attention vector.

$t = 1, 2, \dots, L$ is:

$$\begin{aligned}
 L_R(X, Y) &= \mathbb{E}_l \mathbb{E}_\pi [-R_l(X, Y) \log \pi_l(a_l)] \\
 &= \mathbb{E}_l \mathbb{E}_\pi [-R_l(X, Y) \log \prod_{i=1}^{c_l} p_{l,i} a_{l,i} \\
 &\quad + (1 - p_{l,i})(1 - a_{l,i})] \quad (4) \\
 &= \mathbb{E}_l \mathbb{E}_\pi [-R_l(X, Y) \sum_{i=1}^{c_l} \log [p_{l,i} a_{l,i} \\
 &\quad + (1 - p_{l,i})(1 - a_{l,i})]].
 \end{aligned}$$

Although the agents along with this loss ensure sparse and dynamic responses from the DNN, they do not explicitly impose any specialization of compositional neural modules seen in humans. As the channel-wise "modules" activated in the DNN are directly dependent on the channel-wise attention vectors, we finally apply a specialization loss that we call prototype loss to them. Concretely, for classification, in any batch of inputs, we pull the vectors belonging to the same class together while pushing those from different classes away. This would cause different subsets of channel-wise modules to be used for inputs of different classes. When combined with a sufficiently high "keep ratio", this will encourage overlap and therefore, reuse of relevant previously learned information (for example, reusing channels corresponding to a learned class for a newly observed class) and, consequently, learning of general-purpose features by the modules. For an input batch X with corresponding labels Y , and the corresponding batch of concatenated channel-wise attention vectors V (Equation 2), the prototype loss is given by:

$$L_P(X, Y) = \frac{1 + \sum_{(V_1, V_2) \in V^2: Y_1 = Y_2} \text{MSE}(V_1, V_2)}{1 + \sum_{(V_1, V_2) \in V^2: Y_1 \neq Y_2} \text{MSE}(V_1, V_2)}, \quad (5)$$

where *MSE* refers to the Mean Squared Error estimator. Note that we only apply this loss to samples for which the predictions were correct.

3.2 Multi-Scale Associations

As discussed earlier, one of the mechanisms employed by humans to mitigate forgetting is multi-scale associations between current and previous experiences.

With this goal in mind, we follow recent rehearsal-based approaches (Buzzega et al., 2020; Riemer et al., 2019) that comply with GCL and use a memory buffer during training to store previously seen examples and responses. The buffer is updated using reservoir sampling (Vitter, 1985), which helps to approximate the distribution of the samples seen so far (Isele and Cosgun, 2018). However, we only consider the subset of batch samples on which the prediction was made correctly for addition to the memory buffer. These buffer samples are replayed through the DNN alongside new samples with losses that associate the current response with the stored previous response, resulting in consistent responses over time.

Let M denote the memory buffer and D_T denote the current task stream, from which we sample batches (X_M, Y_M, Z_M, V_M) and (X_l, Y_l) , respectively. Here, Z_M and V_M are the saved logits and channel-wise attention vectors corresponding to X_M when it was initially observed. The consistency losses associated with current and previous responses are obtained

during the task T as follows:

$$\begin{aligned} Z'_M, V'_M &= F(X_M) \\ L_C(Z_M, Z'_M) &= \mathbb{E}_{X_M} [\|Z_M - Z'_M\|_2^2] \\ L_C(V_M, V'_M) &= \mathbb{E}_{X_M} [\|V_M - V'_M\|_2^2]. \end{aligned} \quad (6)$$

In addition to consistency losses, we also enforce accuracy, and dynamic sparsity and modularity on the memory samples. Therefore, we have four sets of losses:

- Task performance loss on current and memory samples to ensure correctness on current and previous tasks. For classification, we use cross-entropy loss (L_{CE}).
- Reward losses (Equation 4) on current and memory samples to ensure dynamic modularity and sparsity on current and previous tasks.
- Prototype losses (Equation 5) on current and memory samples to ensure the specialization of modules on current and previous tasks.
- Consistency losses (Equation 6) for multi-scale associations between current and previous samples.

Putting everything together, the total loss becomes:

$$\begin{aligned} L_{\text{total}} &= L_{CE}(X_B, Y_B) + \gamma L_r(X_B) \\ &+ \beta [L_{CE}(X_M, Y_M) + \gamma L_r(X_M)] \\ &+ \alpha L_C(Z_M, Z'_M) + \alpha_p L_C(V_M, V'_M) \\ &+ w_p [L_P(X_B, Y_B) + L_P(X_M, Y_M)]. \end{aligned} \quad (7)$$

The weights given to the losses - α , α_p , β , w_p , and γ , and the penalty for misclassification (λ) and keep ratio (kr) in Equation 3, are hyperparameters. Note that we employ a warm-up stage at the beginning of training, where neither the memory buffer nor the agents are employed. This is equivalent to training using only the cross-entropy loss for this period, while the agents are kept frozen. This gives agents a better search space when they start searching for a solution. We call our method as described above *Dynamos*.

4 EXPERIMENT DETAILS

Datasets. We show results on sequential variants of MNIST (LeCun et al., 1998) and SVHN: Seq-MNIST and Seq-SVHN (Netzer et al., 2011), respectively. Seq-MNIST and Seq-SVHN divide their respective datasets into 5 tasks, with 2 classes per task. Furthermore, to test the applicability of *Dynamos* under general continual learning, we also use the MNIST-360 dataset (Buzzega et al., 2020).

Architecture. We use a network based on the ResNet-18 (He et al., 2016) structure by removing the later two of its four blocks and reducing the number of filters per convolutional layer from 64 to 32. The initial convolution is reduced to 3×3 to work with smaller image sizes. For the baseline experiments, we did not use any agents. For our method, while agents can be used for all convolutional layers, we only use agents in the second block. We make this choice based on recent studies that observe that earlier layers undergo minimal forgetting (Davari et al., 2022), are highly transferrable (Yosinski et al., 2014), and are used for most examples even when learned with dynamic modularity (Abati et al., 2020). We use a sigmoid with a temperature layer as the probability layer in the agents and a probability of 0.5 as a threshold for picking actions, i.e., channels during inference. The temperature serves the purpose of tuning the range of outputs of the self-attention layers, ensuring that the probabilities being sampled to choose the actions are not too small and that enough activations are chosen to enable learning. The exact network structure used for each experiment, including the self-attention networks of the agents, can be found in Appendix, in Table 3 and Table 4.

Settings. All methods are implemented in the Mammoth repository¹ in PyTorch 1.6 and were trained on Nvidia V100 GPUs. The hyperparameters corresponding to each experiment can be found in Appendix, Table 5. We always maintain a keep ratio higher than $1/Num_tasks$ to allow the learning of overlapping, reusable, and general-purpose modules. The temperature of the Sigmoid activation of the probability layers is kept at 0.15 unless mentioned otherwise.

5 RESULTS

We will evaluate *Dynamos* under two standard evaluation protocols that adhere to the core desiderata of GCL.

5.1 Class-Incremental Learning (CIL)

Class-incremental learning (CIL) refers to the evaluation protocol in which mutually exclusive sets of classes are presented sequentially to the network, and the identity of the task is not provided at the test time, which meets the core desiderata of GCL (Farquhar and Gal, 2018). We compare against Condi-

¹<https://github.com/aimagelab/mammoth/>

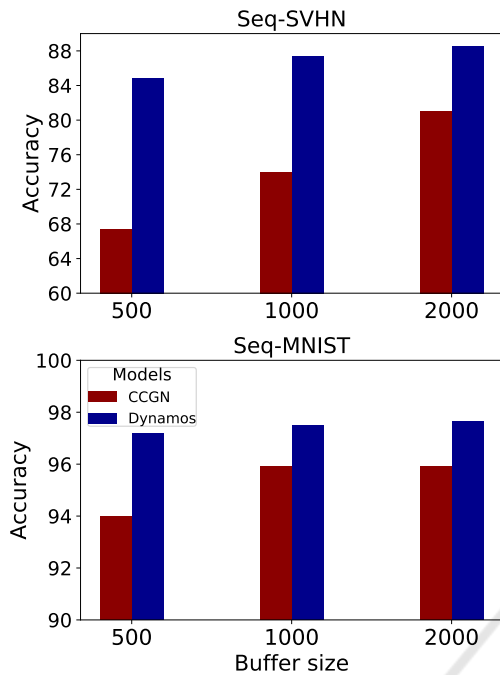


Figure 2: Quantitative results under Class-Incremental Learning protocol. Results are averaged across three seeds. CCGN values taken from the original paper. The precise accuracies can be found in Table 2.

tional Convolutional Gated Network (CCGN) (Abati et al., 2020), which also dynamically composes convolutional filters for continual learning. We observe in Figure 2 that *Dynamos* shows higher accuracies on both the Seq-MNIST and Seq-SVHN datasets under all buffer sizes. However, CCGN requires a separate task vector for every task per convolutional layer, resulting in unrestricted growth during training, whereas we maintain a bounded memory through training. Furthermore, unlike CCGN, we do not leverage the task boundaries or the validation set during training. Therefore, *Dynamos* outperforms the previous state-of-the-art for dynamic compositional continual learning in class-incremental learning, while showing bounded memory consumption during training.

5.2 General Continual Learning (GCL)

So far, we have observed *Dynamos* under the CIL protocol. Unlike CIL, real-world data streams without clear task boundaries, where the same data may reappear under different distributions (e.g. different poses). Following (Buzzega et al., 2020), we approximate this setting using MNIST-360, where tasks overlap in digits (i.e. classes), reappear under different rotations (i.e. distributions), and each example is seen

exactly once during training. This serves as a verification of the adherence to the GCL desiderata (Farquhar and Gal, 2018; Delange et al., 2021). We study the impact of both dynamic modularity as well as multi-scale associations by removing them incrementally from *Dynamos*. When neither is used, the learning is done using vanilla gradient-based training, with no strategy to counter forgetting. When dynamic modularity is removed, the learning strategy forms our baseline, where no agents are used, simplifying the total training loss from Equation 7 to:

$$L_{\text{base}} = L_{CE}(X_B, Y_B) + \beta L_{CE}(X_M, Y_M) + \alpha L_C(Z_M, Z'_M). \quad (8)$$

Table 1 shows that *Dynamos* outperforms the baseline in all buffer sizes, proving that dynamic modularity is advantageous in GCL. Furthermore, when multi-scale associations are also removed, no buffer is used, and the DNN undergoes catastrophic forgetting. Thus, *Dynamos* is applicable to general continual learning, with dynamic modularity improving over the baseline. We hypothesize that dynamic modularity makes dealing with the blurred task boundaries of GCL easier by adaptively reusing relevant previously learned information, which in this case corresponds to learned filters.

6 MODEL CHARACTERISTICS

We now analyze some of the characteristics and advantages of *Dynamos*. For all experiments in this section, we use our model trained on Sequential-MNIST with buffer size 500.

6.1 Dynamic Modularity and Compositionality

Humans show modular and specialized responses to stimuli (Meunier et al., 2010) with dynamic and sparse response to inputs (Graham and Field, 2006) - a capability that we instilled in our DNN while learning a sequence of tasks by dynamically removing channel activations of convolutional layers. Therefore, we examine the task- and class-wise tendencies of the firing rates of each neuron (filter) in Figure 3.

It can be seen that *Dynamos* learns a soft separation of both tasks and classes, as evidenced by the per-task and per-class firing rates, respectively, of each filter. This is in contrast to static methods, where all filters react to all examples. Figure 3a further shows that this allows learning of similar activation patterns for similar examples. For example, MNIST

Table 1: General continual learning results for multiple buffer sizes. All results are averaged across five seeds.

Multi-Scale Associations	Dynamic Modularity	Buffer Size		
		100	200	500
✓	✓	64.418 ± 4.095	79.638 ± 2.853	90.519 ± 0.737
✓	✗	61.192 ± 3.072	75.364 ± 1.259	88.150 ± 0.888
✗	✗	18.712 ± 0.690		

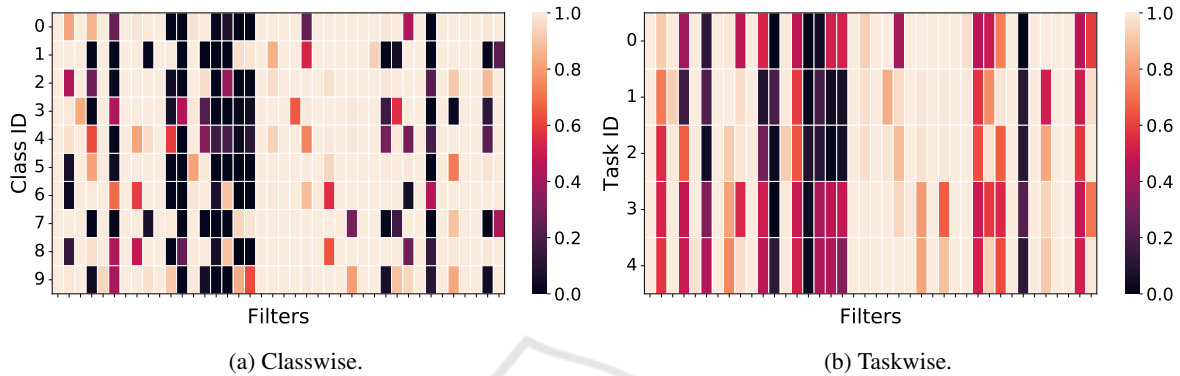


Figure 3: Filter activation rates on the test set for each filter with respect to tasks and classes. For ease of visualization, we only look at the last 40 filters. Full visualizations can be found in Appendix (Figure 7).

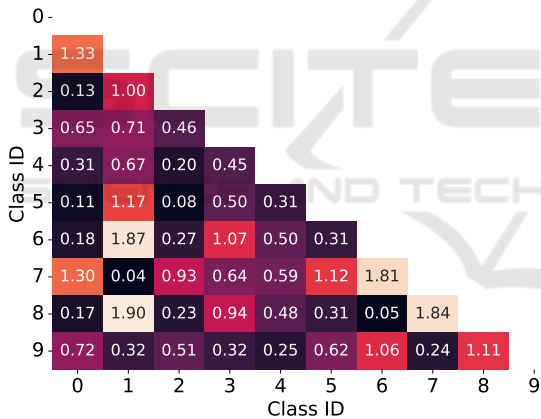


Figure 4: Jensen-Shanon Divergences ($\times 100$) of the activation rates of class pairs on the test set.

digit pairs 1 and 7, and 6 and 8, which share some shape similarities, also share similarities in their activation patterns/rates. This could be attributed to being able to reuse and drop learned filters dynamically, which causes the DNN to react similarly to similar inputs, partitioning its responses based on example similarities. Additionally, the ability to dynamically reuse filters allows DNNs to learn overlapping activation patterns for dissimilar examples and classes, instead of using completely disparate activation patterns. This also facilitates the learning of sequences of tasks without having to grow the DNN capacity or having a larger capacity at initialization, as opposed

to the static parameter isolation methods for continual learning.

Following (Abbasi et al., 2022), we quantify the overlap between the activation rates for each class pair in the final layer using the Jensen-Shanon divergence (JSD) between them in Figure 4. Lower JSDs signify higher overlap. The JSD is lowest for the class pair (1, 7) (both digits look like vertical lines), and is $\sim \frac{1}{15^{\text{th}}}$ the average JSD across class pairs, and $\sim \frac{1}{42^{\text{th}}}$ that of the least overlapping class pair (1, 8) (1 is a line, 8 is formed of loops). Now, as per Equation 1, filters in the layer are activated based on the channel-wise attention vector v_L (see Equation 2), which are pushed together for examples of the same classes, and pushed away from each other for examples of different classes using prototype loss (Equation 5). We visualize the t-SNEs of these v_L s on the test set in Figure 5 and observe that the samples belonging to the same classes are clustered, confirming the effectiveness of our prototype loss. Moreover, the clusters of visually similar classes are close together, which is concomitant with the JSDs and class-wise activation rates seen earlier. Class similarities are also reflected through multiple clusters for the digit 9, indicating its similarity with the digits 6 (loop) and 1 (line) in one cluster, but also with 7 (line) and 4 (line + loop) in another cluster. Finally, we observe that there are examples that are scattered away from their class clusters and overlap with other clusters, probably indicating that these particular examples are visually closer

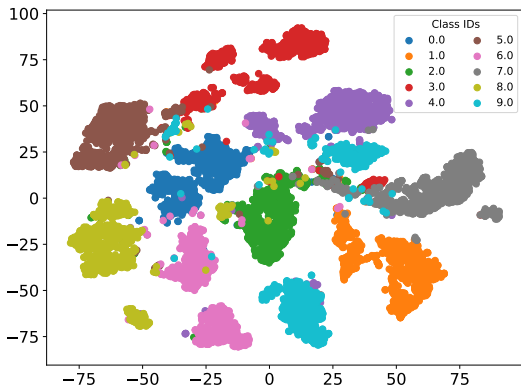


Figure 5: t-SNEs on the test set of class prototypes learned from channel-wise self-attention vectors for all classes.

to other digits. Note, however, that these similar examples and classes are distributed across tasks, which explains the lower similarities in activation patterns between task pairs in Figure 3b compared to the class pairs in Figure 3a.

Therefore, *Dynamos* is capable of learning modular and specialized units that result in input-adaptive dynamic separation and overlap of activations, based on the extent of similarities with previously learned examples. We also contend that the overlapping activations for digits of similar shape suggest the learning of general-purpose features.

6.2 Trial-to-Trial Variability

The brain is known to show variability in response across trials (Faisal et al., 2008; Werner and Mountcastle, 1963). For the same stimulus, the precise neuronal response could differ between trials, a behavior absent in most conventional DNNs. In our method, this aspect of brains can be mimicked by using Bernoulli sampling instead of thresholding to pick keep/drop decisions at each convolutional layer. In Figure 6, we plot the response variability in the last convolutional layer of our DNN with the same example in four trials. We only pick responses for which the predictions were correct. It can be seen that each trial evoked a different response from the DNN. Furthermore, despite the differences, there are also some similarities in the response. There are some filters that are repeatedly left unused, as well as some filters that are used in every trial. This demonstrates that *Dynamos* can additionally simulate the trial-to-trial variability observed in brains.

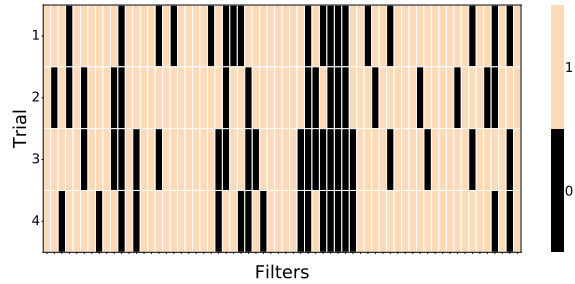


Figure 6: Trial-to-trial variability of responses to same input in *Dynamos*.

7 CONCLUSION AND FUTURE WORK

We propose *Dynamos*, a method for general continual learning, that simulates the dynamically modular and sparse response to stimuli observed in the brain. *Dynamos* rewards the input-adaptive removal of channel activations of convolutional layers using policy gradients for dynamic and sparse responses. To further induce modularity, channel-wise self-attention vectors corresponding to each convolutional layer are pulled together for examples from same classes, and are pushed apart for examples from different classes; these vectors are then used to sample the keep/drop decision for the corresponding channel. Using a memory buffer, we enforce multi-scale consistency between previous and current responses to prevent forgetting. *Dynamos* outperforms previous baselines on multiple datasets when evaluated using class-incremental learning (CIL) and general continual learning (GCL) protocols. *Dynamos* exhibits similar and overlapping responses for similar inputs, yet distinct responses to dissimilar inputs by utilizing subsets of learned filters in an adaptive manner. We quantified the extent of class-wise overlaps and showed that the semantic similarity of classes (digits in MNIST, e.g. 1 and 7) are reflected in higher representation overlaps. We additionally visualized the channel-wise attention vectors and observed that they are clustered by the classes and the clusters of semantically similar classes lie together or overlap. Finally, we also demonstrated the ability of our method to mimic the trial-to-trial variability seen in the brain, where same inputs achieve same outputs through different “responses”, i.e. activations. Thus, we consider our work as a step toward achieving dynamically modular and general-purpose continual learning.

REFERENCES

- Abati, D., Tomczak, J., Blankevoort, T., Calderara, S., Cucchiara, R., and Bejnordi, B. E. (2020). Conditional channel gated networks for task-aware continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3931–3940.
- Abbasi, A., Nooralinejad, P., Braverman, V., Pirsiavash, H., and Kolouri, S. (2022). Sparsity and heterogeneous dropout for continual learning in the null space of neural activations. *arXiv preprint arXiv:2203.06514*.
- Adel, T., Zhao, H., and Turner, R. E. (2020). Continual learning with adaptive weights (CLAW). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Bengio, E., Bacon, P.-L., Pineau, J., and Precup, D. (2015). Conditional computation in neural networks for faster models. *arXiv preprint arXiv:1511.06297*.
- Bolukbasi, T., Wang, J., Dekel, O., and Saligrama, V. (2017). Adaptive neural networks for efficient inference. In *International Conference on Machine Learning*, pages 527–536. PMLR.
- Buzzega, P., Boschini, M., Porrello, A., Abati, D., and Calderara, S. (2020). Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930.
- Chaudhry, A., Ranzato, M., Rohrbach, M., and Elhoseiny, M. (2019). Efficient lifelong learning with A-GEM. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Chen, H.-J., Cheng, A.-C., Juan, D.-C., Wei, W., and Sun, M. (2020). Mitigating forgetting in online continual learning via instance-aware parameterization. *Advances in Neural Information Processing Systems*, 33:17466–17477.
- Collier, M., Kokopoulou, E., Gesmundo, A., and Berent, J. (2020). Routing networks with co-training for continual learning. *arXiv preprint arXiv:2009.04381*.
- Davari, M., Asadi, N., Mudur, S., Aljundi, R., and Belilovsky, E. (2022). Probing representation forgetting in supervised and unsupervised continual learning. *arXiv preprint arXiv:2203.13381*.
- Delange, M., Aljundi, R., Masana, M., Parisot, S., Jia, X., Leonardis, A., Slabaugh, G., and Tuytelaars, T. (2021). A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Faisal, A. A., Selen, L. P., and Wolpert, D. M. (2008). Noise in the nervous system. *Nature reviews neuroscience*, 9(4):292–303.
- Farquhar, S. and Gal, Y. (2018). Towards Robust Evaluations of Continual Learning. *Lifelong Learning: A Reinforcement Learning Approach Workshop at ICML*.
- French, R. M. (1999). Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135.
- Goyal, A. and Bengio, Y. (2020). Inductive biases for deep learning of higher-level cognition. *arXiv preprint arXiv:2011.15091*.
- Graham, D. J. and Field, D. J. (2006). Sparse coding in the neocortex. *Evolution of nervous systems*, 3:181–187.
- Greenwald, N. F., Miller, G., Moen, E., Kong, A., Kagel, A., Dougherty, T., Fullaway, C. C., McIntosh, B. J., Leow, K. X., Schwartz, M. S., et al. (2021). Whole-cell segmentation of tissue images with human-level performance using large-scale data annotation and deep learning. *Nature biotechnology*, pages 1–11.
- Hadsell, R., Rao, D., Rusu, A. A., and Pascanu, R. (2020). Embracing change: Continual learning in deep neural networks. *Trends in cognitive sciences*, 24(12):1028–1040.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Hu, T., Chen, T., Wang, H., and Wang, Z. (2020). Triple wins: Boosting accuracy, robustness and efficiency together by enabling input-adaptive inference. In *ICLR*.
- Isele, D. and Cosgun, A. (2018). Selective experience replay for lifelong learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al. (2017). Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.
- Li, X., Liu, Z., Luo, P., Change Loy, C., and Tang, X. (2017). Not all pixels are equal: Difficulty-aware semantic segmentation via deep layer cascade. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3193–3202.
- McCloskey, M. and Cohen, N. J. (1989). Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier.
- Meunier, D., Lambiotte, R., and Bullmore, E. T. (2010). Modular and hierarchically modular organization of brain networks. *Frontiers in neuroscience*, 4:200.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. (2011). Reading digits in natural images with unsupervised feature learning.
- Nose, Y., Kojima, A., Kawabata, H., and Hironaka, T. (2019). A study on a lane keeping system using cnn for online learning of steering control from real time images. In *2019 34th International Technical Conference on Circuits/Systems, Computers and Communications (ITC-CSCC)*, pages 1–4. IEEE.
- Ostapenko, O., Rodriguez, P., Caccia, M., and Charlin, L. (2021). Continual learning via local module composition. *Advances in Neural Information Processing Systems*, 34:30298–30312.

- Rajasegaran, J., Hayat, M., Khan, S. H., Khan, F. S., and Shao, L. (2019). Random path selection for continual learning. *Advances in Neural Information Processing Systems*, 32.
- Riemer, M., Cases, I., Ajemian, R., Liu, M., Rish, I., Tu, Y., and Tesauro, G. (2019). Learning to learn without forgetting by maximizing transfer and minimizing interference. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Schwarz, J., Czarnecki, W., Luketina, J., Grabska-Barwinska, A., Teh, Y. W., Pascanu, R., and Hadsell, R. (2018). Progress & compress: A scalable framework for continual learning. In *International Conference on Machine Learning*, pages 4528–4537. PMLR.
- Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Taigman, Y., Yang, M., Ranzato, M., and Wolf, L. (2014). Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1701–1708.
- Thrun, S. and Mitchell, T. M. (1995). Lifelong robot learning. *Robotics and autonomous systems*, 15(1-2):25–46.
- Veit, A. and Belongie, S. (2018). Convolutional networks with adaptive inference graphs. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–18.
- Veniat, T., Denoyer, L., and Ranzato, M. (2021). Efficient continual learning with modular networks and task-driven priors. In *International Conference on Learning Representations*.
- Vitter, J. S. (1985). Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*, 11(1):37–57.
- Wang, X., Yu, F., Dou, Z.-Y., Darrell, T., and Gonzalez, J. E. (2018). Skipnet: Learning dynamic routing in convolutional networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 409–424.
- Wang, Z., Zhang, Z., Lee, C.-Y., Zhang, H., Sun, R., Ren, X., Su, G., Perot, V., Dy, J., and Pfister, T. (2022). Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 139–149.
- Werner, G. and Mountcastle, V. B. (1963). The variability of central neural activity in a sensory system, and its implications for the central reflection of sensory events. *Journal of Neurophysiology*, 26(6):958–977.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256.
- Wu, Z., Nagarajan, T., Kumar, A., Rennie, S., Davis, L. S., Grauman, K., and Feris, R. (2018). Blockdrop: Dynamic inference paths in residual networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8817–8826.
- Yosinski, J., Clune, J., Bengio, Y., and Lipson, H. (2014). How transferable are features in deep neural net-

works? *Advances in neural information processing systems*, 27.

- Zenke, F., Poole, B., and Ganguli, S. (2017). Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pages 3987–3995. PMLR.

APPENDIX

Table 2: Class-Incremental learning accuracies for CCGN and Dynamos corresponding to Figure 2.

Buffer Size	Model	Seq-SVHN	Seq-MNIST
500	CCGN	67.45	94.01
	Dynamos	84.815	97.19
1000	CCGN	73.99	95.94
	Dynamos	87.38	97.51
2000	CCGN	81.02	95.94
	Dynamos	88.54	97.57

Table 3: Agent architecture with input features of shape $B \times C_{in} \times H \times W$, output features of shape $B \times C_{out} \times 1 \times 1$, where B is the batch size, C_{in} is the number of channels in the input, and C_{out} is the number of channels expected in the output, which is same as the number of keep/drop actions required for the corresponding convolutional layer.

Operations	Input size	Output size
Pointwise Conv	$B \times C_{in} \times h \times w$	$B \times C_{out} \times h \times w$
Average Pooling	$B \times C_{out} \times h \times w$	$B \times C_{out} \times 1 \times 1$
Reshape	$B \times C_{out} \times 1 \times 1$	$B \times C_{out}$
Linear	$B \times C_{out}$	$B \times C_{out}/16$
ReLU	$B \times C_{out}/16$	$B \times C_{out}/16$
Linear	$B \times C_{out}/16$	$B \times C_{out}$
Reshape	$B \times C_{out}$	$B \times C_{out} \times 1 \times 1$
Sigmoid $_{\tau}$	$B \times C_{out} \times 1 \times 1$	$B \times C_{out} \times 1 \times 1$

Table 4: Architectures used in our experiments. For baseline experiments without dynamic compositionality, we do not use the "Agent" branch. $\text{Conv}(k, n, s, p)$ refers to convolutional layer with kernel size k , number of filters n , stride s , and padding p . BN refers to Batch Normalization. $\text{Linear}(M, N)$ refers to a linear layer with M -dimensional input and N -dimensional output. $\text{Agent}(C_{in}, C_{out}, \tau)$ refers to the Agent subnetwork with C_{in} input channels, C_{out} output channels, and τ temperature of the sigmoid in the probability layer (See Figure 1, Section 4). The elementwise multiplication of the actions from the agents with the output of a convolutional layer is done after the application of batch normalization, if present, but before the ReLU activation function, if present. For complete description of Agent architecture, refer to Table 3. $num_classes$ refers to the number of classes to be predicted.

Component	Main Branch	Residual branch	Agent branch
Conv1	$\text{Conv}(3, 32, 1, 1), \text{BN}, \text{ReLU}$	–	–
Block1	$\begin{bmatrix} \text{Conv}(3, 32, 1, 1), \text{BN}, \text{ReLU} \\ \text{Conv}(3, 32, 1, 1), \text{BN}, \text{ReLU} \end{bmatrix}$	$\begin{bmatrix} \text{Identity} \\ \text{Identity} \end{bmatrix}$	$\begin{bmatrix} \text{Agent}(64, 64, 0.15) \\ \text{Agent}(64, 64, 0.15) \end{bmatrix}$
Block2	$\begin{bmatrix} \text{Conv}(3, 32, 1, 1), \text{BN}, \text{ReLU} \\ \text{Conv}(3, 32, 1, 1), \text{BN}, \text{ReLU} \end{bmatrix}$	$\begin{bmatrix} \text{Conv}(1, 64, 2, 0), \text{BN} \\ \text{Conv}(1, 64, 2, 0), \text{BN} \end{bmatrix}$	$\begin{bmatrix} \text{Agent}(64, 64, 0.15) \\ \text{Agent}(64, 64, 0.15) \end{bmatrix}$
Classifier	$\text{Linear}(64, num_classes)$	–	–

Table 5: Hyperparameters for all the datasets for Dynamos.

Dataset	Buffer Size	lr	#Epochs	#Warmup Ep./Itr.	Batch Size	Memory Batch Size	α	β	α_p	λ	γ	w_p	kr
Seq-MNIST	500	0.07	1	10 it	10	10	0.2	2.0	0.2	500	0.5	0.3	0.7
	1000	0.07	1	10 it	10	10	0.1	2.5	0.2	200	0.7	0.5	0.7
	2000	0.07	1	10 it	10	10	0.5	3.0	0.2	200	0.5	0.5	0.7
SVHN	500	0.07	70	10 ep	16	16	2.0	3.0	1.0	500	1.0	0.5	0.7
	1000	0.07	70	10 ep	16	16	2.5	2.0	0.2	500	1.0	0.5	0.7
	2000	0.07	70	10 ep	16	16	2.5	2.0	0.2	500	1.0	0.5	0.7
MNIST-360	100	0.07	1	10 it	16	16	0.2	1.0	0.1	200	0.5	0.5	0.7
	200	0.07	1	10 it	16	16	0.2	1.5	0.1	200	1.0	0.5	0.7
	500	0.07	1	10 it	16	16	0.1	1.5	0.1	200	0.3	0.3	0.7

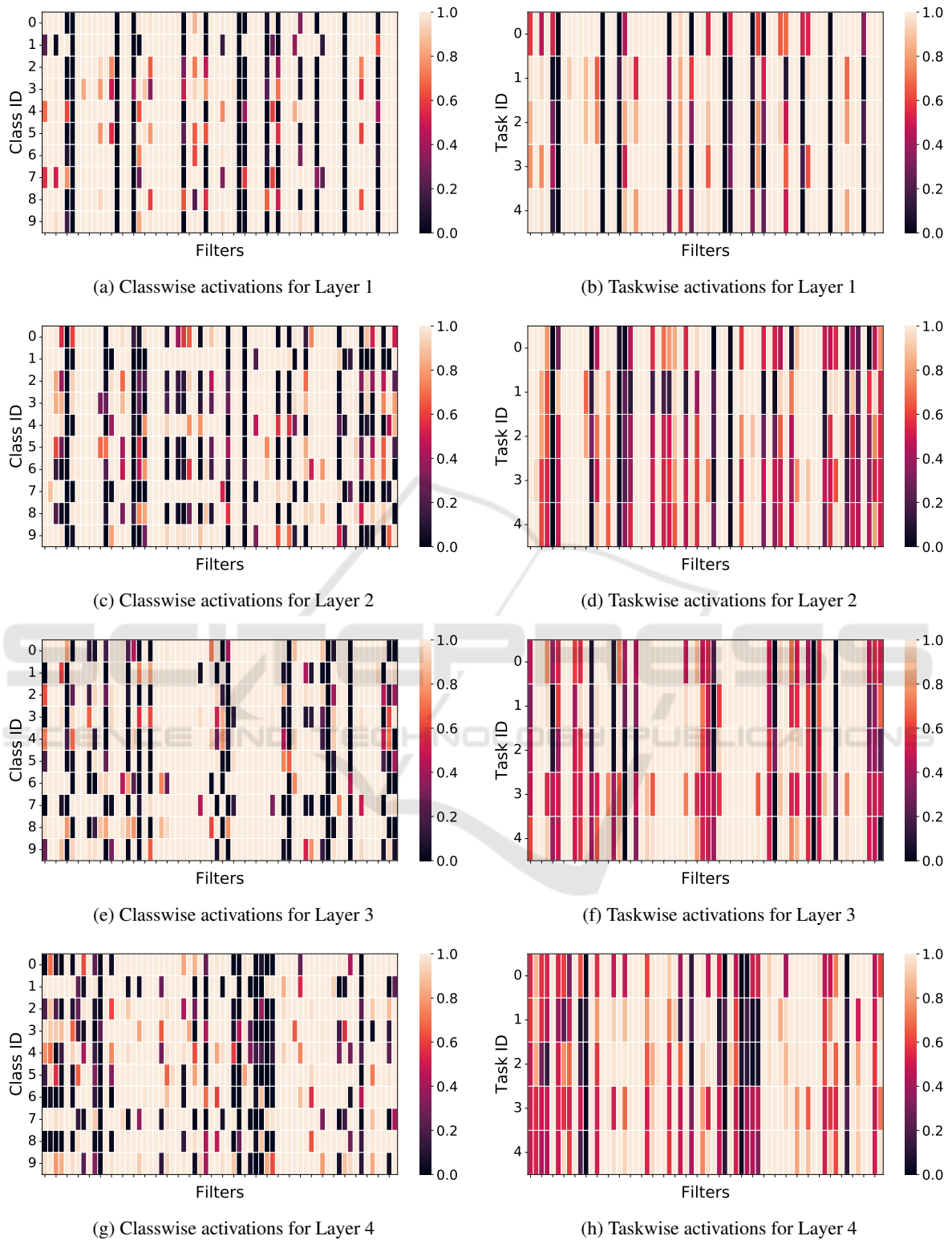


Figure 7: Filter activation rates for each filter in each convolutional layer of Block 2 with respect to MNIST tasks and classes. Overlapping activations of tasks and classes indicative of similarities between them can still be observed. For e.g. 1 and 7 still show very similar responses.