# Continual learning using dendritic modulations on view-invariant feedforward weights
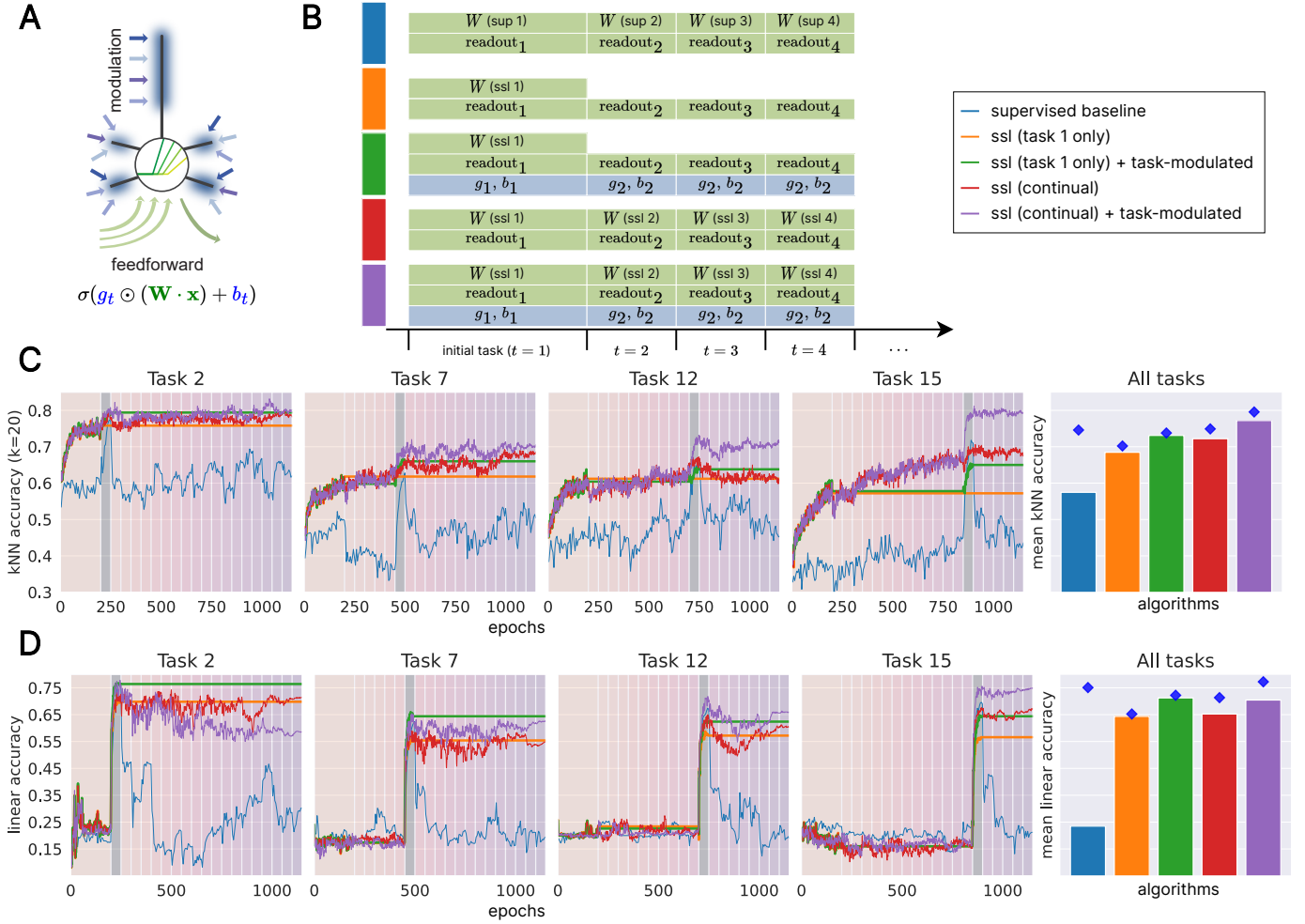
## Summary

The brain is remarkably adept at learning from a continuous stream of data without significantly forgetting previously learnt skills. Conventional machine learning models struggle at continual learning, as weight updates that optimize the current task interfere with previously learnt tasks. A simple remedy to catastrophic forgetting is freezing a network pretrained on a set of base tasks, and training task-specific readouts on this shared trunk. However, this assumes that representations in the frozen network are separable under new tasks, therefore leading to sub-par performance. To continually learn on novel task data, previous methods suggest weight consolidation – preserving weights that are most impactful for the performance of previous tasks – and memory-based approaches – where the network is allowed to see a subset of images from previous tasks.

For biological networks, prior work showed that dendritic top-down modulations provide a powerful mechanism to learn novel tasks while initial feedforward weights solely extract generic view-invariant features. Therefore, we propose a continual learner that optimizes the feedforward weights towards view-invariant representations while training task-specific modulations towards separable class clusters. In a task-incremental setting, we train feedforward weights using a self-supervised algorithm, while training the task-specific modulations and readouts in a supervised fashion, both exclusively through current-task data. We show that this simple approach avoids catastrophic forgetting of class clusters, as opposed to training the whole network in a supervised manner, while also outperforming (a) task-specific readout without modulations and (b) frozen feedforward weights. This suggests that (a) top-down modulations are necessary and sufficient to shift the representations towards separable clusters and that (b) the SSL objective learns novel features based on the newly presented objects while maintaining features relevant to previous tasks, without requiring specific synaptic consolidation mechanisms.

## Additional details

A common theory of biological visual processing is that the initial feedforward pass generates representations that encode task-independent features, while complex visual tasks require top-down inputs (Serre et al., 10.1073/pnas.0700622104). Indeed, with suitable feed-forward weights, biophysical simulations of L5PC neurons show that task-specific modulations are sufficient to learn different tasks (Fig. A, Wybo et al., 10.1073/pnas.2300558120). Therefore, it is biologically implausible to optimize the feedforward weights to solve classification tasks directly, as such objectives encourage visual representations to be invariant to task-irrelevant features, which might turn out to be relevant for other tasks. Instead, biology suggests that view-invariance is a more suitable objective for the feedforward weights. In fact, it has been shown that rats perform visual object recognition tasks well despite changes to size, view and lighting (Zoccolan et al., 10.1073/pnas.0811583106). For newborn chicks in a virtual environment, view-invariant representations emerged only if reared in environments with smoothly moving visual objects, whereas presenting temporally non-smooth objects led to representations that are sensitive to viewpoint changes (Wood and Wood, 10.1111/cogs.12595). We hence focus on self-supervised learning (SSL) algorithms for visual representations that train a network to map distorted (e.g. random shifts or color jitter) views of the same image to the same vector representation. Whereas current continual learning approaches make use of weight consolidation (e.g. Kirkpatrick et al., 10.1073/pnas.1611835114) or replay (Rebuffi et al., 10.1109/CVPR.2017.587), we show that we can avoid the former by training the feedforward weights on an objective that produces representations useful for all view-invariant tasks, while the latter is only necessary to restore the readout weights.

We use a standard continual learning benchmark, where CIFAR-100 is split into 20 balanced tasks with 5 classes each. The first task is trained for 200 epochs, the following tasks are trained for 50 epochs each. At any point during training, only the image-label pairs corresponding to the classes of the current task can be used. We use the "tiny" variant of the vision transformer architecture (Dosovitsky et al., 10.48550/arXiv.2010.11929) with a patch size of 4, $d_{\text{model}} = 192$ and $d_{\text{ff}} = 768$, otherwise following

the ViT-Base configuration. All parameters of the original architecture are considered feedforward (referred to as $\mathbf{W}$ in the figure). The task-specific modulations are implemented as additional affine transformations (gain and bias, Fig. A) to the attentional query, key and value projections as well as to the first layer of each MLP. We implement supervised learning as cross-entropy loss minimization on task-specific readouts, and we use Barlow Twins (Zbontar et al., 10.48550/arXiv.2103.03230) to implement SSL. We describe 5 algorithms (Fig. B): the supervised baseline (blue) trains the whole network and the task-specific readout on all data available for the respective task. Then, we consider a pretraining-esque algorithm (orange, green) – training the feedforward weights with SSL only during task 1 and freezing the weights afterwards – and a continual algorithm (red, purple) – applying the SSL objective during all tasks on the available data. For both cases, we assess how task-modulations (green and purple) affect task performance compared to the readout-only approach (orange and red). Using k-nearest neighbor (kNN, $k = 20$) evaluation, we find that continual SSL outperforms its frozen counterparts, while not suffering from catastrophic forgetting as is the case with the supervised baseline (Fig. C). Furthermore, we find that performance increases with later tasks, suggesting that the continual algorithm learns to extract novel features. Averaged over all tasks, combining continual SSL with task-modulations yields the highest accuracy with minor forgetting (blue diamond: best performance obtained throughout training, averaged over all tasks). Forgetting occurs more prominently with linear readout (Fig. D), and the continual approach performs comparably to the pretraining approach. Still, the high kNN accuracy of the continual SSL with task-modulations suggests that the obtained representations form class-specific clusters. These representational drifts are also present in biological networks, and biologically plausible replay-based learning rules have been proposed to restore linear readout accuracy (Rule and O'Leary, 10.1073/pnas.2106692119).

In conclusion, continual learning emerges from training task-specific modulations on top of view-invariant feedforward weights, mitigating catastrophic forgetting while improving performance over time.