


PHI-NETS: BRAIN-INSPIRED NON-CONTRASTIVE LEARNING BASED ON TEMPORAL PREDICTION HYPOTHESIS

Satoki Ishikawa^{1,2,*}, Makoto Yamada^{2,*,†}, Han Bao^{3,2}, Yuki Takezawa^{2,3}

¹Science Tokyo, ²Okinawa Institute of Science and Technology, ³Kyoto University

ABSTRACT

Predictive coding has been established as a promising neuroscientific theory to describe the mechanism of information processing in the retina or cortex. This theory hypothesises that cortex predicts sensory inputs at various levels of abstraction to minimise prediction errors. Inspired by predictive coding, [Chen et al. \(2024\)](#) proposed another theory, *temporal prediction hypothesis*, to claim that sequence memory residing in hippocampus has emerged through predicting input signals from the past sensory inputs. Specifically, they supposed that the CA3 predictor in hippocampus creates synaptic delay between input signals, which is compensated by the following CA1 predictor. Though recorded neural activities were replicated based on the temporal prediction hypothesis, its validity has not been fully explored. In this work, we aim to explore the temporal prediction hypothesis from the perspective of self-supervised learning (SSL). Specifically, we focus on non-contrastive learning, which generates two augmented views of an input image and predicts one from another. Non-contrastive learning is intimately related to the temporal prediction hypothesis because the synaptic delay is implicitly created by StopGradient. Building upon a popular non-contrastive learner, SimSiam, we propose *PhiNet*, an extension of SimSiam with two predictors explicitly corresponding to the CA3 and CA1, respectively. Through studying the PhiNet model, we discover two findings. First, meaningful data representations emerge in PhiNet more stably than in SimSiam. This is initially supported by our learning dynamics analysis: PhiNet is more robust to the representational collapse. Second, PhiNet adapts more quickly to newly incoming patterns in online and continual learning scenarios. For practitioners, we additionally propose an extension called X-PhiNet integrated with a momentum encoder, excelling in continual learning. All in all, our work reveals that the temporal prediction hypothesis is a reasonable model in terms of the robustness and adaptivity.

 <https://github.com/riverstone496/PhiNets>

1 INTRODUCTION

How does learning and adaptivity emerge in a biological system? It has been a long-standing question in both neuroscience and machine learning. In the neuroscience community, predictive coding has been a promising hypothesis to support the flexibility of biological brains. While predictive coding was initially proposed to explain cortical functions, [Chen et al. \(2024\)](#) recently extended predictive coding to propose the *temporal prediction hypothesis*, which claims that the hippocampus predicts future sensory inputs based on past experiences ([Mumford, 1992](#); [Rao and Ballard, 1999](#); [Friston, 2005](#)). Specifically, [Chen et al. \(2024, Figure 1\)](#) modelled the hippocampus with the CA3 predictor followed by the CA1 predictor—while the former yields synaptic delay to input signals, the latter compensates the time difference between the past and incoming signals by temporal prediction. This is the first attempt to explain the mechanism of sequence (short-term) memory from the viewpoint of temporal prediction. While they tested the temporal prediction hypothesis by using recorded neural activities, the validity of the hypothesis has not been explored sufficiently.

*Equal contribution

†Corresponding author, e-mail: makoto.yamada@oist.jp

This work is aimed at exploring a learning model built upon the temporal prediction hypothesis to see when the hypothesis is reasonable. To this end, we shed light on self-supervised learning (SSL). SSL is a paradigm to train a learner from input sensory patterns without supervised signals, which aligns to biological learning more closely. Over the past decade, machine learning researchers have developed a number of SSL models. Popular SSL models are SimCLR (Chen et al., 2020a) and MoCo (Chen et al., 2020b), which are contrastive learning methods that learn data representations with two augmented views generated from an input image by minimizing the InfoNCE loss (van den Oord et al., 2018), requiring a tremendous number of negative samples to stably obtain representations. Thus, we focus on another SSL model, non-contrastive learning, which learns data representations from only the two augmented views without requiring negative samples. Specifically, SimSiam (Chen and He, 2021) is a natural model to study the temporal prediction hypothesis—SimSiam predicts one augmented view of an input image from another view, introducing an implicit time difference through the StopGradient operation. For this reason, we choose SimSiam, unlike the other non-contrastive models such as Barlow Twins (Zbontar et al., 2021). This implicit connection between SimSiam and the temporal prediction hypothesis is an appealing test bed to computationally verify how memory-based prediction processes in the brain behave in different scenarios.

Building upon SimSiam, we propose the brain-inspired SSL model *PhiNet* for investigating the effectivity of the temporal prediction hypothesis in the context of machine learning. PhiNet extends SimSiam by incorporating an additional predictor after the original predictor. We associate the original and additional predictors with the CA3 and CA1 regions in the hippocampal model (see Figure 1 in Chen et al. (2024) and Figure 1b). We leverage PhiNet as a computational model to implement the temporal prediction hypothesis and study when it effectively learns sensory inputs.

Our first discovery is that PhiNet is less prone to the representational collapse, which leads to stable learning. Non-contrastive learning intrinsically faces the challenge of collapsing into a trivial representation because it eliminates explicit negative signals. In Section 4, we theoretically analyse the learning dynamics of PhiNet to reveal that PhiNet is less sensitive to initialization and the weight decay hyperparameter, and has a wider retraction basin to a non-trivial representation (in C1), relative to SimSiam. This supports the empirically better linear probing performance and hyperparameter robustness of PhiNet across different image datasets. Our second discovery is that PhiNet empirically performs better in online and continual learning, in particular. We tested PhiNet and baseline non-contrastive learners by using the CIFAR-5m dataset (Nakkiran et al., 2021), exposing learners to a gigantic amount of input images but with only significantly fewer epochs. In this scenario, effective memory functions are necessary to lead the learning to success. As a result, PhiNet exhibits better accuracy with less forgetting than SimSiam. Therefore, the effectiveness of the temporal prediction hypothesis is witnessed from the perspective of the robustness and adaptivity.

For practically better performance, we extend PhiNet to additionally propose *X-PhiNet*, which incorporates a momentum encoder, inspired by the Complementary Learning Systems (CLS) theory (McClelland et al., 1995). This extra momentum encoder represent long-term memory in the neocortex, storing information derived from the hippocampal model. X-PhiNet maintains good performances especially in online and continual learning scenarios.

Contributions.

- Section 3: we propose a new non-contrastive learning model called PhiNet, which is inspired by a hippocampal model (Chen et al., 2024).
- Section 4: we compare the learning dynamics (Tian et al., 2021) of PhiNet and SimSiam. Consequently, it elucidates that PhiNet can avoid the complete collapse of representations (Liu et al., 2023; Bao, 2023) more easily than SimSiam with the aid of the additional predictor.
- Section 5.1: we investigate the image classification performance of PhiNet using CIFAR and ImageNet datasets. We show that PhiNet performs comparably to SimSiam but is more robust against weight decay.
- Section 5.2: we further extend PhiNet by proposing X-PhiNet to integrate the neocortex model based on the Complementary Learning Systems (CLS) theory (McClelland et al., 1995). Experimentally, X-PhiNet works effectively in online and continual learning.

Limitations One major limitation of our approach is the use of backpropagation, which differs from the mechanisms in biological neural networks. Our long-term goal is to eliminate backpropagation

to better imitate brain function, but this work focuses on structural aspects of network architectures. Currently, backpropagation-free predictive coding mechanisms for complex architectures like ResNet are in the early stages of development, with most research being limited to simple CNNs. Future research should explore if the proposed structure can enable effective learning with backpropagation-free predictive coding. Another key difference between PhiNet and brains is the presence of recurrent structures. Making the data into time series data and adding a recurrent structure to the model remains as future work.

2 RELATED WORK

Brain-inspired methods. Predictive coding, initially introduced as a theory of the retina (Srinivasan et al., 1982), has gained attention as a unifying theory of cortical functions (Mumford, 1992; Rao and Ballard, 1999; Friston, 2005). They suggest that brains operate by predicting sensory inputs at various levels of abstraction to minimise prediction errors. Recent studies have leveraged these ideas for contrastive learning (van den Oord et al., 2018; Henaff, 2020). Chen et al. (2024) extended the predictive coding theory to the hippocampus with the temporal prediction hypothesis. Specifically, the temporal prediction hypothesis supposes that prediction errors are calculated with the CA1 model and used to update the CA3 model. Some studies have attempted to apply the hippocampal model to representation learning (Pham et al., 2021; 2023). Among them, DualNet refines representation learning based on CLS theory (McClelland et al., 1995; Kumaran et al., 2016), which supposes that the interplay between slow (self-supervised) and fast (supervised) architectures is the basis of brain learning. Pham et al. (2021) examined supervised learning tasks alongside self-supervised training.

Self-supervised learning. Current mainstream approaches to self-supervised learning (SSL) often rely on cross-view prediction frameworks (Becker and Hinton, 1992), with contrastive learning emerging as a prominent SSL paradigm. In contrastive learning like SimCLR (Chen et al., 2020a), a network contrasts positive (similar) and negative (dissimilar) samples to learn data representations. One limitation of SimCLR is its empirical reliance on gigantic negative samples. Theoretically, contrastive learning essentially requires huge number of negative samples (Bao et al., 2022; Awasthi et al., 2022). To address this issue, recent research has focused on approaches free from negative sampling (Grill et al., 2020; Caron et al., 2020; 2021). For instance, BYOL (Grill et al., 2020) trains representations by aligning online and target networks, where the target network is created by maintaining a moving average of the online network parameters. SimSiam (Chen and He, 2021) utilises a Siamese network to align two augmented views of an input by fixing one of the networks with StopGradient. While the lack of negative samples may easily yield collapsed representations, namely, constant representations, Tian et al. (2021) analysed the BYOL/SimSiam dynamics with a two-layer network and found that complete collapse is prevented unless weight decay is excessively strong. We partially leverage their analysis framework to explain the mechanism of our PhiNet. In recent years, many studies have leveraged SimSiam for continual learning (Smith et al., 2021; Madaan et al., 2022) and reinforcement learning (Tang et al., 2023). RM-SimSiam (Fu et al., 2024) and CaSSLe (Fini et al., 2022) enhance the performance of continual learning by incorporating a memory block into SimSiam, while its architecture has not been neuroscientifically grounded. In video self-supervised learning, built on predictive coding principles, a common strategy is to train a predictor that takes a frame or clip from one time-step to generate a distinct representation at a different time-step (Han et al., 2019; 2020; Tan et al., 2023; Bardes et al., 2024).

Note that our aim is to bridge the temporal prediction hypothesis and self-supervised learning. To this end, *non-contrastive* learning provides a better model because both hippocampus and neocortex do not have any mechanism corresponding to negative sample generation. Specifically, SimSiam is a simple yet powerful learning model, and we can benefit from its StopGradient to effectively draw a connection to predictive coding. Thus, we focus on SimSiam as a backbone model in this work.

3 PHINETs (Φ -NETS)

In this paper, we propose PhiNets, which are non-contrastive methods based on CLS theory (McClelland et al., 1995) and the temporal prediction hypothesis (Chen et al., 2024). Chen et al. (2024, Figure 1) provides a hippocampal model, where the entorhinal cortex (EC) serves as an input signal layer, the CA3 region serves as the predictor, and the CA1 region measures the prediction error.

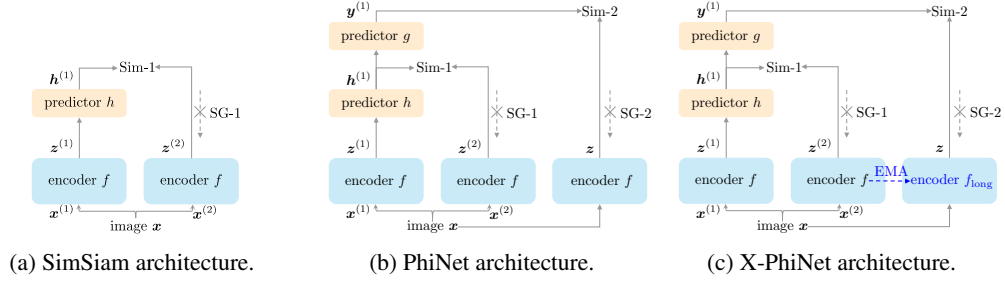


Figure 1: The architecture of SimSiam (Chen and He, 2021) and PhiNets. EMA in the X-PhiNet model stands for the exponential moving average. The architecture originates from a single input, branches out into three paths, and then compares the similarity of all paths in Sim-2. Thus, we call it PhiNet (Φ -Net) because the shape of the architecture resembles the Greek letter Phi (Φ).

The CA3 region receives an input signal from the EC and recurrently forecasts future signals. The prediction output of CA3 is propagated to the CA1 region, which computes the discrepancy between the CA3 prediction and the EC input and refines the internal model stored in CA3. Compensating for the time differences between EC-CA3 and EC-CA1 is hypothesised to facilitate the learning and replay of time sequences in the hippocampus.

Whereas Chen et al. (2024) tested this model to replicate recorded neural activities through simulation, we develop a self-supervised learner PhiNet based on this hypothesis as follows:

- We use deep encoders f and/or f_{long} to represent cortex. See Section 3.2 for more details.
- We model CA3 by a predictor network.
- We model CA1 by combining a loss function and another predictor.
- We train the model by jointly minimizing the loss for the hippocampus and the neocortex models.
- The long-term memory is implemented by an exponential moving average.

Figures 1a and 1b depict the architecture of SimSiam and PhiNet, respectively. Figure 2 illustrates how PhiNet can be interpreted as a hippocampal model (Chen et al., 2024) under the temporal prediction hypothesis.

Note that our approach diverges from the temporal prediction hypothesis method proposed by Chen et al. (2024). Specifically, while they assume an image sequence as input, we consider an original input image and two augmented images as input and feedback signals with time difference (thanks to the StopGradient operation), expanding the applicability of hippocampal models to standard vision tasks.

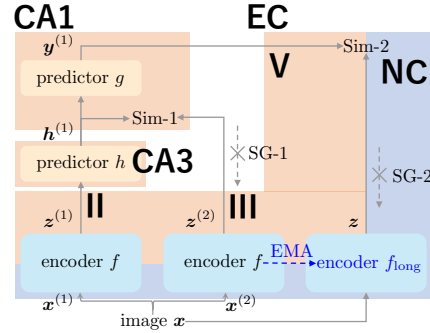


Figure 2: The interpretation as a hippocampal model. NC stands for NeoCortex.

3.1 FAST LEARNING BASED ON TEMPORAL PREDICTION HYPOTHESIS

We provide detailed implementation of the hippocampal model, which serves as a fast learner. The model consists of EC, CA3, and CA1, and we describe each of them below.

Modelling of EC layer. The entorhinal cortex (EC) is the main input and output cortex of the hippocampus (Chen et al., 2024). Let us denote the original input as $x_i \in \mathbb{R}^d$. We model that the hippocampal model has two augmented signals from the original input as $x_i^{(1)} \in \mathbb{R}^d$ and $x_i^{(2)} \in \mathbb{R}^d$, in addition to the original input x_i . Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ denote the encoder. Then, the cortical representation in the EC is given as follows:

$$z_i^{(1)} = f(x_i^{(1)}), \quad z_i^{(2)} = f(x_i^{(2)}), \quad \text{and} \quad z_i = f(x_i).$$

We regard each corresponding to the layers II, III, and V of the EC in [Chen et al. \(2024, Figure 1\)](#). For self-supervised training, we have the triplet dataset $\mathcal{D} = \{(\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, \mathbf{x}_i)\}_{i=1}^n$. The hippocampal learning can be characterized as a learning problem of the encoder f from the training dataset \mathcal{D} .

Modelling of CA3 region. The CA3 region is responsible for predicting future signals:

$$\mathbf{h}_i^{(1)} = h(\mathbf{z}_i^{(1)}), \quad \mathbf{h}_i^{(2)} = h(\mathbf{z}_i^{(2)}),$$

where $h : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is the predictor network. We implement the predictor with a two-layer neural network with the ReLU activation and batch normalization.

Modelling of CA1 region. CA1 measures the difference between the predicted signal and its future signal. In this paper, we model CA1 by a mixture of a loss function and a predictor, while [Chen et al. \(2024\)](#) uses only the MSE loss for modelling CA1. For the Sim-1 of CA1 depicted in Figure 2, we use the symmetric negative cosine loss function to measure the temporally distant signal $\mathbf{z}_i^{(2)}$ (layer III of EC) and the predicted representation from CA3 $\mathbf{h}_i^{(1)} = h(\mathbf{z}_i^{(1)})$:

$$L_{\text{Cos}}(\theta) = -\frac{1}{2n} \sum_{i=1}^n \frac{(\mathbf{h}_i^{(1)})^\top \tilde{\mathbf{z}}_i^{(2)}}{\|\mathbf{h}_i^{(1)}\|_2 \|\tilde{\mathbf{z}}_i^{(2)}\|_2} - \frac{1}{2n} \sum_{i=1}^n \frac{(\tilde{\mathbf{z}}_i^{(1)})^\top \mathbf{h}_i^{(2)}}{\|\tilde{\mathbf{z}}_i^{(1)}\|_2 \|\mathbf{h}_i^{(2)}\|_2},$$

where θ represent the entire model parameter and $\tilde{\mathbf{z}}_i^{(1)} := \text{SG}(\mathbf{z}_i^{(1)}) \in \mathbb{R}^m$ is a latent variable with StopGradient, in which the gradient update shall not be executed.

Remark that StopGradient yields a “time difference” during backpropagation, for which we can interpret PhiNet as a hippocampal model. Let us look closely at Sim-1 in Figure 1b. We let f_t denote the encoder f at the t -th gradient update. The left path of Sim-1 can then be expressed as $\mathbf{h}^{(1)} = h(f_t(\mathbf{x}))$. As the right path of Sim-1 is adopted with StopGradient, it can be written as $\mathbf{z}^{(2)} = \text{SG}(f_t(\mathbf{x})) = f_{t-1}(\mathbf{x})$. Eventually, Sim-1 aligns $f_t(\mathbf{x})$ and $f_{t-1}(\mathbf{x})$ by the predictor h . This Sim-1 interpretation indicates that PhiNet predicts *past* signals, which slightly deviates from the original temporal prediction hypothesis ([Chen et al., 2024](#)) supposing that CA3 is in charge of predicting *future* signals.

In addition to measuring the difference, the CA1 region outputs the signal to the EC (V layer). Thus, we model the output of CA1 as follows:

$$\mathbf{y}_i^{(1)} = g(\mathbf{h}_i^{(1)}), \quad \mathbf{y}_i^{(2)} = g(\mathbf{h}_i^{(2)}),$$

where $g : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is another predictor network. As we will see soon, CLS theory supposes that this feedback from CA1 to the EC eventually propagates to the neocortex (NC), which is stored in long-term memory.

3.2 INCORPORATING SLOW LEARNING MECHANISM

The hippocampus and neocortex play crucial roles in brain cognition. For effective long-term memory storage, it is essential to transfer information from the hippocampus to the NC. We first aim to formulate the joint learning of the hippocampus and NC models. Then, we propose using the exponential moving average (EMA) to transfer model parameters from short-term to long-term memory, with the goal of compressing the original input signal.

In the EC layer, we model the update of the encoder function by using the output of CA1 and the representation $\mathbf{z}_i = f(\mathbf{x}_i)$ (V layer of EC). Then, the loss function can be given as follows:

$$L_{\text{NC}}(\theta) = \frac{1}{2n} \sum_{i=1}^n \|\mathbf{y}_i^{(1)} - \text{SG}(\mathbf{z}_i)\|_2^2 + \frac{1}{2n} \sum_{i=1}^n \|\mathbf{y}_i^{(2)} - \text{SG}(\mathbf{z}_i)\|_2^2.$$

This corresponds to Sim-2 in Figure 2 and is regarded as slow learning. Finally, the whole objective function of PhiNet is given as

$$L(\theta) = \underbrace{L_{\text{Cos}}(\theta)}_{\text{Hippocampus loss / Sim-1}} + \underbrace{L_{\text{NC}}(\theta)}_{\text{Neocortex loss / Sim-2}}.$$

We then minimise $L(\theta)$ to learn the hippocampus and the NC models. The optimisation can be efficiently performed using backpropagation. It is worth noting that we can utilise different loss functions for Sim-1 and/or Sim-2 in PhiNet. In this paper, we set Sim-1 to negative cosine similarity and Sim-2 to either MSE or the negative cosine similarity.

X-PhiNet: Slow learning via stable encoder. The original PhiNet formulation employs the same encoder for both the short-term and long-term memories for simplicity (Section 3.2). To further enhance slow learning, the input representation in EC-V z_i should maintain long-term signals. Thus, we introduce the following stable encoder for long-term memory:

$$z_i = f_{\text{long}}(x_i).$$

Then, we solve the PhiNet optimisation problem by minimising both $L_{\text{NC}}(\theta)$ and $L_{\text{Cos}}(\theta)$ using the exponential moving average (EMA) of the model parameters of f and f_{long} as

$$\xi_{\text{long}} \leftarrow \beta \xi_{\text{long}} + (1 - \beta) \xi,$$

where ξ and ξ_{long} are the model parameters of f and f_{long} , respectively, and $\beta \in [0, 1]$ is a hyperparameter. Model parameters persist in f_{long} more stably than the original encoder f , which facilitates slow learning. We call this method as X-PhiNet.

4 WHAT WE BENEFIT FROM ADDITIONAL CA1 PREDICTOR: LEARNING DYNAMICS PERSPECTIVE

When PhiNet is compared with SimSiam, the additional predictor g in CA1 is peculiar. We study the learning dynamics of PhiNet with a toy model. Despite its simplicity, dynamics analysis is beneficial in showcasing how the predictor g effectively prevents complete collapse.

Analysis model. Let us specify the analysis model, following Tian et al. (2021). The d -dimensional input is sampled from the isotropic normal $x \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and augmented by the isotropic normal $x^{(1)}, x^{(2)} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})$, where σ^2 indicates the strength of data augmentation. The encoder f and predictors g and h are modelled by linear networks without bias: $f(x) := \mathbf{W}_f x$, $g(h) := \mathbf{W}_g h$, and $h(z) := \mathbf{W}_h z$, where $\mathbf{W}_f \in \mathbb{R}^{m \times d}$ and $\mathbf{W}_g, \mathbf{W}_h \in \mathbb{R}^{m \times m}$. The predictors h and g transform latents $z^{(1)}, h^{(1)} \in \mathbb{R}^m$ into $h^{(1)}, y^{(1)} \in \mathbb{R}^m$ with the same dimension m to predict the other noisy latent $z^{(2)}$ and the noise-free latent z , respectively (see Figure 1b).

Unlike L_{Cos} introduced in Section 3, we focus on the (not symmetrised) MSE loss for measuring the discrepancy between $h^{(1)}$ and $z^{(2)}$ for the transparency of analysis. Interested readers may refer to Halvagal et al. (2023) and Bao (2023) for further extension to incorporate the cosine loss into the SimSiam dynamics. Consequently, the expected loss function of PhiNet $\bar{L}(\mathbf{W}_f, \mathbf{W}_g, \mathbf{W}_h)$ is given as follows:

$$\bar{L} := \frac{1}{2} \mathbb{E}_{x, x^{(1)}, x^{(2)} | x} \left[\|\mathbf{W}_h \mathbf{W}_f x^{(1)} - \text{SG}(\mathbf{W}_f x^{(2)})\|^2 + \|\mathbf{W}_g \mathbf{W}_h \mathbf{W}_f x^{(1)} - \text{SG}(\mathbf{W}_f x)\|^2 \right].$$

We will analyse the gradient flow $\dot{\mathbf{W}}_{\{f,g,h\}} = -\nabla \bar{L} - \rho \mathbf{W}_{\{f,g,h\}}$ ($\rho > 0$: weight decay intensity) subsequently. The gradient flows are derived as follows (see Appendix B.1):

$$\begin{aligned} \dot{\mathbf{W}}_f &= -\mathbf{W}_h^\top \{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) \mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\} \mathbf{W}_f - \rho \mathbf{W}_f, \\ \dot{\mathbf{W}}_g &= -\{(1 + \sigma^2) \mathbf{W}_h - \mathbf{I}\} \mathbf{W}_f \mathbf{W}_f^\top \mathbf{W}_h^\top - \rho \mathbf{W}_g, \\ \dot{\mathbf{W}}_h &= -\{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) \mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\} \mathbf{W}_f \mathbf{W}_f^\top - \rho \mathbf{W}_h. \end{aligned}$$

Eigenvalue dynamics. The matrix dynamics we have derived are rigorous but not amenable to further analysis. Here, we decouple the matrix dynamics into the eigenvalue dynamics. Let $\Phi := \mathbf{W}_f \mathbf{W}_f^\top \in \mathbb{R}^{m \times m}$. Following Tian et al. (2021, Theorem 3) and Bao (2023, Proposition 1), we can show that the eigenspaces of Φ , \mathbf{W}_g , and \mathbf{W}_h quickly align as t increases (see Appendix B.2). Therefore, we assume the following conditions:

- (A1) \mathbf{W}_g and \mathbf{W}_h are symmetric.
- (A2) The eigenspaces of Φ , \mathbf{W}_g , and \mathbf{W}_h align for every time step t .

Under these assumptions, Φ , \mathbf{W}_g , and \mathbf{W}_h are simultaneously diagonalizable and can be written as $\Phi = \mathbf{U} \Lambda_\Phi \mathbf{U}^\top$, $\mathbf{W}_g = \mathbf{U} \Lambda_g \mathbf{U}^\top$, and $\mathbf{W}_h = \mathbf{U} \Lambda_h \mathbf{U}^\top$, where \mathbf{U} is the (time-dependent) common orthogonal eigenvectors. Here, $\Lambda_\Phi = \text{diag}[\phi_1, \dots, \phi_m]$, $\Lambda_g = \text{diag}[\gamma_1, \dots, \gamma_m]$, and

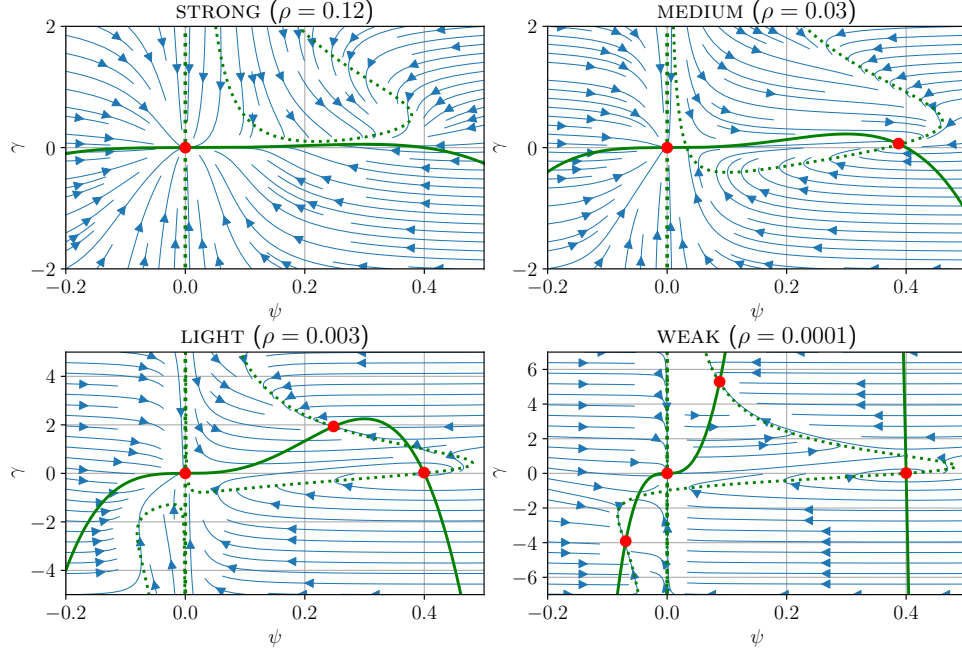


Figure 3: State space diagrams of PhiNet dynamics with different levels of weight decay: STRONG ($\rho = 0.12$), MEDIUM ($\rho = 0.03$), LIGHT ($\rho = 0.003$), and WEAK ($\rho = 0.0001$). The vector fields are numerically computed with $\sigma^2 = 1.5$. The state space bifurcates at the boundary of each level. The nullclines are shown with the green real ($\dot{\psi} = 0$) and dotted ($\dot{\gamma} = 0$) lines. The red dots are sinks.

$\Lambda_h = \text{diag}[\psi_1, \dots, \psi_m]$ are the corresponding eigenvalues. Noting that the dynamics quickly falls on to $\phi(t) = \psi(t)^2$, we can decouple the matrix dynamics into the eigenvalue dynamics of (ψ, γ) only (shown in Appendix B.3 and B.4):

$$\text{(PhiNet-dynamics)} \quad \begin{cases} \dot{\psi} &= \{(1 + \gamma) - (1 + \sigma^2)\psi\}\psi^2 - \rho\psi, \\ \dot{\gamma} &= \{1 - (1 + \sigma^2)\psi\}\psi^3 - \rho\gamma. \end{cases} \quad (1)$$

From the (ψ, γ) -dynamics, it is easy to see that $(\psi, \gamma) = (0, 0)$ is one of the equilibrium points. Can the eigenvalues escape from this collapsed solution?

Bifurcation of PhiNet dynamics. The state space diagrams of dynamics (1) are shown in Figure 3. In this figure, the nullclines $\dot{\psi} = 0$ and $\dot{\gamma} = 0$ are shown in the green real and dotted lines, respectively. Noting that intersecting points of nullclines are equilibrium points (Hirsch et al., 2012), we observe saddle-node bifurcation of PhiNet dynamics parametrized by weight decay $\rho > 0$.

- STRONG: Weight decay ρ is too strong that the collapsed point $(\psi, \gamma) = (0, 0)$ is a unique sink.
- MEDIUM: A new sink (ψ, γ) such that $\psi \gg 0$ and $\gamma \approx 0$ emerges. The number of sinks is two.
- LIGHT: Another non-trivial sink (ψ, γ) such that $\psi, \gamma \gg 0$ emerges. There are three sinks.
- WEAK: The last sink emerges such that $\psi < 0$ and $\gamma \ll 0$. The number of sinks is four.

Comparison with SimSiam dynamics. Tian et al. (2021) derived the SimSiam dynamics under the same setup as above. Specifically, they modelled the encoder f and the predictor h with linear networks $\mathbf{W}_f \mathbf{x}$ and $\mathbf{W}_h \mathbf{z}$, respectively, and defined the gradient flow dynamics with the MSE loss $\|\mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x}^{(2)})\|^2$ (without the additional predictor g). By decoupling the matrix dynamics into the eigenvalues with the same adiabatic elimination $\phi = \psi^2$, we can derive the SimSiam dynamics solely with respect to ψ -dynamics as follows:

$$\text{(SimSiam-dynamics)} \quad \dot{\psi} = \{1 - (1 + \sigma^2)\psi\}\psi^2 - \rho\psi. \quad (2)$$

We set $\tau = 1$ (ablating the exponential moving average used in BYOL) in Tian et al. (2021, Eq. (16)) to obtain this dynamics. SimSiam is free from the additional predictor g , so the dynamics (2) is univariate, unlike the bivariate system (1). Figure 4 shows the dynamics (2).

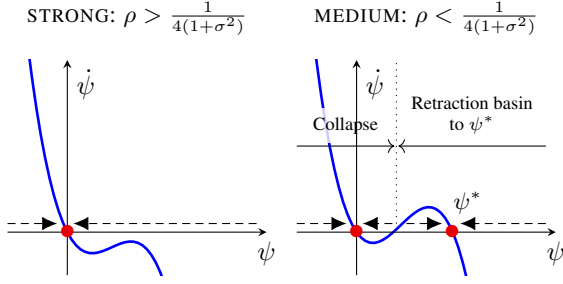


Figure 4: Illustration of SimSiam dynamics (2). Unlike the bivariate PhiNet dynamics shown in Figure 3, SimSiam dynamics is univariate, shown in the ψ -axis. The red dots are sinks.

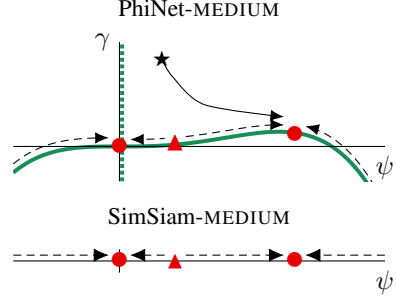


Figure 5: The SimSiam-MEDIUM flow is conjugate with the flow on the nullcline $\dot{\psi} = 0$ (green real line) in PhiNet-MEDIUM.

Table 1: **PhiNet is comparable to SimSiam.** We trained the models for 100 epochs and then validated them on the test sets using linear probing on the head. We trained with three seeds and calculated means and variances (subscripts). Both are unstable when the weight decay is small, but PhiNet still achieves high accuracy.

	Accuracy by Linear Probing (w.r.t. weight decay)					
	0.0	0.00001	0.00002	0.00005	0.0001	0.0002
SimSiam	25.41 _{0.02}	2.63 _{0.18}	60.82 _{1.57}	44.51 _{33.64}	68.17 _{0.18}	67.12 _{0.13}
PhiNet (MSE)	49.89 _{0.35}	55.90 _{1.57}	33.92 _{7.42}	66.73 _{0.03}	68.25 _{0.21}	67.83 _{0.15}

The SimSiam dynamics bifurcates into STRONG and MEDIUM at $\rho = 1/4(1 + \sigma^2)$. These two modes correspond to STRONG and MEDIUM of PhiNet in that ψ -axis of Figure 4 and the nullcline $\dot{\psi} = 0$ (green real line) in Figure 3 are topologically conjugate. The other LIGHT and WEAK are peculiar to the PhiNet dynamics. By comparing Figures 3 and 4, we have the following observations:

- (C1) *The retraction basin to non-collapsed solutions is wider:* Since SimSiam dynamics is univariate, ψ cannot avoid collapse once $\psi(0)$ is initialized outside the retraction basin to the non-collapse point $\psi^* \neq 0$ (namely, smaller than the source point \blacktriangle in Figure 5). By contrast, PhiNet avoids collapse even if $\psi(0)$ is close to zero, as long as $\gamma(0)$ is sufficiently large (see the initial point \star in Figure 5).
- (C2) *Even negative initialization ψ can avoid collapse:* In SimSiam-MEDIUM, ψ cannot be attracted to the non-collapsed solution if $\psi(0)$ is initialized to negative. By contrast, PhiNet-WEAK has a negative sink (at the bottom left in Figure 3), which attracts negative initialization $\psi(0) < 0$.

To sum it up, we have witnessed with a toy model that PhiNet is advantageous over SimSiam because the collapsed solution can be avoided more easily. This is why another predictor g is beneficial.

Remark 1. *The learning dynamics analysis in this section reveals that smaller weight decay ρ brings us benefits only regarding the stability of non-collapsed solutions. Indeed, we may benefit from larger ρ to accelerate convergence to the invariant parabola and eigenspace alignment of $(\Phi, \mathbf{W}_h, \mathbf{W}_g)$ (Appendices B.4 and B.2), each of which corresponds to the positive effects #3 and #7 in Tian et al. (2021), respectively. Moreover, moderately large ρ often yields good generalization in non-contrastive learning (Cabannes et al., 2023). Thus, smaller ρ may not be a silver bullet.*

5 EXPERIMENTS

We first test the robustness of PhiNets against the design choice and weight decay hyperparameter. We then discuss the effectiveness of X-PhiNets in online and continual learning.

5.1 LINEAR PROBING ANALYSIS

Figure 6 and Table 1 show the sensitivity analysis using CIFAR10 (Krizhevsky, 2009) and ImageNet (Krizhevsky et al., 2012), respectively, by changing the weight decay parameter. First, we

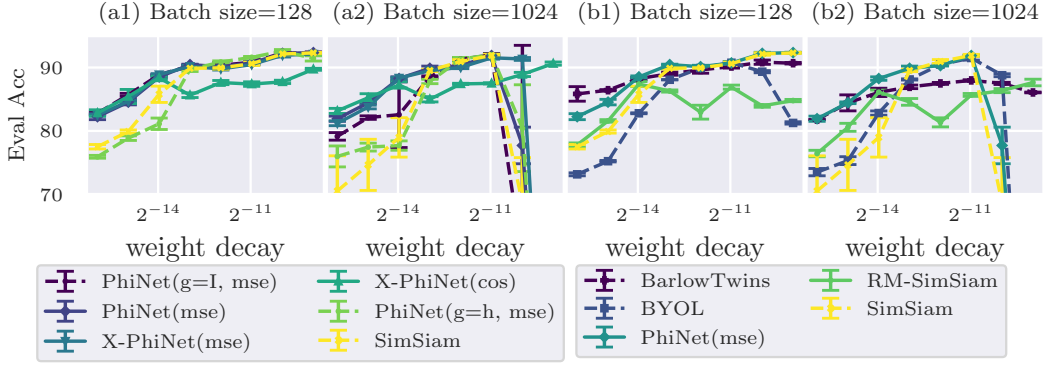


Figure 6: **PhiNet and X-PhiNet are robust against weight decay.** We evaluated PhiNet variants in (a1)-(a2) and compared the existing non-contrastive methods with PhiNet in (b1)-(b2) on CIFAR10. We evaluated the performance using linear probing. The loss function in brackets represents neocortex loss $L_{NC}(\theta)$. PhiNets perform particularly better than the baselines when weight decay is small.

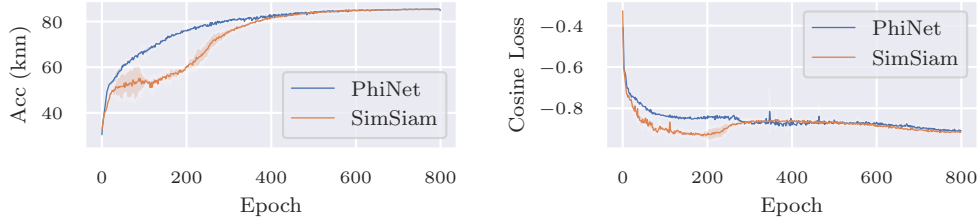


Figure 7: **PhiNet is stable in the early stages of learning.** We trained PhiNet and SimSiam with a batch size of 1024 and the weight decay of $1e-4$ on STL10. SimSiam is unstable in the early stages of learning. This may be due to the cosine loss being too small in SimSiam.

emphasise the improvement of PhiNet over SimSiam for most of the setups, supporting the importance of the CA1 predictor and Sim-2 loss. Subsequently, we closely look at the results.

PhiNet improves SimSiam. We observed weight decay significantly impacts the final model performance. When the MSE loss is used for Sim-2, PhiNet consistently outperforms the original SimSiam or other baselines (BYOL, RM-SimSiam) regardless of weight decay value, shown in Figure 6 (right). Moreover, as shown in Figure 7, PhiNet have a stabilizing effect during the early stages of training. This can likely be attributed to PhiNet’s regularization effect, which prevents the cosine loss from becoming too small at the early phase of training.

Bless of additional CA1 predictor. To see whether the additional predictor g besides h is beneficial, we test variants of predictor g : $g = h$ (reminiscent of the recurrent structure in CA3) and $g = \mathbf{I}$ (identity predictor). For CIFAR10 with batch size = 128, Figure 6 (left) indicates that the predictor choice slightly affects the final model performance if we properly set the weight decay. However, if we set the batch size as 1024, the separate predictor performs more stably over other choices.

Sim-2 loss should be MSE. Based on Figure 6 and Table 11 in the appendix, we found that the MSE loss used for Sim-2 generally improves model performance across most weight decay parameters, while the negative cosine loss performs comparably to the MSE loss with smaller weight decay but degrades it with larger weight decay.

Overall, our sensitivity study on CIFAR10 revealed that PhiNets are robust to the choice of the weight decay parameter, which supports the importance of the CA1 predictor and Sim-2 loss. See Appendix E.1 for more detailed sensitivity studies. In addition, the results for different batch sizes and datasets (STL10 (Coates et al., 2011)) can be found in Appendix E.4.

5.2 ONLINE LEARNING AND CONTINUAL LEARNING

SimSiam and other non-contrastive methods typically require up to 800 epochs of training on CIFAR10, which is quite different from the online nature of brains. To address this, we conducted

Table 2: **X-PhiNet performs good results when memorization is important.** We trained PhiNets on CIFAR-5m and Split CIFAR-5m. In Split CIFAR-5m, Acc is the average of the final accuracy (higher is better), and Fg is Forgetting (smaller is better). We present the results for two different weight decay ($5e-4$ and $2e-5$) in Split CIFAR-5m.

		BYOL	SimSiam	Barlow Twins	PhiNet (MSE)	RM-SimSiam	X-PhiNet (MSE)	X-PhiNet (Cos)
CIFAR-5m		81.05 _{0.04}	77.71 _{1.97}	85.32 _{0.10}	76.74 _{1.82}	82.09 _{0.22}	87.30 _{0.13}	87.46_{0.14}
Split C-5m	Acc	90.44 _{0.28}	90.84 _{0.31}	90.30 _{0.17}	90.69 _{0.11}	90.04 _{0.11}	91.02 _{0.36}	92.83_{0.12}
(wd= $5e-4$)	Fg	1.61 _{0.50}	2.45 _{0.42}	3.36 _{1.10}	2.96 _{0.23}	2.44 _{0.22}	3.44 _{0.36}	1.95 _{0.19}
Split C-5m	Acc	86.87 _{0.12}	88.20 _{0.35}	89.86 _{0.34}	88.60 _{0.15}	87.07 _{0.36}	90.90_{0.38}	90.72_{0.23}
(wd= $2e-5$)	Fg	-0.16 _{0.23}	0.36 _{0.51}	1.29 _{0.94}	0.05 _{0.19}	0.82 _{0.14}	-1.03 _{0.41}	0.43 _{0.17}

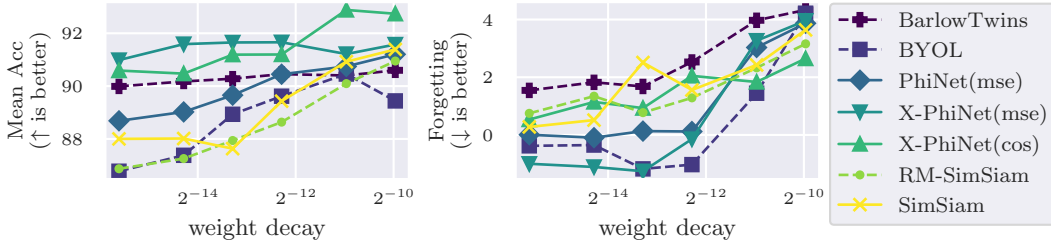


Figure 8: **X-PhiNet is also robust to weight decay in continual learning.** We measured the mean accuracy and forgetting at different weight decay on Split CIFAR-5m.

experiments using the CIFAR-5m dataset, which has six million synthetic CIFAR10-like images generated by the DDPM generative model (Nakkiran et al., 2021). Instead of training CIFAR10 with 50k samples for 800 epochs, we trained CIFAR-5m with 5m samples for 8 epochs. Although this is not exactly online learning, it seems closer to online learning compared to CIFAR10 due to the restriction on the training epochs. Table 2 shows that X-PhiNet has higher accuracy than SimSiam and PhiNet. The superior performance of X-PhiNet compared to PhiNet suggests that long-term memory with EMA is important in online learning. Sensitivity to weight decay and results for one-epoch online learning are given in Appendix D.1.

X-PhiNet draws inspiration from CLS theory, which proposes a framework for understanding continual learning processes in human brains. To evaluate the effectiveness of X-PhiNet in continual learning, we created a split CIFAR-5m dataset from CIFAR-5m, dividing it into five tasks, each with two classes. We trained on each task for one epoch and evaluated performance by the average accuracy across all tasks and the average forgetting, which is the difference between the peak accuracy and the final accuracy of each task. Table 2 shows that X-PhiNet has higher performance than SimSiam while maintaining minimal forgetting. Figure 8 further demonstrates that X-PhiNet consistently outperforms other methods like SimSiam in continual learning, regardless of weight decay. X-PhiNet also demonstrates high performance on Split CIFAR10 and Split CIFAR100, as well as when using replay methods (Appendix D.2).

6 CONCLUSION

In this paper, we proposed PhiNets based on non-contrastive learning with the temporal prediction hypothesis. Specifically, we leveraged StopGradient to artificially simulate the synaptic delay, and the prediction errors are modelled via Sim-1 and Sim-2 losses. Through theoretical analysis of learning dynamics, we showed that the proposed PhiNets have an advantage over SimSiam by more easily avoiding collapsed solutions. We empirically validated that the proposed PhiNets are robust with respect to weight decay and favorably comparable with SimSiam in terms of final classification performance. Experimental results also show that X-PhiNet performs better than SimSiam in online and continual learning, where memory function matters. These findings corroborate the effectiveness of the temporal prediction hypothesis when robustness and adaptivity are important.

ACKNOWLEDGEMENT

The authors thank Kenji Doya and Tomoki Fukaki for their helpful comments. We also thank Rio Yokota for providing computational resources. M.Y. and H.B. were supported by JSPS KAKENHI Grant Number 24K03004, H.B. was supported by JST PRESTO Grant Number JPMJPR24K6, and Y.T. was supported by JSPS KAKENHI Grant Number 23KJ1336.

REFERENCES

- P. Awasthi, N. Dikkala, and P. Kamath. Do more negative samples necessarily hurt in contrastive learning? In *ICML*, 2022.
- H. Bao. Feature normalization prevents collapse of non-contrastive learning dynamics. *arXiv preprint arXiv:2309.16109*, 2023.
- H. Bao, Y. Nagano, and K. Nozawa. On the surrogate gap between contrastive and supervised losses. In *ICML*, 2022.
- A. Bardes, Q. Garrido, J. Ponce, X. Chen, M. Rabbat, Y. LeCun, M. Assran, and N. Ballas. Revisiting feature prediction for learning visual representations from video. *arXiv:2404.08471*, 2024.
- S. Becker and G. E. Hinton. Self-organizing neural network that discovers surfaces in random-dot stereograms. *Nature*, 355(6356):161–163, 1992.
- V. Cabannes, B. Kiani, R. Balestriero, Y. LeCun, and A. Bietti. The SSL interplay: Augmentations, inductive bias, and generalization. In *ICML*, 2023.
- M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020.
- M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021.
- T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020a.
- X. Chen and K. He. Exploring simple Siamese representation learning. In *CVPR*, 2021.
- X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.
- Y. Chen, H. Zhang, M. Cameron, and T. Sejnowski. Predictive sequence learning in the hippocampal formation. *Neuron*, 112:2645–2658, 2024.
- A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. In *AISTATS*, 2011.
- E. Fini, V. G. T. Da Costa, X. Alameda-Pineda, E. Ricci, K. Alahari, and J. Mairal. Self-supervised models are continual learners. In *CVPR*, 2022.
- K. Friston. A theory of cortical responses. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 360(1456):815–836, 2005.
- F. Fu, Y. Gao, Z. Lu, H. Wu, and S. Zhao. Unsupervised continual learning of image representation via rememory-based simsiam. In *ICASSP*, 2024.
- J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, B. Piot, k. kavukcuoglu, R. Munos, and M. Valko. Bootstrap your own latent-a new approach to self-supervised learning. In *NeurIPS*, 2020.
- M. S. Halvagal, A. Laborieux, and F. Zenke. Implicit variance regularization in non-contrastive SSL. In *NeurIPS*, 2023.
- T. Han, W. Xie, and A. Zisserman. Video representation learning by dense predictive coding. In *CVPR Workshops*, 2019.

- T. Han, W. Xie, and A. Zisserman. Memory-augmented dense predictive coding for video representation learning. In *ECCV*, 2020.
- K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- O. Henaff. Data-efficient image recognition with contrastive predictive coding. In *ICML*, 2020.
- M. W. Hirsch, S. Smale, and R. L. Devaney. *Differential Equations, Dynamical Systems, and an Introduction to Chaos*. Academic Press, 2012.
- Y.-C. Hsu, Y.-C. Liu, A. Ramasamy, and Z. Kira. Re-evaluating continual learning scenarios: A categorization and case for strong baselines. *arXiv:1810.12488*, 2018.
- A. Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- D. Kumaran, D. Hassabis, and J. L. McClelland. What learning systems do intelligent agents need? complementary learning systems theory updated. *Trends in Cognitive Sciences*, 20:512–534, 2016.
- Z. Lin, Y. Wang, and H. Lin. Continual contrastive learning for image classification. In *ICME*, 2022.
- Z. Liu, E. S. Lubana, M. Ueda, and H. Tanaka. What shapes the loss landscape of self supervised learning? In *ICLR*, 2023.
- D. Madaan, J. Yoon, Y. Li, Y. Liu, and S. J. Hwang. Representational continuity for unsupervised continual learning. In *ICLR*, 2022.
- J. L. McClelland, B. L. McNaughton, and R. C. O’Reilly. Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological Review*, 102(3):419, 1995.
- D. Mumford. On the computational architecture of the neocortex: II the role of cortico-cortical loops. *Biological Cybernetics*, 66(3):241–251, 1992.
- P. Nakkiran, B. Neyshabur, and H. Sedghi. The deep bootstrap framework: Good online learners are good offline generalizers. In *ICLR*, 2021.
- Q. Pham, C. Liu, and S. Hoi. DualNet: Continual learning, fast and slow. In *NeurIPS*, 2021.
- Q. Pham, C. Liu, and S. C. Hoi. Continual learning, fast and slow. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- R. P. Rao and D. H. Ballard. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1):79–87, 1999.
- O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115:211–252, 2015.
- F. Sarnthein, G. Bachmann, S. Anagnostidis, and T. Hofmann. Random teachers are good teachers. In *ICML*, 2023.
- J. Smith, C. Taylor, S. Baer, and C. Dovrolis. Unsupervised progressive learning and the STAM architecture. In *IJCAI*, 2021.
- M. V. Srinivasan, S. B. Laughlin, and A. Dubs. Predictive coding: a fresh view of inhibition in the retina. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, 216(1205): 427–459, 1982.
- R. Tan, M. De Lange, M. Iuzzolino, B. A. Plummer, K. Saenko, K. Ridgeway, and L. Torresani. Multiscale video pretraining for long-term activity forecasting. *arXiv:2307.12854*, 2023.

- Y. Tang, Z. D. Guo, P. H. Richemond, B. Avila Pires, Y. Chandak, R. Munos, M. Rowland, M. Gheshlaghi Azar, C. Le Lan, C. Lyle, A. György, S. Thakoor, W. Dabney, B. Piot, D. Calandriello, and M. Valko. Understanding self-predictive learning for reinforcement learning. In *ICML*, 2023.
- Y. Tian, X. Chen, and S. Ganguli. Understanding self-supervised learning dynamics without contrastive pairs. In *ICML*, 2021.
- G. M. Van de Ven, H. T. Siegelmann, and A. S. Tolias. Brain-inspired replay for continual learning with artificial neural networks. *Nature Communications*, 2020.
- A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv:1807.03748*, 2018.
- N. Vyas, A. Atanasov, B. Bordelon, D. Morwani, S. Sainathan, and C. Pehlevan. Feature-learning networks are consistent across widths at realistic scales. In *NeurIPS*, 2023.
- J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow Twins: Self-supervised learning via redundancy reduction. In *ICML*, 2021.

A LIMITATIONS AND FUTURE WORK (EXTENDED VERSION)

One major limitation of our approach is the use of backpropagation, which differs from the mechanisms in biological neural networks. Our long-term goal is to eliminate backpropagation to better imitate brain function, but this work focuses on the model’s structural aspects. Currently, backpropagation-free predictive coding mechanisms for complex architectures like ResNet are in the early stages of development, with most research limited to simple CNNs. Conversely, non-contrastive methods like SimSiam require more advanced models than ResNet. Future research should explore if the proposed structure can enable effective learning with backpropagation-free predictive coding. Another key difference between PhiNet and brains is the presence of recurrent structures. However, in this PhiNet, only one time step is considered, so it is possible that the recurrent structure required to predict time series data was not necessary. Making the data into time series data and adding a recurrent structure to the model remains as future work.

It is also unclear whether cosine loss or MSE loss is more suitable for the Sim-2 in PhiNets. Cosine loss performs better when weight decay is small or online and continual learning where the additional predictor of PhiNets is important. However, MSE loss is preferable when weight decay is large on CIFAR10. This is likely because using cosine loss in sim-2 has a stronger impact on learning dynamics compared to MSE loss. Analyzing gradient norms could be useful for this kind of evaluation, but is left for future work.

B DETAILS OF LEARNING DYNAMICS ANALYSIS

In this section, we complement the missing details of learning dynamics analysis provided in Section 4. In our analysis, we will use the gradient flow $\dot{\mathbf{W}}_{\{f,g,h\}} = -\nabla \bar{L} - \rho \mathbf{W}_{\{f,g,h\}}$ ($\rho > 0$), which is the continuous limit of gradient descent. This corresponds to considering the following gradient descent in discrete updates and taking the limit as $\eta \rightarrow 0$.

$$\mathbf{W}_{\{f,g,h\}}(t + \eta) = \mathbf{W}_{\{f,g,h\}}(t) - \eta \nabla \bar{L} - \eta \rho \mathbf{W}_{\{f,g,h\}} \quad (3)$$

$$\frac{1}{\eta} (\mathbf{W}_{\{f,g,h\}}(t + \eta) - \mathbf{W}_{\{f,g,h\}}(t)) = -\nabla \bar{L} - \rho \mathbf{W}_{\{f,g,h\}} \quad (4)$$

B.1 DERIVATION OF MATRIX DYNAMICS

Recall the PhiNet loss function:

$$\bar{L} := \frac{1}{2} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}^{(1)}, \mathbf{x}^{(2)} | \mathbf{x}} \left[\|\mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x}^{(2)})\|^2 + \|\mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x})\|^2 \right].$$

Let us derive its matrix gradient.

$$\begin{aligned} \nabla_{\mathbf{W}_f} \bar{L} &= \frac{1}{2} \nabla_{\mathbf{W}_f} \mathbb{E} \left[(\mathbf{x}^{(1)\top} \mathbf{W}_f^\top \mathbf{W}_h^\top - \text{SG}(\mathbf{x}^{(2)\top} \mathbf{W}_f^\top)) (\mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x}^{(2)})) \right. \\ &\quad \left. + (\mathbf{x}^{(1)} \mathbf{W}_f^\top \mathbf{W}_h^\top \mathbf{W}_g^\top - \text{SG}(\mathbf{x}^\top \mathbf{W}_f^\top)) (\mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x})) \right] \\ &= \left\{ \mathbf{W}_h^\top \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] - \mathbf{W}_h^\top \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(2)} \mathbf{x}^{(1)\top}] \right\} \\ &\quad + \left\{ \mathbf{W}_h^\top \mathbf{W}_g^\top \mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] - \mathbf{W}_h^\top \mathbf{W}_g^\top \mathbf{W}_f \mathbb{E}[\mathbf{x} \mathbf{x}^{(1)\top}] \right\} \\ &= \mathbf{W}_h^\top \left\{ (\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] - \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(2)} \mathbf{x}^{(1)\top}] - \mathbf{W}_g^\top \mathbf{W}_f \mathbb{E}[\mathbf{x} \mathbf{x}^{(1)\top}] \right\} \\ &= \mathbf{W}_h^\top \left\{ (1 + \sigma^2) (\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) \mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top) \right\} \mathbf{W}_f, \end{aligned}$$

where the last line is derived from our assumption on the data distributions:

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}^{(1)} | \mathbf{x}} [\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] &= \mathbb{E}_{\mathbf{x}} [\mathbf{x} \mathbf{x}^\top] + \sigma^2 \mathbf{I} = (1 + \sigma^2) \mathbf{I}, \\ \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}^{(1)}, \mathbf{x}^{(2)} | \mathbf{x}} [\mathbf{x}^{(2)} \mathbf{x}^{(1)\top}] &= \mathbb{E}_{\mathbf{x}} [\mathbf{x} \mathbf{x}^\top] = \mathbf{I}, \\ \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}^{(1)} | \mathbf{x}} [\mathbf{x} \mathbf{x}^{(1)\top}] &= \mathbb{E}_{\mathbf{x}} [\mathbf{x} \mathbf{x}^\top] = \mathbf{I}. \end{aligned}$$

Similarly, we derive $\nabla_{\mathbf{W}_g} \bar{L}$ and $\nabla_{\mathbf{W}_h} \bar{L}$.

$$\begin{aligned}\nabla_{\mathbf{W}_g} \bar{L} &= \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top \mathbf{W}_h^\top - \mathbf{W}_f \mathbb{E}[\mathbf{x} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top \mathbf{W}_h^\top \\ &= \{(1 + \sigma^2) \mathbf{W}_h - \mathbf{I}\} \mathbf{W}_f \mathbf{W}_f^\top \mathbf{W}_h^\top.\end{aligned}$$

$$\begin{aligned}\nabla_{\mathbf{W}_h} \bar{L} &= \left\{ \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top - \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(2)} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top \right\} \\ &\quad + \left\{ \mathbf{W}_g^\top \mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top - \mathbf{W}_g^\top \mathbf{W}_f \mathbb{E}[\mathbf{x} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top \right\} \\ &= \{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) \mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\} \mathbf{W}_f \mathbf{W}_f^\top.\end{aligned}$$

From these, we obtain the matrix dynamics.

B.2 EIGENSPACE ALIGNMENT

Our aim is to show that the three matrices Φ , \mathbf{W}_g , and \mathbf{W}_h share a common eigenspace, i.e., simultaneously diagonalizable, asymptotically in time t . Let

$$\mathbf{C}_1 := [\Phi, \mathbf{W}_g], \quad \mathbf{C}_2 := [\Phi, \mathbf{W}_h], \quad \text{and} \quad \mathbf{C}_3 := [\mathbf{W}_g, \mathbf{W}_h],$$

where $[\mathbf{A}, \mathbf{B}] := \mathbf{AB} - \mathbf{BA}$ is the commutator (matrix). By noting that commutative matrices are simultaneously diagonalizable, we show that the time-dependent commutators $\mathbf{C}_1(t)$, $\mathbf{C}_2(t)$, and $\mathbf{C}_3(t)$ asymptotically converges to \mathbf{O} as $t \rightarrow \infty$.

Hereafter, we assume the symmetry assumption (A1) on \mathbf{W}_g and \mathbf{W}_h , and heavily use the following formulas on commutators implicitly:

- $[\mathbf{A}, \mathbf{A}] = \mathbf{O}$.
- $[\mathbf{A}, \mathbf{B}] = -[\mathbf{B}, \mathbf{A}]$.
- $[\mathbf{A}, \mathbf{BC}] = [\mathbf{A}, \mathbf{B}]\mathbf{C} + \mathbf{B}[\mathbf{A}, \mathbf{C}]$.
- $[\mathbf{AB}, \mathbf{C}] = \mathbf{A}[\mathbf{B}, \mathbf{C}] + [\mathbf{A}, \mathbf{C}]\mathbf{B}$.

First, compute $\dot{\mathbf{C}}_1$ based on the matrix dynamics of Φ (can be found in Appendix B.3) and \mathbf{W}_g :

$$\begin{aligned}\dot{\mathbf{C}}_1 &= \dot{\Phi} \mathbf{W}_g + \Phi \dot{\mathbf{W}}_g - \dot{\mathbf{W}}_g \Phi - \mathbf{W}_g \dot{\Phi} \\ &= \{-(1 + \sigma^2)(\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2) \mathbf{W}_h \Phi + \Phi \mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2) \mathbf{W}_h) + (\mathbf{W}_h \Phi + \Phi \mathbf{W}_h) \\ &\quad + (\mathbf{W}_h \mathbf{W}_g \Phi + \Phi \mathbf{W}_g \mathbf{W}_h) - 2\rho \Phi\} \mathbf{W}_g + \Phi \{-(1 + \sigma^2) \mathbf{W}_h - \mathbf{I}\} \Phi \mathbf{W}_h - \rho \mathbf{W}_g\} \\ &\quad + \{(1 + \sigma^2) \mathbf{W}_h - \mathbf{I}\} \Phi \mathbf{W}_h + \rho \mathbf{W}_g\} \Phi \\ &\quad + \mathbf{W}_g \{(1 + \sigma^2)(\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2) \mathbf{W}_h \Phi + \Phi \mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2) \mathbf{W}_h) - (\mathbf{W}_h \Phi + \Phi \mathbf{W}_h) \\ &\quad - (\mathbf{W}_h \mathbf{W}_g \Phi + \Phi \mathbf{W}_g \mathbf{W}_h) + 2\rho \Phi\} \\ &= -3\rho \mathbf{C}_1 + (1 + \sigma^2)[\mathbf{W}_g, \mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2) \mathbf{W}_h \Phi] + (1 + \sigma^2)[\mathbf{W}_g, \Phi \mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2) \mathbf{W}_h] \\ &\quad + [\mathbf{W}_h \Phi + \Phi \mathbf{W}_h, \mathbf{W}_g] + [\mathbf{W}_h, \mathbf{W}_g \Phi \mathbf{W}_g] + [\Phi, \mathbf{W}_g \mathbf{W}_h \mathbf{W}_g] \\ &\quad + [(1 + \sigma^2) \mathbf{W}_h - \mathbf{I}] \Phi \mathbf{W}_h, \Phi \\ &= -3\rho \mathbf{C}_1 + (1 + \sigma^2)\{(\mathbf{C}_3 \mathbf{W}_h \Phi + \Phi \mathbf{W}_h \mathbf{C}_3) + (\mathbf{W}_h \mathbf{C}_3 \Phi + \Phi \mathbf{C}_3 \mathbf{W}_h) \\ &\quad + (\mathbf{C}_3 \mathbf{W}_g^2 \mathbf{W}_h + \mathbf{W}_h \mathbf{W}_g^2 \mathbf{C}_3) \Phi - (\mathbf{W}_h \mathbf{W}_g^2 \mathbf{W}_h \mathbf{C}_1 + \mathbf{C}_1 \mathbf{W}_h \mathbf{W}_g^2 \mathbf{W}_h) \\ &\quad - (\mathbf{W}_h^2 \mathbf{C}_1 + \mathbf{C}_1 \mathbf{W}_h^2) + \Phi(\mathbf{C}_3 \mathbf{W}_g^2 \mathbf{W}_h + \mathbf{W}_h \mathbf{W}_g^2 \mathbf{C}_3)\} \\ &\quad + (\mathbf{W}_h \mathbf{C}_1 + \mathbf{C}_1 \mathbf{W}_h) - (\mathbf{C}_3 \Phi + \Phi \mathbf{C}_3) - (\mathbf{C}_3 \Phi \mathbf{W}_g + \mathbf{W}_g \Phi \mathbf{C}_3) \\ &\quad + (\mathbf{C}_1 \mathbf{W}_h \mathbf{W}_g + \mathbf{W}_g \mathbf{W}_h \mathbf{C}_1) - (1 + \sigma^2)(\mathbf{W}_h \Phi \mathbf{C}_2 + \mathbf{C}_2 \Phi \mathbf{W}_h) + \Phi \mathbf{C}_2.\end{aligned}$$

Similarly, $\dot{\mathbf{C}}_2$ and $\dot{\mathbf{C}}_3$ are computed:

$$\begin{aligned}\dot{\mathbf{C}}_2 &= -3\rho\mathbf{C}_2 + (\mathbf{C}_2\mathbf{W}_h + \mathbf{W}_h\mathbf{C}_2) - \mathbf{C}_1\Phi + (\mathbf{W}_g\mathbf{C}_2 + \mathbf{C}_2\mathbf{W}_g) + (\mathbf{C}_3\Phi + \Phi\mathbf{C}_3) \\ &\quad - (1 + \sigma^2)\mathbf{C}_2\Phi - (1 + \sigma^2)(\mathbf{W}_g\mathbf{C}_1 + \mathbf{C}_1\mathbf{W}_g) + (1 + \sigma^2)\mathbf{W}_h(\mathbf{C}_3\mathbf{W}_g + \mathbf{W}_g\mathbf{C}_3) \\ &\quad - (1 + \sigma^2)(\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h\mathbf{C}_2 + \mathbf{C}_2\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h) \\ &\quad + (1 + \sigma^2)\Phi\mathbf{W}_h(\mathbf{C}_3\mathbf{W}_g + \mathbf{W}_g\mathbf{C}_3)\mathbf{W}_h, \\ \dot{\mathbf{C}}_3 &= -3\rho\mathbf{C}_3 + (\mathbf{I} - (1 + \sigma^2)\mathbf{W}_h)\mathbf{C}_2\mathbf{W}_h + (1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^2)(\mathbf{W}_h\mathbf{C}_1 - \mathbf{C}_3\Phi) \\ &\quad + (\mathbf{I} + \mathbf{W}_g)\mathbf{C}_1.\end{aligned}$$

Next, we vectorize the commutator matrices—for $\mathbf{C} \in \mathbb{R}^{h \times h}$, $\text{vec}(\mathbf{C}) \in \mathbb{R}^{h^2}$ indicates a (column) vector stacking the columns of \mathbf{C} . For the commutators \mathbf{C}_1 , \mathbf{C}_2 , and \mathbf{C}_3 , let us write $\xi := \text{vec}(\mathbf{C}_1)$, $\eta := \text{vec}(\mathbf{C}_2)$, and $\zeta := \text{vec}(\mathbf{C}_3)$. In what follows, we heavily leverage the vectorization formula:

- $\text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \otimes \mathbf{A})\text{vec}(\mathbf{B}) = (\mathbf{I} \otimes \mathbf{AB})\text{vec}(\mathbf{C}) = (\mathbf{C}^\top \mathbf{B}^\top \otimes \mathbf{I})\text{vec}(\mathbf{A})$
- $\text{vec}(\mathbf{AB}) = (\mathbf{I} \otimes \mathbf{A})\text{vec}(\mathbf{B}) = (\mathbf{B}^\top \otimes \mathbf{I})\text{vec}(\mathbf{A})$

Here, $\mathbf{A} \otimes \mathbf{B}$ denotes the Kronecker product of two matrices. We write $\mathbf{A} \oplus \mathbf{B} := \mathbf{A} \otimes \mathbf{B} + \mathbf{B} \otimes \mathbf{A}$ for notational convenience. We derive the ODE of $\xi = \text{vec}(\mathbf{C}_1)$ as follows:

$$\begin{aligned}\dot{\xi} &= -3\rho\mathbf{I}\xi + (1 + \sigma^2)((\Phi\mathbf{W}_h \oplus \mathbf{I})\zeta + (\Phi \oplus \mathbf{W}_h)\zeta + (\Phi \otimes \mathbf{I})(\mathbf{W}_h\mathbf{W}_g^2 \oplus \mathbf{I})\zeta \\ &\quad - (\mathbf{W}_h\mathbf{W}_g^2\mathbf{W}_h \oplus \mathbf{I})\xi - (\mathbf{I} \oplus \mathbf{W}_h^2)\xi + (\mathbf{I} \otimes \Phi)(\mathbf{W}_h\mathbf{W}_g^2 \oplus \mathbf{I})\zeta + (\mathbf{I} \oplus \mathbf{W}_h)\xi \\ &\quad - (\Phi \oplus \mathbf{I})\zeta - (\mathbf{W}_g\Phi \oplus \mathbf{I})\zeta + (\mathbf{W}_g\mathbf{W}_h \oplus \mathbf{I})\xi - (1 + \sigma^2)(\mathbf{I} \oplus \mathbf{W}_h\Phi)\eta + (\mathbf{I} \otimes \Phi)\eta \\ &= -\{3\rho\mathbf{I} + \mathbf{I} \oplus ((1 + \sigma^2)\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h) + \mathbf{W}_h(\mathbf{I} + \mathbf{W}_g)\}\xi \\ &\quad - \{(1 + \sigma^2)(\mathbf{I} \otimes \mathbf{W}_h\Phi) - \mathbf{I} \otimes \Phi\}\eta \\ &\quad - \{(\mathbf{I} + \mathbf{W}_g)\Phi \oplus \mathbf{I} - (1 + \sigma^2)(\Phi\mathbf{W}_h \oplus \mathbf{I} + (\mathbf{I} \oplus \Phi)(\mathbf{W}_h\mathbf{W}_g^2 \oplus \mathbf{I}))\}\zeta \\ &= -(3\rho\mathbf{I} + \mathbf{K}_{11})\xi - \mathbf{K}_{12}\eta - \mathbf{K}_{13}\zeta,\end{aligned}$$

where

$$\begin{aligned}\mathbf{K}_{11} &= \mathbf{I} \oplus ((1 + \sigma^2)\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h) + \mathbf{W}_h(\mathbf{I} + \mathbf{W}_g), \\ \mathbf{K}_{12} &= (1 + \sigma^2)(\mathbf{I} \otimes \mathbf{W}_h\Phi) - \mathbf{I} \otimes \Phi, \\ \mathbf{K}_{13} &= (\mathbf{I} + \mathbf{W}_g)\Phi \oplus \mathbf{I} - (1 + \sigma^2)(\Phi\mathbf{W}_h \oplus \mathbf{I} + (\mathbf{I} \oplus \Phi)(\mathbf{W}_h\mathbf{W}_g^2 \oplus \mathbf{I})).\end{aligned}$$

Similarly, we derive the ODEs of $\eta = \text{vec}(\mathbf{C}_2)$ and $\zeta = \text{vec}(\mathbf{C}_3)$.

$$\begin{aligned}\dot{\eta} &= -\mathbf{K}_{21}\xi - (3\rho\mathbf{I} + \mathbf{K}_{22})\eta - \mathbf{K}_{23}\zeta, \\ \dot{\zeta} &= -\mathbf{K}_{31}\xi - \mathbf{K}_{32}\eta - (3\rho\mathbf{I} + \mathbf{K}_{33})\zeta,\end{aligned}$$

where

$$\begin{aligned}\mathbf{K}_{21} &= \Phi \otimes \mathbf{I} + (1 + \sigma^2)(\mathbf{W}_g \oplus \mathbf{I}), \\ \mathbf{K}_{22} &= (1 + \sigma^2)(\Phi \otimes \mathbf{I}) + \mathbf{I} \oplus \{(1 + \sigma^2)(\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h - (\mathbf{W}_g + \mathbf{W}_h))\}, \\ \mathbf{K}_{23} &= -\mathbf{I} \oplus \Phi - (1 + \sigma^2)\{(\mathbf{I} \otimes \mathbf{W}_h) + (\mathbf{W}_h \otimes \Phi\mathbf{W}_h)\}(\mathbf{I} \oplus \mathbf{W}_g), \\ \mathbf{K}_{31} &= -(1 + \sigma^2)(\mathbf{I} \otimes (\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h) - \mathbf{I} \otimes (\mathbf{I} + \mathbf{W}_g), \\ \mathbf{K}_{32} &= -\mathbf{W}_h \otimes (\mathbf{I} - (1 + \sigma^2)\mathbf{W}_h), \\ \mathbf{K}_{33} &= (1 + \sigma^2)(\Phi \otimes (\mathbf{I} + \mathbf{W}_g^2)).\end{aligned}$$

By combining all the above, we obtain a single ODE for (ξ, η, ζ) :

$$\frac{d}{dt} \begin{bmatrix} \xi \\ \eta \\ \zeta \end{bmatrix} = - \underbrace{\begin{bmatrix} 3\rho\mathbf{I} + \mathbf{K}_{11} & \mathbf{K}_{12} & \mathbf{K}_{13} \\ \mathbf{K}_{21} & 3\rho\mathbf{I} + \mathbf{K}_{22} & \mathbf{K}_{23} \\ \mathbf{K}_{31} & \mathbf{K}_{32} & 3\rho\mathbf{I} + \mathbf{K}_{33} \end{bmatrix}}_{:= 3\rho\mathbf{I} + \mathbf{K}} \underbrace{\begin{bmatrix} \xi \\ \eta \\ \zeta \end{bmatrix}}_{:= \Xi},$$

or alternatively, $\dot{\Xi} = -(3\rho\mathbf{I} + \mathbf{K})\Xi$. Note that $\mathbf{K}(t)$ is time-dependent. Finally, we can obtain the desired result by invoking [Tian et al. \(2021, Lemma 2\)](#).

Lemma 1 (Tian et al. (2021, Lemma 2)). *Let $\mathbf{H}(t)$ be time-varying positive semidefinite matrices whose minimal eigenvalues are bounded away from zero:*

$$\inf_{t \geq 0} \lambda_{\min}(\mathbf{H}(t)) \geq \lambda_0 > 0.$$

Then, the following dynamics

$$\frac{d\mathbf{w}(t)}{dt} = -\mathbf{H}(t)\mathbf{w}(t)$$

satisfies $\|\mathbf{w}(t)\|_2 \leq \exp(-\lambda_0 t)\|\mathbf{w}(0)\|_2$, which means that $\mathbf{w}(t) \rightarrow \mathbf{0}$.

When minimal eigenvalues of $3\rho\mathbf{I} + \mathbf{K}(t)$ are always bounded away from zero, we immediately see $\Xi(t) \rightarrow \mathbf{0}$, namely, $(\mathbf{C}_1(t), \mathbf{C}_2(t), \mathbf{C}_3(t)) \rightarrow (\mathbf{0}, \mathbf{0}, \mathbf{0})$ as $t \rightarrow \infty$. The strict positive-definiteness of $3\rho\mathbf{I} + \mathbf{K}(t)$ would not be necessarily satisfied; however, larger weight decay $\rho > 0$ induces it more easily. The convergence of the commutators is faster with larger $\rho > 0$ as well.

B.3 DECOUPLING INTO EIGENVALUE DYNAMICS

We have obtained the following matrix dynamics:

$$\begin{aligned}\dot{\mathbf{W}}_f &= -\mathbf{W}_h^\top \{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g)\mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\}\mathbf{W}_f - \rho\mathbf{W}_f, \\ \dot{\mathbf{W}}_g &= -\{(1 + \sigma^2)\mathbf{W}_h - \mathbf{I}\}\mathbf{W}_f\mathbf{W}_f^\top\mathbf{W}_h^\top - \rho\mathbf{W}_g, \\ \dot{\mathbf{W}}_h &= -\{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g)\mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\}\mathbf{W}_f\mathbf{W}_f^\top - \rho\mathbf{W}_h.\end{aligned}$$

Our aim is to decouple the matrix dynamics into their eigenvalue counterparts. Beforehand, let us execute the change-of-variable $\Phi = \mathbf{W}_f\mathbf{W}_f^\top$:

$$\begin{aligned}\dot{\Phi} &= \dot{\mathbf{W}}_f\mathbf{W}_f^\top + \mathbf{W}_f\dot{\mathbf{W}}_f^\top \\ &= -\mathbf{W}_h^\top \{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g)\mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\}\Phi \\ &\quad - \Phi \{(1 + \sigma^2)\mathbf{W}_h^\top(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) - (\mathbf{I} + \mathbf{W}_g)\}\mathbf{W}_h - 2\rho\Phi.\end{aligned}$$

By the symmetry assumption (A1), $(\Phi, \mathbf{W}_g, \mathbf{W}_h)$ -dynamics can be simplified as follows:

$$\begin{aligned}\dot{\Phi} &= -(1 + \sigma^2)\{\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h, \Phi\} + \{\mathbf{W}_h, \Phi\} + (\mathbf{W}_h\mathbf{W}_g\Phi + \Phi\mathbf{W}_g\mathbf{W}_h) - 2\rho\Phi, \quad (5) \\ \dot{\mathbf{W}}_g &= -\{(1 + \sigma^2)\mathbf{W}_h - \mathbf{I}\}\Phi\mathbf{W}_h - \rho\mathbf{W}_g, \\ \dot{\mathbf{W}}_h &= -\{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h + (\mathbf{I} + \mathbf{W}_g)\}\Phi - \rho\mathbf{W}_h,\end{aligned}$$

where $\{\mathbf{A}, \mathbf{B}\} := \mathbf{AB} + \mathbf{BA}$ is the anticommutator for two symmetric matrices \mathbf{A} and \mathbf{B} with the same dimension.

Next, we decouple them into the corresponding eigenvalues. The parameter matrices are simultaneously diagonalized by $\Phi = \mathbf{U}\Lambda_\Phi\mathbf{U}^\top$, $\mathbf{W}_g = \mathbf{U}\Lambda_g\mathbf{U}^\top$, and $\mathbf{W}_h = \mathbf{U}\Lambda_h\mathbf{U}^\top$, with the aid of the symmetry assumption (A1) and common eigenspace assumption (A2). Here, we can easily show that the eigenspace is time-independent, namely, $\dot{\mathbf{U}} = \mathbf{0}$, using the same argument of Tian et al. (2021, Appendix B.1). By multiplying \mathbf{U}^\top and \mathbf{U} from left and right, respectively, Φ -dynamics can be written as follows:

$$\dot{\Lambda}_\Phi = -2(1 + \sigma^2)(\mathbf{I} + \Lambda_g^2)\Lambda_h^2\Lambda_\Phi + 2\Lambda_h\Lambda_\Phi + 2\Lambda_h\Lambda_g\Lambda_\Phi - 2\rho\Lambda_\Phi,$$

where all matrices are diagonal, and thus, we can write down the dynamics in terms of j -th diagonal element (but the index j is omitted for simplicity):

$$\dot{\phi} = -2(1 + \sigma^2)(1 + \gamma^2)\psi^2\phi + 2\psi\phi + 2\psi\phi\gamma - 2\rho\phi.$$

We can decouple \mathbf{W}_g - and \mathbf{W}_h -dynamics similarly:

$$\begin{aligned}\dot{\gamma} &= -(1 + \sigma^2)(\psi - 1)\phi\psi - \rho\gamma, \\ \dot{\psi} &= -\{(1 + \sigma^2)(1 + \gamma^2)\psi + (1 + \gamma)\}\phi - \rho\psi.\end{aligned}$$

Note that ψ and γ correspond to (one of the) eigenvalues of the linear networks h and g , respectively. Intuitively, we can regard ψ and γ as ‘‘scalarization’’ of the predictor networks.

To sum it up, we decouple the dynamics of $(\Phi, \mathbf{W}_h, \mathbf{W}_g)$ into the following dynamics of (ϕ, ψ, γ) :

$$(\Phi\text{-dynamics}) \quad \dot{\phi} = -2\psi\phi\{(1+\sigma^2)(1+\gamma^2)\psi - (1+\gamma)\} - 2\rho\phi, \quad (6)$$

$$(\mathbf{W}_h\text{-dynamics}) \quad \dot{\psi} = -\phi\{(1+\sigma^2)(1+\gamma^2)\psi - (1+\gamma)\} - \rho\psi, \quad (7)$$

$$(\mathbf{W}_g\text{-dynamics}) \quad \dot{\gamma} = -\psi\phi\{(1+\sigma^2)\psi - 1\} - \rho\gamma, \quad (8)$$

B.4 ADIABATIC ELIMINATION

The eigenvalue dynamics obtained in Appendix B.3 is jointly with respect to (ϕ, ψ, γ) . Here, we eliminate ϕ by confirming that ϕ and ψ are asymptotically bound on an *invariant parabola*.

By combining (6) and (7), we have $2\psi\dot{\psi} - \dot{\phi} = -2\rho(\psi^2 - \phi)$. This can be integrated, and we obtain the following solution:

$$\psi(t)^2 - \phi(t) = C \exp(-2\rho t) \xrightarrow{t \rightarrow \infty} 0, \quad (9)$$

where C is a constant of integration. Thus, $(\phi(t), \psi(t))$ converges to this invariant parabola (9) exponentially quickly, which we suppose is much faster than the dynamics stabilization. On this invariant parabola $\phi = \psi^2$, the eigenvalue dynamics can be further simplified as follows by eliminating ϕ :

$$\begin{cases} \dot{\psi} &= \{(1+\gamma) - (1+\sigma^2)(1+\gamma^2)\psi\}\psi^2 - \rho\psi, \\ \dot{\gamma} &= \{1 - (1+\sigma^2)\psi\}\psi^3 - \rho\gamma. \end{cases}$$

Note that the convergence to the invariant parabola is faster when weight decay ρ is more intense.

C PSEUDOCODE FOR PHINET AND X-PHINET

The pseudo codes for PhiNet and X-PhiNet are shown in Listing.1 and Listing.2.

```

1 # f: backbone + projection mlp
2 # h: prediction mlp
3 # g: prediction mlp
4
5 for x in loader: # load a minibatch x with n samples
6     x1, x2 = aug(x), aug(x) # random augmentation
7     z0, z1, z2 = f(x), f(x1), f(x2) # projections, n-by-d
8     p1, p2 = h(z1), h(z2) # predictions, n-by-d
9     y1, y2 = g(p1), g(p2) # predictions, n-by-d
10    z0 = z0.detach()
11    Lcos = D(p1, z2)/2 + D(p2, z1)/2 # loss
12    Lcor = mse_loss(y1, z0)/2 + mse_loss(y2, z0)/2
13    L = Lcos + Lcor
14    L.backward() # back-propagate
15    update(f, h) # SGD update
16
17 def D(p, z): # negative cosine similarity
18     z = z.detach() # stop gradient
19     p = normalize(p, dim=1) # l2-normalize
20     z = normalize(z, dim=1) # l2-normalize
21     return -(p*z).sum(dim=1).mean()
```

Listing 1: PhiNet Pseudocode (PyTorch-like)

```

1 # f: backbone + projection mlp
2 # h: prediction mlp
3 # g: prediction mlp
4
5 for x in loader: # load a minibatch x with n samples
6     x1, x2 = aug(x), aug(x) # random augmentation
7     z0, z1, z2 = f_long(x), f(x1), f(x2) # projections, n-by-d
8     p1, p2 = h(z1), h(z2) # predictions, n-by-d
9     y1, y2 = g(p1), g(p2) # predictions, n-by-d
10    z0 = z0.detach()
```


Table 3: X-PhiNet performs robustly for different weight decays on CIFAR-5m.

	Accuracy by Linear Probing (w.r.t. weight decay)			
	0.0001	5e-05	2e-05	1e-05
BYOL	67.88 _{0.58}	75.71 _{0.34}	81.05 _{0.04}	80.70 _{0.84}
SimSiam	77.69 _{0.67}	75.02 _{5.92}	76.87 _{3.13}	77.71 _{1.97}
PhiNet	76.43 _{2.12}	77.57 _{2.01}	77.64 _{1.44}	77.74 _{0.79}
RM-SimSiam	74.24 _{0.56}	77.52 _{0.88}	82.09 _{0.22}	79.38 _{0.38}
X-PhiNet with Aug (mse)	65.96 _{16.24}	83.21 _{0.15}	86.45 _{0.25}	85.17 _{0.54}
X-PhiNet with Aug (cos)	84.31 _{0.29}	86.40 _{0.31}	86.96 _{0.17}	84.85 _{0.99}
X-PhiNet (mse)	69.02 _{14.25}	84.24 _{0.37}	87.30 _{0.13}	85.11 _{0.17}
X-PhiNet (cos)	85.80 _{0.34}	87.29 _{0.22}	87.46 _{0.19}	85.03 _{0.19}
X-PhiNet (cos, $g = I$)	84.72 _{0.16}	86.41 _{0.08}	86.71 _{0.27}	83.83 _{0.36}

Table 4: X-PhiNet performs robustly well for one epoch training.

	Accuracy by Linear Probing (w.r.t. weight decay)			
	0.0001	5e-05	2e-05	1e-05
BYOL	63.86 _{0.77}	59.80 _{0.50}	58.04 _{0.52}	57.65 _{0.38}
SimSiam	68.50 _{0.19}	69.65 _{0.29}	69.41 _{1.06}	69.60 _{0.97}
PhiNet	66.27 _{1.60}	64.26 _{1.82}	64.34 _{1.13}	62.68 _{1.60}
RM-SimSiam	62.90 _{1.11}	63.30 _{1.86}	63.45 _{0.92}	63.05 _{1.76}
X-PhiNet (mse)	74.25 _{0.80}	72.65 _{0.85}	71.20 _{0.48}	71.95 _{0.65}
X-PhiNet (cos)	74.76 _{0.52}	72.89 _{0.66}	72.10 _{0.15}	71.93 _{0.51}

```

11 Lcos = D(p1, z2)/2 + D(p2, z1)/2 # loss
12 Lcor = mse_loss(y1, z0)/2 + mse_loss(y2, z0)/2
13 L = Lcos + Lcor
14 L.backward() # back-propagate
15 update(f, h) # SGD update
16 f_long = beta * f_long + (1-beta) * f # EMA for projection
17
18 def D(p, z): # negative cosine similarity
19     z = z.detach() # stop gradient
20     p = normalize(p, dim=1) # l2-normalize
21     z = normalize(z, dim=1) # l2-normalize
22     return -(p*z).sum(dim=1).mean()

```

Listing 2: X-PhiNet Pseudocode (PyTorch-like)

D ADDITIONAL EXPERIMENTS ON ONLINE AND CONTINUAL LEARNING

D.1 CIFAR-5M (ONLINE LEARNING)

Table 3 shows the accuracy of linear probing for different weight decay values in CIFAR-5m. X-PhiNet consistently demonstrates high performance. In Table 4, we trained for only one epoch on CIFAR-5m. Also in this case, X-PhiNet performs better than SimSiam. Additionally, X-PhiNet with $g = I$ achieves lower accuracy compared to the standard X-PhiNet (cos). This represents a major difference from the results shown in Figure 6, where both achieved similar accuracy with batch size = 128.

In Table 3, "X-PhiNet with Aug (mse)" and "X-PhiNet with Aug (cos)" represent cases where data augmentation is also applied to the input x of f_{long} . In this scenario, all inputs to the model are augmented. While X-PhiNet with Aug outperforms other baselines such as SimSiam and RM-SimSiam, its performance is still inferior to our standard X-PhiNet. This discrepancy might be attributed to an imbalance in regularization strength, although the exact reason remains unclear and is left for future work.

D.2 CONTINUAL LEARNING

Epochs per task. In Table 2, we trained on each task for one epoch. However, in Madaan et al. (2022), 200 epochs were trained for each task on Split CIFAR10, and the number of iterations differs from this case. The effect of early stopping may be apparent when the number of iterations is different. Thus, we trained on each task for two epochs to match the number of iterations. The result is shown in Table 5. The performance of X-PhiNet is still high even when the number of epochs per task is set to 2 epochs.

Table 5: **X-PhiNet shows higher accuracy when the number of epochs per task is increased.** We trained X-PhiNet on Split CIFAR-5m. Unlike Table 2, this table presents results obtained from training 2 epochs for each task.

		BYOL	SimSiam	Barlow Twins	PhiNet	RM- SimSiam	X-PhiNet (MSE)	X-PhiNet (Cos)
Split C-5m (2epoch)	Acc	91.36 _{0.25}	92.25 _{0.10}	90.73 _{0.28}	92.22 _{0.09}	90.11 _{0.34}	92.33 _{0.15}	92.83 _{0.06}
	Fg	4.10 _{0.25}	3.88 _{0.33}	5.25 _{0.67}	4.01 _{0.26}	6.12 _{0.50}	3.79 _{0.47}	3.71 _{0.14}

Replay with Mixup. Replay is one of the most promising methods for improving the performance of continual learning though it requires additional memory costs (Hsu et al., 2018; Van de Ven et al., 2020; Madaan et al., 2022; Lin et al., 2022). We thus examined the performance of our method in combination with the mixup-based replay method proposed in (Madaan et al., 2022). Table 6 shows that when only one epoch is trained for each task, X-PhiNet shows considerably higher accuracy than the other methods. On the other hand, when we train two epochs for each task, the accuracy of other methods such as BYOL, BarlowTwins and RM-SimSiam also increases, showing an accuracy comparable to that of X-PhiNet.

Table 6: **X-PhiNet performs higher or comparable results for Split-CIFAR5m even with Mixup.** We trained X-PhiNet on Split CIFAR-5m with replay methods.

		BYOL	SimSiam	Barlow Twins	PhiNet	RM- SimSiam	X-PhiNet (MSE)	X-PhiNet (Cos)
Split C-5m (1epoch)	Acc	90.18 _{0.65}	91.51 _{0.42}	90.74 _{0.63}	91.66 _{0.21}	91.92 _{0.17}	91.78 _{0.29}	92.43_{0.14}
	Fg	0.36 _{3.07}	-1.70 _{0.09}	-0.97 _{3.05}	-2.21 _{0.54}	-1.14 _{0.27}	-0.73 _{0.07}	-0.69 _{0.65}
Split C-5m (2epoch)	Acc	92.36_{0.03}	91.77 _{0.01}	92.36_{0.70}	91.00 _{1.78}	92.48_{0.12}	92.12_{0.28}	92.26_{0.35}
	Fg	0.64 _{0.01}	2.24 _{0.01}	-0.07 _{0.64}	1.97 _{1.41}	1.05 _{0.31}	2.02 _{0.71}	1.82 _{0.71}

Split CIFAR10 and Split CIFAR100. Up to this point, we have experimented with continual learning using the CIFAR-5m-based dataset. Now, we test on the standard benchmarks, Split CIFAR10 and Split CIFAR100. Table 7 shows that in both Split CIFAR10 and Split CIFAR100, X-PhiNet outperforms SimSiam. However, PhiNet sometimes shows higher accuracy than X-PhiNet. Note that PhiNet is a special case of X-PhiNet, and we have set the momentum of X-PhiNet to 0.99 in this study. If we carefully select the momentum value, X-PhiNet’s performance might improve, surpassing PhiNet. When using mixup for replay, X-PhiNet shows significantly higher accuracy compared to other methods.

Effect of exponential moving average X-PhiNet has an additional hyperparameter, the exponential moving average. We set $\beta = 0.99$ in all the experiments in this paper. As shown in Table 8, in tasks such as continual learning, where it is important to apply a strong exponential moving average, accuracy increases as β increases and then decreases again from a certain point.

E ADDITIONAL EXPERIMENTS ON THE ROBUSTNESS OF PHINET

E.1 ADDITIONAL ABLATION STUDY WITH CIFAR10

Usage of original input: We first investigate what is the best way to input the original signal to the model. To this end, we first replace one of the augmented signals in SimSiam as an original

Table 7: **X-PhiNet produces good results when memorization is important.** We trained X-PhiNet on Split CIFAR10. Acc is the average of the final Acc (higher is better), and Fg is Forgetting (smaller is better).

		SimSiam	RM-SimSiam	PhiNet (MSE)	X-PhiNet ($g = I$, MSE)	X-PhiNet (MSE)	X-PhiNet (Cos)
Split C10 (FineTune)	Acc	91.05 _{0.29}	89.35 _{0.08}	91.25_{0.09}	90.65 _{0.43}	90.90 _{0.50}	90.97_{0.49}
	Fg	5.31 _{0.65}	3.70 _{0.16}	4.86 _{0.50}	1.02 _{0.31}	5.72 _{0.82}	3.95 _{0.37}
Split C100 (FineTune)	Acc	77.93 _{0.64}	78.19 _{0.41}	78.50_{0.25}	78.31_{0.16}	77.50 _{0.04}	77.44 _{0.28}
	Fg	7.06 _{1.00}	-0.57 _{0.98}	6.51 _{0.31}	5.49 _{1.27}	8.46 _{0.21}	4.46 _{0.34}
Split C10 (Mixup)	Acc	90.68 _{0.89}	91.14 _{0.84}	89.89 _{0.69}	90.69 _{0.21}	90.49 _{0.32}	91.56_{0.12}
	Fg	0.85 _{0.16}	1.08 _{0.60}	1.04 _{0.22}	-2.11 _{1.61}	1.80 _{0.09}	1.36 _{0.15}
Split C100 (Mixup)	Acc	81.77 _{0.14}	82.47 _{0.70}	80.76 _{0.18}	82.16 _{0.76}	83.32 _{0.04}	83.88_{0.26}
	Fg	1.23 _{0.81}	-1.35 _{0.05}	1.18 _{1.27}	-1.23 _{2.11}	1.28 _{0.34}	-0.07 _{0.30}

Table 8: **X-PhiNet produces good results when memorization is important.** We trained X-PhiNet on Split CIFAR100 with mixup with different exponential moving average value.

		EMA β					
		0.999	0.997	0.99	0.97	0.9	0.7
Split C100 (Mixup)	Acc	83.19 _{0.22}	82.62 _{0.06}	83.88 _{0.26}	82.72 _{0.98}	81.76 _{0.64}	81.73 _{0.06}
	Fg	-0.22 _{0.34}	0.33 _{0.06}	-0.07 _{0.30}	0.73 _{0.64}	2.15 _{0.11}	1.82 _{0.69}

input. Then, we found that comparing the augmented images and original input significantly degrades the model performance. This indicates that the original SimSiam performs pretty well even if we do not use the original inputs, and naively adding additional input hurts the model performance significantly. In contrast, the PhiNet with MSE loss and StopGradient compares favorably with the original SimSiam model.

StopGradient-2: We analysed the impact of the StopGradient-2 technique, as shown in the table. The StopGradient operator effectively prevents mode collapse. Interestingly, while the StopGradient operator is not essential for avoiding mode collapse, models without it perform worse compared to those with it. Thus, the StopGradient operator contributes to improved stability when using the MSE loss. On the other hand, mode collapse still occurs with the negative cosine loss function.

E.2 COMPARISON OF FAST LEARNER WITH SLOW LEARNER

We can observe that EMA plays as a slow learner through experiments. In Figure 9, we conducted continual learning experiments, where linear probing of the EMA encoder (dashed lines) performs consistently worse than the encoder without EMA (solid lines). This indicates that EMA does not quickly adapt to the most recent samples and learns more stable features as a slow learner. In this sense, we believe it is natural to think of slow learning as serving as a regularization for past samples, similar to other continual learning techniques such as elastic weight consolidation.

E.3 ADDITIONAL ABLATION STUDY WITH IMAGENET

Table 10 shows the ablation of additional predictors in ImageNet. In this case, we used a higher weight decay of $1e-3$. $g = h$ has a lower accuracy than other methods, which is consistent with the results in CIFAR10.

E.4 FOR DIFFERENT DATASETS AND EVALUATION METRICS WITH DIFFERENT BATCH SIZES

Table 11 and Table 12 present the CIFAR10 experiment results with varying batch sizes and weight decay. PhiNet shows equal or better performance than SimSiam across different batch sizes. The evaluation trends from KNN classification and linear probing are also consistent. It is also a consistent result that training on CIFAR10 performs poorly when cosine loss is used as the cortex loss function. Table 13 shows the results for STL10, where PhiNet performs comparably to SimSiam. However, as

Table 9: Ablation study for PhiNet using CIFAR10 data. We use SGD with momentum as an optimiser and set the base learning rate as 0.03 and run 800 epochs. We evaluated the performance using KNN classification with $K = 200$. See Table 19 for further details.

Method	Sim-2	SG-2	Pred-2	Acc (w.r.t. weight decay)		
				0.0	0.0005	0.001
SimSiam	–	–	–	74.12 _{0.39}	90.39 _{0.10}	90.98 _{0.02}
SimSiam (Orig-In)	–	–	–	72.82 _{0.18}	76.67 _{1.13}	69.03 _{11.37}
PhiNet	MSE	✓	g	77.77 _{1.13}	90.77 _{0.22}	91.38 _{0.19}
	MSE	✓	$g = \mathbf{I}$	77.63 _{0.11}	91.01 _{0.12}	91.50 _{0.07}
	MSE		g	62.80 _{0.47}	91.40 _{0.23}	89.01 _{0.55}
	MSE	✓	h	74.87 _{0.58}	91.23 _{0.12}	91.18 _{0.34}
	Cos	✓	g	80.06 _{0.47}	87.73 _{0.26}	88.27 _{0.24}
	Cos	✓	$g = \mathbf{I}$	75.34 _{3.27}	87.38 _{0.17}	87.90 _{0.10}
	Cos		g	27.57 _{4.41}	9.98 _{0.00}	9.98 _{0.00}
	Cos	✓	h	75.99 _{0.28}	85.97 _{0.23}	85.04 _{0.11}

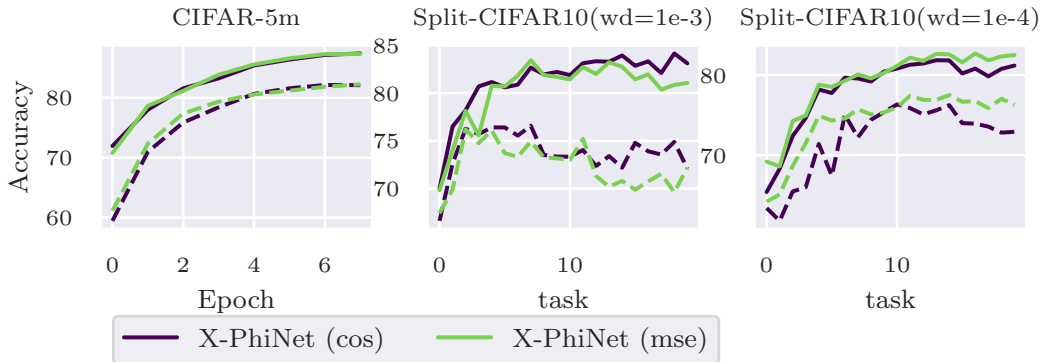


Figure 9: **The accuracy for slow weight is lower than the accuracy for fast weight.** The solid line represents the accuracy of the fast weights (the encoder without momentum), while the dashed line represents the accuracy of the slow weights (the encoder with momentum). Note that the accuracy for split-CIFAR10 represents the average accuracy.

Table 10: $g = h$ **has a low accuracy on ImageNet.** We trained the models for 100 epochs and then validated them on the test sets using linear probing on the head. Unlike Table 1, we train linear probing for 40 epochs to save computational costs.

	Model			
	SimSiam	PhiNet (MSE)	PhiNet ($g = I$)	PhiNet ($g = h$)
Linear Probing Acc	66.35	66.64	66.47	55.12

illustrated in Figure 7, PhiNet demonstrates better convergence in the early learning stages. In the early stages of training, when SimSiam is not stable, the cosine loss is smaller than that of the PhiNet, while as the training continues, the cosine loss increases again, in agreement with the PhiNet. This suggests that something close to mode collapse occurs in the early stages of SimSiam training, while PhiNet may suppress this collapse.

E.5 PERFORMANCE ON TRANSFER LEARNING

We conducted experiments for transfer learning using object detection on VOC. In Table 14, following the original SimSiam paper, we conducted pre-training experiments with two different settings for learning rate and weight decay. This table demonstrates that our X-PhiNet produces comparable performance to MoCo across various tasks. Furthermore, as shown in Table 15, we can find that

Table 11: **PhiNet shows equal or better performance than SimSiam.** Sensitivity analysis for PhiNet using CIFAR10 data. We use SGD with momentum as an optimiser, set the base learning rate as 0.03, and run 800 epochs. We evaluated the performance using KNN classification with $K = 200$.

weight decay		Acc (w.r.t. batch size)			
		128	256	512	1024
0.0001	SimSiam	86.35 _{2.28}	88.05 _{0.66}	88.34 _{0.28}	85.96 _{2.93}
	PhiNet	88.90_{0.23}	88.92_{0.10}	88.93_{0.33}	88.91_{0.07}
	X-PhiNet (MSE)	88.67_{0.39}	88.67_{0.23}	88.71_{0.20}	88.44 _{0.32}
	X-PhiNet (Cos)	82.57 _{5.67}	79.31 _{10.65}	72.35 _{19.83}	84.79 _{0.56}
0.0005	SimSiam	90.15_{0.15}	90.36_{0.15}	90.39 _{0.10}	90.89 _{0.08}
	PhiNet	90.40_{0.16}	90.48_{0.34}	90.57_{0.15}	91.15_{0.04}
	X-PhiNet (MSE)	90.05_{0.29}	90.13 _{0.02}	90.39 _{0.12}	90.70 _{0.10}
	X-PhiNet (Cos)	86.35 _{0.23}	86.36 _{0.16}	86.42 _{0.11}	86.65 _{0.65}
0.001	SimSiam	91.23_{0.11}	91.30 _{0.05}	90.98 _{0.02}	76.68 _{12.83}
	PhiNet	91.23_{0.07}	91.44_{0.08}	91.50_{0.03}	73.97 _{7.04}
	X-PhiNet (MSE)	91.04 _{0.13}	91.09 _{0.14}	91.11 _{0.13}	90.08_{0.37}
	X-PhiNet (Cos)	86.54 _{0.18}	86.33 _{1.04}	86.79 _{0.75}	87.47 _{1.16}

Table 12: **Linear probing shows similar trends to knn classification.** Sensitivity analysis for PhiNet using CIFAR10 data. We use SGD with momentum as an optimiser, set the base learning rate as 0.03, and run 800 epochs. We evaluated the performance using linear probing on the head.

weight decay		Acc (w.r.t. batch size)			
		128	256	512	1024
0.0001	SimSiam	88.64 _{1.73}	89.44 _{0.35}	89.39 _{0.49}	88.01 _{1.80}
	PhiNet	90.27 _{0.13}	89.83 _{0.35}	89.79 _{0.24}	89.70 _{0.21}
	X-PhiNet (MSE)	88.67 _{0.39}	88.67 _{0.23}	88.71 _{0.20}	88.44 _{0.32}
	X-PhiNet (Cos)	83.29 _{4.10}	84.15 _{0.46}	72.35 _{19.84}	83.23 _{2.52}
0.0005	SimSiam	90.39 _{0.07}	90.68 _{0.05}	91.15 _{0.12}	91.65 _{0.06}
	PhiNet	90.68 _{0.10}	90.87 _{0.34}	91.11 _{0.08}	91.93 _{0.13}
	X-PhiNet (MSE)	90.05 _{0.29}	90.13 _{0.02}	90.39 _{0.12}	90.70 _{0.10}
	X-PhiNet (Cos)	86.52 _{0.12}	86.47 _{0.11}	86.45 _{0.16}	87.00 _{0.18}
0.001	SimSiam	92.09 _{0.22}	92.36 _{0.11}	92.46 _{0.29}	77.71 _{12.73}
	PhiNet	92.18 _{0.06}	92.44 _{0.21}	92.63 _{0.08}	75.18 _{6.69}
	X-PhiNet (MSE)	91.04 _{0.13}	91.09 _{0.14}	91.11 _{0.13}	90.08 _{0.37}
	X-PhiNet (Cos)	87.04 _{0.33}	87.08 _{0.20}	87.24 _{0.19}	88.15 _{0.09}

Table 13: **SimSham and PhiNet show comparable performance.** Sensitivity analysis for PhiNet using STL10 data. We use SGD with momentum as an optimiser, set the base learning rate as 0.03, and run 800 epochs. We evaluated the performance using linear probing on the head.

weight decay		Acc (w.r.t. batch size)			
		128	256	512	1024
0.0001	SimSiam	85.97 _{2.46}	87.26 _{0.26}	87.53 _{0.20}	87.17 _{0.05}
	PhiNet	84.28 _{0.41}	87.32 _{0.16}	87.22 _{0.12}	87.01 _{0.28}
0.0005	SimSiam	88.89 _{0.34}	89.23 _{0.11}	89.39 _{0.12}	88.57 _{0.40}
	PhiNet	89.33 _{0.13}	89.26 _{0.02}	89.36 _{0.27}	88.62 _{0.36}
0.001	SimSiam	89.54 _{0.05}	89.61 _{0.11}	89.37 _{0.03}	<i>nan_{nan}</i>
	PhiNet	89.52 _{0.06}	89.71 _{0.07}	89.28 _{0.23}	10.00 _{0.01}

Pretrained	VOC 07 detection			VOC 07+12 detection		
	AP ₅₀	AP	AP ₇₅	AP ₅₀	AP	AP ₇₅
MoCo v2	73.2	46.6	50.2	82.3	57.1	63.2
SimSiam (lr=0.05, wd=1e-4)	71.7	45.5	49.4	80.6	55.1	61.0
SimSiam (lr=0.5, wd=1e-5)	73.6	46.6	49.8	82.7	57.3	64.6
PhiNet (lr=0.05, wd=1e-4)	72.7	46.35	50.4	81.9	56.8	62.81
PhiNet (lr=0.5, wd=1e-5)	74.4	46.2	49.8	82.6	56.4	62.6
X-PhiNet (lr=0.05, wd=1e-4)	72.9	46.2	49.9	82.3	56.9	63.9
X-PhiNet (lr=0.5, wd=1e-5)	74.9	45.9	50.1	82.7	55.7	62.4

Table 14: **In transfer learning for object detection, X-PhiNet is comparable to MoCo (He et al., 2020) and SimSiam.** PhiNet is pre-trained by two training recipes similar to those in the SimSiam paper.

wd	2e-3	1e-3	5e-4	2.5e-4	1.25e-4	6.25e-5	3.125e-5	1.5625e-5
SimSiam	10.00	63.53	90.80	90.05	89.06	79.70	78.14	76.35
MoCo	87.47	88.11	87.76	87.01	85.79	83.45	81.34	80.44
PhiNet	10.00	78.33	91.19	90.35	89.34	86.25	83.21	81.60

Table 15: **PhiNet is robust to weight decay in transfer learning.** Performance comparison of SimSiam, MoCo, and PhiNet at different weight decay values.

PhiNet demonstrates a higher sensitivity to weight decay. Thus, it seems that our PhiNet can be extended to object detection without any modifications.

E.6 ON THE AUGMENTATION FOR x

We use the unaugmented view for the Sim-2 loss to simulate a “time difference” between different views, which is partially supported by the temporal prediction hypothesis. This architecture does slightly increase the performance. See Table 16 and Figure 10, where “with aug” performs slightly worse than our proposed architecture while robustness to weight decay is still higher than SimSiam.

E.7 COMPUTATIONAL COSTS

Table 17 shows the memory consumption when training on CIFAR10. There is little overhead for PhiNet and X-PhiNet over SimSiam, as the maximum memory consumption during training is not only related to weights, but also to gradients and activation state. In fact, the GPU consumption is highly dependent on batch size, indicating that the gradient and activation state, which are dependent on batch size, are dominant in this setting. Additionally, Table 18 includes a comparison of training times, demonstrating that PhiNet can be trained in time comparable to SimSiam.

	Accuracy by Linear Probing (w.r.t. weight decay)			
	0.0001	5e-05	2e-05	1e-05
SimSiam	77.69 _{0.67}	75.02 _{5.92}	76.87 _{3.13}	77.71 _{1.97}
X-PhiNet with Aug (mse)	65.96 _{16.24}	83.21 _{0.15}	86.45 _{0.25}	85.17 _{0.54}
X-PhiNet with Aug (cos)	84.31 _{0.29}	86.40 _{0.31}	86.96 _{0.17}	84.85 _{0.99}
X-PhiNet (mse)	69.02 _{14.25}	84.24 _{0.37}	87.30 _{0.13}	85.11 _{0.17}
X-PhiNet (cos)	85.80 _{0.34}	87.29 _{0.22}	87.46 _{0.19}	85.03 _{0.19}

Table 16: **Even when x is augmented, X-PhiNet performs better than SimSiam (CIFAR-5m).** We used the same setting as in Table 3 in the original paper. “with aug” performs data augmentation for x , which is not augmented in Table 3.

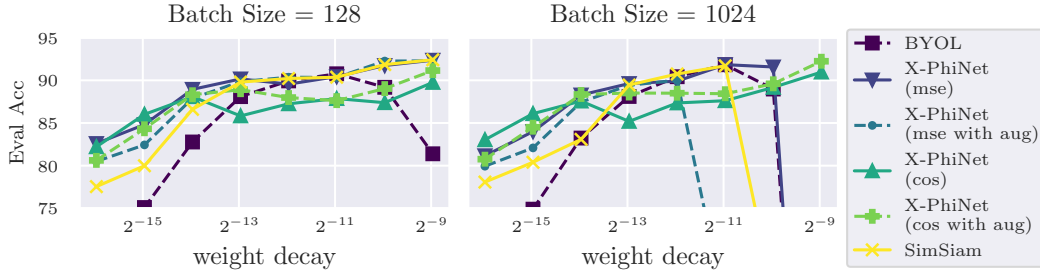


Figure 10: **Even when x is augmented, PhiNet is more robust to weight decay than SimSiam (CIFAR10).** We used the same setting as in Figure 6 in the original paper. “with aug” performs data augmentation for x , which is not data augmented in Figure 6.

Table 17: **Comparison of GPU memory costs.** This is a comparison of memory consumption when training on CIFAR10. We report batch sizes of 128 and 1024.

Batch Size	BYOL	SimSiam	Barlow Twins	PhiNet	RM-SimSiam	X-PhiNet (MSE)	X-PhiNet (Cos)
BS=128	4.3 (GB)	3.26 (GB)	2.79 (GB)	3.25 (GB)	4.06 (GB)	3.44 (GB)	3.44 (GB)
BS=1024	22.28 (GB)	17.10 (GB)	12.23 (GB)	17.18 (GB)	21.96 (GB)	17.11 (GB)	17.11 (GB)

E.8 STABLE RANK OF ADDITIONAL PREDICTOR LAYER

Figure 11 shows the rank for 2 linear layers in additional predictor blocks of PhiNet. We used stable rank in this figure and it is defined as $\text{srnk}(M) = \|M\|_F^2 / \|M\|^2$, which is the lower rank of the standard rank and is more stable to the small eigenvalues of M . According to this figure, the rank of the additional layer remains large when the weight decay is small, suggesting that the additional layer may play a more important role in learning when the weight decay is small.

F EXPERIMENTAL SETTINGS

F.1 SETTINGS FOR TRAINING WITH CIFAR10, CIFAR100 AND STL10

Table 19 shows the model and experimental setup for Figure 6, Table 9, Table 11, Table 12 and Table 13. Note that in the graph of sensitivity with respect to weight decay, we explored a wider range of values. For linear probing evaluation, we trained the head layer by SGD for 100 epochs. For both CIFAR10 and STL10, we used 50,000 samples for training and 10,000 samples for testing. We have implemented it based on code that is already publicly available¹.

F.2 SETTINGS FOR TRAINING ON IMAGENET

In our ImageNet (Russakovsky et al., 2015) experiments, we follow the formal implementation of SimSiam by Pytorch² (Chen and He, 2021). Table 20 shows the model and experimental setup for Table 1. For ImageNet, we used 1,281,167 samples for training and 100,000 samples for testing. We trained on the three seeds and obtained the mean and variance.

¹<https://github.com/PatrickHua/SimSiam>

²<https://github.com/facebookresearch/simsiam>

Table 18: **Comparison of different models with varying batch sizes.**

Batch Size	SimSiam	BYOL	Barlow-Twins	RM-SimSiam	PhiNet	X-PhiNet
BS=128	6.89 (h)	7.09 (h)	21.38 (h)	7.76 (h)	6.89 (h)	7.04 (h)
BS=1024	6.74 (h)	6.51 (h)	7.68 (h)	7.03 (h)	6.44 (h)	6.52 (h)

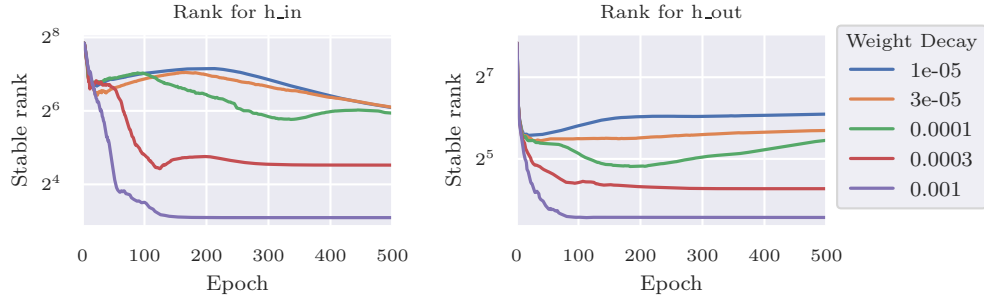


Figure 11: **The smaller the weight decay, the larger the rank of the additional predictor.** We trained PhiNet on CIFAR10 with SGD and evaluate the rank for layers in additional predictor blocks.

Table 19: **The experimental setups of Figure 6, Table 9, Table 11, Table 12 and Table 13.**

Learning	Optimiser	SGD
	Momentum	0.9
	Learning Rate	0.03
	Epochs	800
Encoder	Backbone	ResNet18_cifar_variant1
	Projector output dimension	2048
Predictor h	Latent dimension m	2048
	Hidden dimension h	512
	Activation function	ReLU
	Batch normalization	Yes
Predictor g	Latent dimension m	2048
	Hidden dimension g	512
	Activation function (Hidden)	ReLU
	Activation function (Output)	Tanh
	Batch normalization	Yes
Computational resource	GPUs	V100 or A100

Table 20: **The experimental setups of Table 1.**

Learning	Optimiser	SGD
	Momentum	0.9
	Learning Rate	0.05
	Epochs	100
Encoder	Backbone	ResNet50
	Projector output dimension	2048
Predictor h	Latent dimension m	2048
	Hidden dimension h	512
	Activation function	ReLU
	Batch normalization	Yes
Predictor g	Latent dimension m	2048
	Hidden dimension g	512
	Activation function (Hidden)	ReLU
	Activation function (Output)	Tanh
	Batch normalization	Yes
Computational resource	GPUs	4×V100

F.3 SETTINGS FOR TRAINING ON CIFAR-5M

CIFAR-5m (Nakkiran et al., 2021) is a dataset that is sometimes used as a vision dataset for online learning (Vyas et al., 2023; Sarnthein et al., 2023). We experimented with CIFAR-5m in a setting similar to online learning. Note that CIFAR-5m has 5m samples, but we chose to train CIFAR-5m for 8 epochs, as most of the SimSiam training on CIFAR10 involves training for 800 epochs. We experimented with three learning rates: $\{0.03, 0.01, 0.003\}$, and selected the one that yielded the best results. For Barlow Twins, the learning rate of 0.03 does not converge, so 0.003 is chosen instead. For all other methods, a learning rate of 0.03 is selected. The model architecture is the same as in CIFAR10.

F.4 SETTINGS FOR TRAINING ON SPLIT CIFAR10, SPLIT CIFAR100 AND SPLIT-CIFAR-5M

As a benchmark for evaluating continual learning, we used split CIFAR10 and split CIFAR100 (Krizhevsky, 2009). Additionally, we created split cifar-5m, which is inspired by split-CIFAR10 but uses CIFAR-5m dataset. In split CIFAR10 and split CIFAR5m, we split CIFAR10 and CIFAR-5m into 5 tasks, each of which contains 2 classes. In split CIFAR10, we split CIFAR100 into 10 tasks, each of which contains 2 classes. The model architecture is the same as in CIFAR10. The implementation of continual learning is based on the official implementation of Madaan et al. (2022).

We evaluated the results using Average Accuracy and Average Forgetting. The average accuracy after the model has trained for T tasks is defined as:

$$A_T = \frac{1}{T} \sum_{i=1}^T a_{T,i}, \quad (10)$$

where $a_{t,i}$ is the validation accuracy on task i after the model finished task t . The average forgetting is defined as the difference between the maximum accuracy and the final accuracy of each task. Therefore, average forgetting after the model has trained for T tasks can be defined as:

$$F = \frac{1}{T-1} \sum_{i=1}^{T-1} \max_{t \in \{1, \dots, T-1\}} (a_{t,i} - a_{T,i}) \quad (11)$$