

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 PHINETS: BRAIN-INSPIRED NON-CONTRASTIVE LEARNING BASED ON TEMPORAL PREDICTION HY- POTHESIS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Predictive coding has been established as a promising neuroscientific theory to describe the mechanism of long-term memory residing in the retina. This theory hypothesises that cortex predicts sensory inputs at various levels of abstraction to minimise prediction errors and forms long-term memory. 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 to have 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.

## 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.

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 model 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. A widely used 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 (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 effectiveness 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 hippocampus 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 single and 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)), compared with 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 method 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 the model’s structural aspects. Currently, backpropagation-free predictive coding mechanisms for complex architectures like ResNet are in the

108 early stages of development, with most research limited to simple CNNs. Conversely, non-contrastive  
 109 methods like SimSiam require more advanced models than ResNet. Future research should explore  
 110 if the proposed structure can enable effective learning with backpropagation-free predictive coding.  
 111 Another key difference between PhiNet and brains is the presence of recurrent structures. However,  
 112 in this PhiNet, only one time step is considered, so it is possible that the recurrent structure required  
 113 to predict time series data was not necessary. Making the data into time series data and adding a  
 114 recurrent structure to the model remains as future work.

## 116 2 RELATED WORK

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

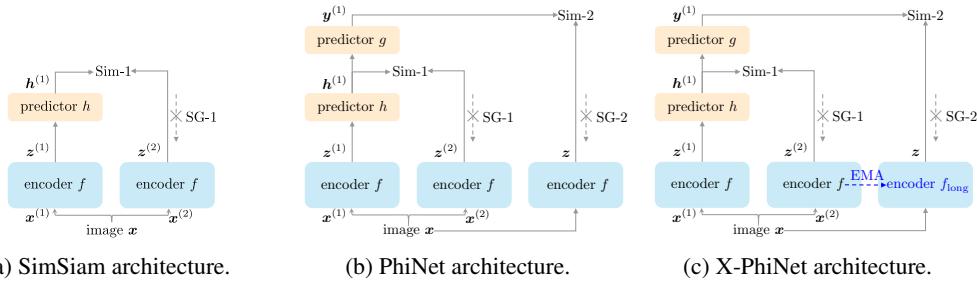
130 **Self-supervised learning.** Current mainstream approaches to self-supervised learning (SSL) often  
 131 rely on cross-view prediction frameworks (Becker and Hinton, 1992), with contrastive learning  
 132 emerging as a prominent SSL paradigm. In contrastive learning like SimCLR (Chen et al., 2020a),  
 133 models contrast positive (similar) and negative (dissimilar) samples to learn data representations.  
 134 One limitation of SimCLR is its empirical reliance on gigantic negative samples, which has been  
 135 theoretically articulated (Bao et al., 2022; Awasthi et al., 2022). To address this issue, recent research  
 136 has focused on approaches free from negative sampling (Grill et al., 2020; Caron et al., 2020; 2021).  
 137 For instance, BYOL (Grill et al., 2020) trains representations by aligning online and target networks,  
 138 where the target network is created by maintaining a moving average of the online network parameters.  
 139 SimSiam (Chen and He, 2021) utilises a Siamese network to align two augmented views of an input  
 140 by fixing one of the networks with StopGradient. While the lack of negative samples may easily  
 141 yield collapsed representations, namely, constant representation, Tian et al. (2021) analysed the  
 142 BYOL/SimSiam dynamics with a two-layer network and found that complete collapse is prevented  
 143 unless weight decay is excessively strong. We partially leverage their analysis framework to explain  
 144 the mechanism of our PhiNet. In recent years, many studies have been leveraging SimSiam for  
 145 continual learning (Smith et al., 2021; Madaan et al., 2022) and reinforcement learning (Tang et al.,  
 146 2023). RM-SimSiam (Fu et al., 2024) and CaSSLe (Fini et al., 2022) enhances the performance of  
 147 continual learning by incorporating a memory block into SimSiam, while its architecture has not  
 148 been neuroscientifically grounded.

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

## 154 3 PHINETS ( $\Phi$ -NETS)

156 In this paper, we propose PhiNets, which are non-contrastive methods based on CLS theory (Mc-  
 157 Clelland et al., 1995) and the temporal prediction hypothesis (Chen et al., 2024). Chen et al. (2024,  
 158 Figure 1) provides the hippocampal model, where the entorhinal cortex (EC) serves as an input signal  
 159 layer, the CA3 region serves as the predictor, and the CA1 region measures the prediction error.  
 160 The CA3 region receives an input signal from the EC and recurrently forecasts future signals. The  
 161 prediction output of CA3 is propagated to the CA1 region, which computes the discrepancy between  
 162 the CA3 prediction and the EC input and refines the internal model stored in CA3. Compensating the

162  
163  
164  
165  
166  
167  
168  
169  
170  
171



172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

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** ( $\Phi$ -Net) because the shape of the architecture resembles the Greek letter Phi ( $\Phi$ ).

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  $f_{\text{long}}$  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, 1b, and 1c depict the architecture of **SimSiam** and **PhiNets**, 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 the two augmented images as an input and feedback signals with time difference (thanks to the StopGradient operation), expanding the applicability of hippocampal models to standard vision tasks.

### 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  $\mathbf{x}_i \in \mathbb{R}^d$ . We model that the hippocampus model has two augmented signals from the original input as  $\mathbf{x}_i^{(1)} \in \mathbb{R}^d$  and  $\mathbf{x}_i^{(2)} \in \mathbb{R}^d$ , in addition to the original input  $\mathbf{x}_i$ . Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$  denote the encoder. Then, the cortical representation in the EC is given as follows:

$$\mathbf{z}_i^{(1)} = f(\mathbf{x}_i^{(1)}), \quad \mathbf{z}_i^{(2)} = f(\mathbf{x}_i^{(2)}), \quad \text{and} \quad \mathbf{z}_i = f(\mathbf{x}_i).$$

We regard each corresponding to the layers II, III, and V of the EC in Chen et al. (2024, Figure 1). Thus, for self-supervised training, we have the size- $n$  triplet dataset  $\mathcal{D} = \{(\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, \mathbf{x}_i)\}_{i=1}^n$ . The

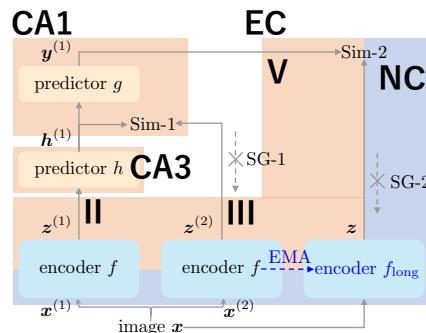


Figure 2: The interpretation as a hippocampal model. NC stands for NeoCortex.

216 hippocampal learning can be characterized as a learning problem of the encoder  $f$  from the training  
 217 dataset  $\mathcal{D}$ .  
 218

219 **Modelling of CA3 region.** The CA3 region is responsible for predicting future signals:  
 220

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

222 where  $h : \mathbb{R}^m \rightarrow \mathbb{R}^m$  is the predictor network. We implement the predictor with a two-layer neural  
 223 network with the ReLU activation and batch normalization.  
 224

225 **Modelling of CA1 region.** CA1 measures the difference between the predicted signal and its future  
 226 signal. In this paper, we model CA1 by a mixture of a loss function and a predictor, while Chen et al.  
 227 (2024) uses only MSE loss for modelling CA1. We use the symmetric negative cosine loss function  
 228 to measure the temporally distant signal  $\mathbf{z}_i^{(2)}$  (layer III of EC) and the predicted representation from  
 229 CA3  $\mathbf{h}_i^{(1)} = h(\mathbf{z}_i^{(1)})$ :  
 230

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

234 where  $\boldsymbol{\theta}$  represent the entire model parameter and  $\bar{\mathbf{z}}_i^{(1)} := \text{SG}(\mathbf{z}_i^{(1)}) \in \mathbb{R}^m$  is a latent variable with  
 235 StopGradient, in which the gradient update shall not be executed.  
 236

Remark that StopGradient yields a “time difference”, for which we can interpret PhiNet as a hippocampus model. Let us look closely at Sim-1 in Figure 1c. We let  $f_t$  denote the encoder  $f$  at time  $t$  **where t denotes the update times**. 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 adapted with StopGradient, it can be written as  $\mathbf{z}^{(2)} = \text{SG}(f(\mathbf{x})) = f_{t-1}(\mathbf{x}(t))$ . 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 X-PHINET: 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:

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

This is regarded as slow learning. Finally, the whole objective function of PhiNet is given as

$$L(\boldsymbol{\theta}) = \underbrace{L_{\text{Cos}}(\boldsymbol{\theta})}_{\text{Hippocampus loss}} + \underbrace{L_{\text{NC}}(\boldsymbol{\theta})}_{\text{Neocortex loss}}.$$

We then minimise  $L(\boldsymbol{\theta})$  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 negative cosine similarity.

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

$$274 \quad z_i = f_{\text{long}}(\mathbf{x}_i).$$

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

$$277 \quad \xi_{\text{long}} \leftarrow \beta \xi_{\text{long}} + (1 - \beta) \xi,$$

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

#### 283 4 WHAT WE BENEFIT FROM ADDITIONAL CA1 PREDICTOR: DYNAMICS 284 PERSPECTIVE

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

288 **Analysis model.** Let us specify the analysis model, following Tian et al. (2021). The  $d$ -dimensional  
 289 input is sampled from the isotropic normal  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and augmented by the isotropic normal  
 290  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \sim \mathcal{N}(\mathbf{x}, \sigma^2 \mathbf{I})$ , where  $\sigma^2$  indicates the strength of data augmentation. The encoder  $f$  and  
 291 predictors  $g$  and  $h$  are modelled by linear networks without bias:  $f(\mathbf{x}) := \mathbf{W}_f \mathbf{x}$ ,  $g(\mathbf{h}) := \mathbf{W}_g \mathbf{h}$ ,  
 292 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  
 293 latents  $z^{(1)}, h^{(1)} \in \mathbb{R}^h$  into  $h^{(1)}, y^{(1)} \in \mathbb{R}^m$  with the same dimension  $m$  to predict the other noisy  
 294 latent  $z^{(2)}$  and the noise-free latent  $z$ , respectively (see Figure 1b).

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

$$300 \quad \bar{L} := \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})} \left[ \|\mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x}^{(2)})\|^2 + \|\mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x})\|^2 \right].$$

301 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)  
 302 subsequently. The gradient flows are derived as follows (see Appendix B.1):

$$303 \quad \begin{aligned} \dot{\mathbf{W}}_f &= -\mathbf{W}_h^\top \{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g)\mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\}\mathbf{W}_f - \rho \mathbf{W}_f, \\ 304 \quad \dot{\mathbf{W}}_g &= -\{(1 + \sigma^2)\mathbf{W}_h - \mathbf{I}\}\mathbf{W}_f \mathbf{W}_f^\top \mathbf{W}_h^\top - \rho \mathbf{W}_g, \\ 305 \quad \dot{\mathbf{W}}_h &= -\{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g)\mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\}\mathbf{W}_f \mathbf{W}_f^\top - \rho \mathbf{W}_h. \end{aligned}$$

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

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

313 Under these assumptions,  $\Phi$ ,  $\mathbf{W}_g$ , and  $\mathbf{W}_h$  are simultaneously diagonalizable and can be written  
 314 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)  
 315 common orthogonal eigenvectors. Here,  $\Lambda_\Phi = \text{diag}[\phi_1, \dots, \phi_m]$ ,  $\Lambda_g = \text{diag}[\gamma_1, \dots, \gamma_m]$ , and

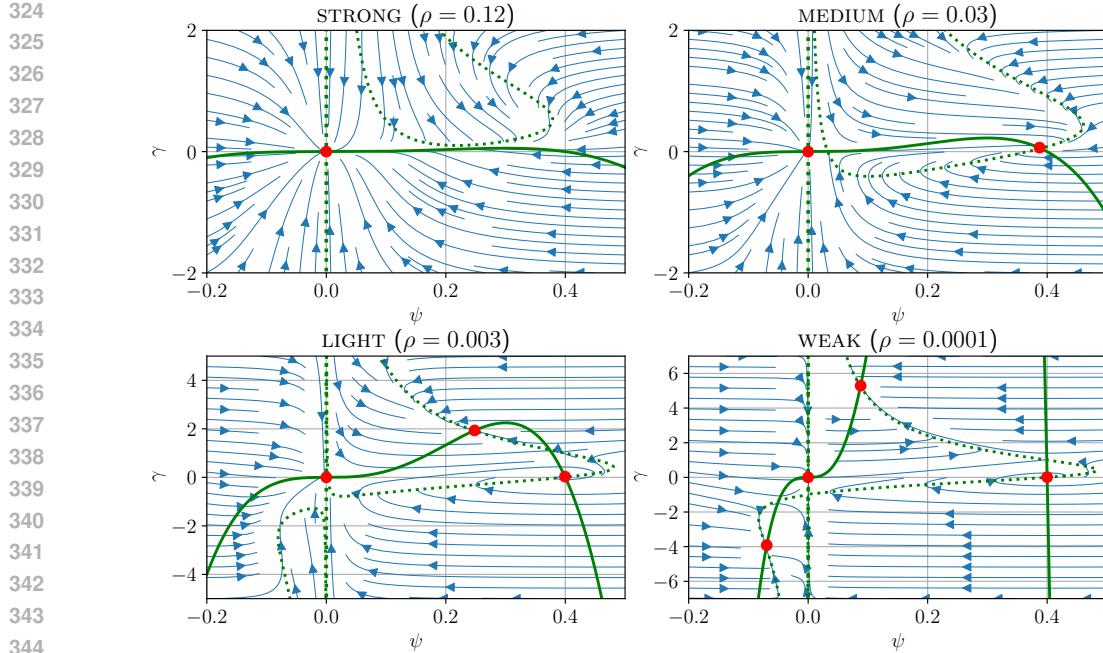


Figure 3: State space diagrams of PhiNet dynamics with different levels of weight decay: STRONG ( $\rho = 0.12$ ), MEDIUM ( $\rho = 0.03$ ), LIGHT ( $\rho = 0.003$ ), and WEAK ( $\rho = 0.0001$ ). The vector fields are numerically computed with  $\sigma^2 = 1.5$ . The state space bifurcates at the boundary 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 fall on to  $\phi(t) = \psi(t)^2$ , we can decouple the matrix dynamics into the eigenvalue dynamics of  $(\psi, \gamma)$  only (shown in Appendix B.3 and B.4):

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

From the  $(\psi, \gamma)$ -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  $\rho$  is too strong that the collapsed point  $(\psi, \gamma) = (0, 0)$  is a unique sink.
- MEDIUM: A new sink  $(\psi, \gamma)$  such that  $\psi \gg 0$  and  $\gamma \approx 0$  emerges. The number of sinks is two.
- LIGHT: Another non-trivial sink  $(\psi, \gamma)$  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  $\mathbf{W}_f x$  and  $\mathbf{W}_h z$ , respectively, and defined the gradient flow dynamics with the MSE loss  $\|\mathbf{W}_h \mathbf{W}_f x^{(1)} - \text{SG}(\mathbf{W}_f x^{(2)})\|^2$  (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  $\psi$ -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).

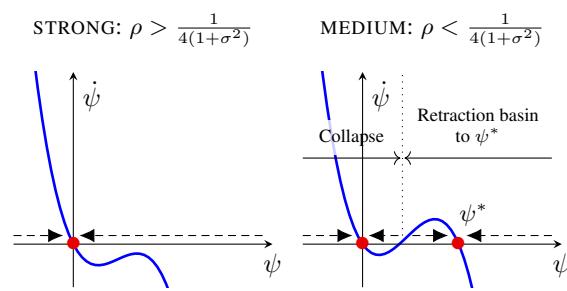


Figure 4: Illustration of SimSiam dynamics (2). Unlike the bivariate PhiNet dynamics shown in Figure 3, SimSiam dynamics is univariate, shown in the  $\psi$ -axis. The red dots are sinks.

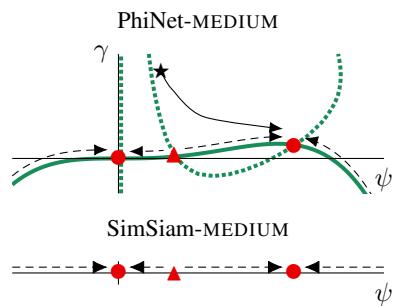


Figure 5: The SimSiam-MEDIUM flow is conjugate with the flow on the nullcline  $\psi = 0$  (green real line) in PhiNet-MEDIUM.

**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 <sub>0.02</sub>	2.63 <sub>0.18</sub>	60.82 <sub>1.57</sub>	44.51 <sub>33.64</sub>	68.17 <sub>0.18</sub>	67.12 <sub>0.13</sub>
PhiNet (MSE)	49.89 <sub>0.35</sub>	55.90 <sub>1.57</sub>	33.92 <sub>7.42</sub>	66.73 <sub>0.03</sub>	68.25 <sub>0.21</sub>	67.83 <sub>0.15</sub>

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  $\psi$ -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,  $\psi$  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  $\psi$  can avoid collapse:* In SimSiam-MEDIUM,  $\psi$  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  $q$  is beneficial.

**Remark 1.** The learning dynamics analysis in this section reveals that smaller weight decay  $\rho$  brings us benefits only regarding the stability of non-collapsed solutions. Indeed, we may benefit from larger  $\rho$  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  $\rho$  often yields good generalization in non-contrastive learning (Cabannes et al., 2023). Thus, smaller  $\rho$  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

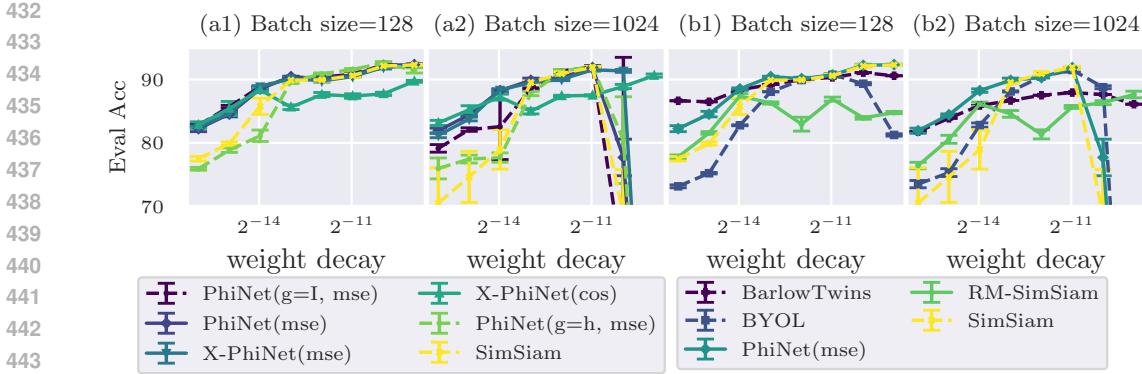


Figure 6: **PhiNet and X-PhiNet are robust against weight decay.** We compared PhiNet variants in (a1)-(a2) and compared the existing non-contrastive methods with PhiNet in (b1)-(b2) on CIFAR10. We evaluated the performance using linear probing. The loss function in brackets represents neocortex loss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

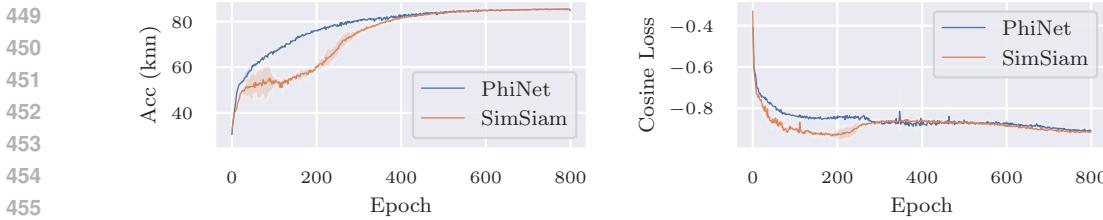


Figure 7: **PhiNet is stable in the early stages of learning.** We trained PhiNet and SimSiam with a batch size of 1024 and the weight decay of  $1e-4$  on STL10. SimSiam is unstable in the early stages of learning. This may be due to the cosine loss being too small in SimSiam.

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.

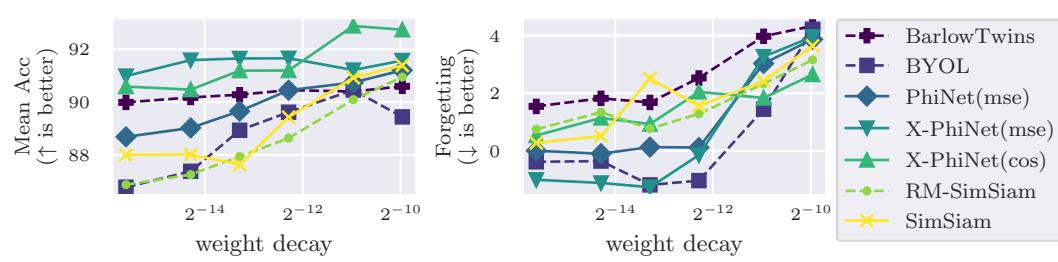
Overall, our sensitivity study on CIFAR10 revealed that PhiNets are robust to the choice of the weight decay parameter, which supports the importance of the CA1 predictor and Sim-2 loss. See Appendix E.1 for more detailed sensitivity studies. In addition, the results for different batch sizes and datasets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

486  
 487 Table 2: **X-PhiNet performs good results when memorization is important.** We trained PhiNets  
 488 on CIFAR-5m and Split CIFAR-5m. In Split CIFAR-5m, Acc is the average of the final accuracy  
 489 (higher is better), and Fg is Forgetting (smaller is better). We present the results for two different  
 490 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 <sub>0.04</sub>	77.71 <sub>1.97</sub>	85.32 <sub>0.10</sub>	76.74 <sub>1.82</sub>	82.09 <sub>0.22</sub>	87.30 <sub>0.13</sub>	<b>87.46<sub>0.14</sub></b>
Split C-5m (wd=5e-4)	Acc	90.44 <sub>0.28</sub>	90.84 <sub>0.31</sub>	90.30 <sub>0.17</sub>	90.69 <sub>0.11</sub>	90.04 <sub>0.11</sub>	91.02 <sub>0.36</sub>	<b>92.83<sub>0.12</sub></b>
	Fg	1.61 <sub>0.50</sub>	2.45 <sub>0.42</sub>	3.36 <sub>1.10</sub>	2.96 <sub>0.23</sub>	2.44 <sub>0.22</sub>	3.44 <sub>0.36</sub>	1.95 <sub>0.19</sub>
Split C-5m (wd=2e-5)	Acc	86.87 <sub>0.12</sub>	88.20 <sub>0.35</sub>	89.86 <sub>0.34</sub>	88.60 <sub>0.15</sub>	87.07 <sub>0.36</sub>	<b>90.90<sub>0.38</sub></b>	<b>90.72<sub>0.23</sub></b>
	Fg	-0.16 <sub>0.23</sub>	0.36 <sub>0.51</sub>	1.29 <sub>0.94</sub>	0.05 <sub>0.19</sub>	0.82 <sub>0.14</sub>	-1.03 <sub>0.41</sub>	0.43 <sub>0.17</sub>



500 Figure 8: **X-PhiNet is also robust to weight decay in continual learning.** We measured the mean  
 501 accuracy and forgetting at different weight decay on Split CIFAR-5m.  
 502

503 experiments using the CIFAR-5m dataset, which has six million synthetic CIFAR-10-like images  
 504 generated by DDPM generative model (Nakkiran et al., 2021). Instead of training CIFAR10 with 50k  
 505 samples for 800 epochs, we trained CIFAR-5m with 5m samples for 8 epochs. Although this is not  
 506 exactly online learning, it seems closer to online learning compared to CIFAR10 due to the restriction  
 507 on the training epochs. Table 2 shows that X-PhiNet has higher accuracy than SimSiam and PhiNet.  
 508 The superior performance of X-PhiNet compared to PhiNet suggests that long-term memory with  
 509 EMA is important in online learning. Sensitivity to weight decay and results for one-epoch online  
 510 learning are given in Appendix D.1.

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

## 6 CONCLUSION

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

540 REFERENCES  
541

- 542 P. Awasthi, N. Dikkala, and P. Kamath. Do more negative samples necessarily hurt in contrastive  
543 learning? In *ICML*, 2022.
- 544 H. Bao. Feature normalization prevents collapse of non-contrastive learning dynamics. *arXiv preprint*  
545 *arXiv:2309.16109*, 2023.
- 546 H. Bao, Y. Nagano, and K. Nozawa. On the surrogate gap between contrastive and supervised losses.  
547 In *ICML*, 2022.
- 548 S. Becker and G. E. Hinton. Self-organizing neural network that discovers surfaces in random-dot  
549 stereograms. *Nature*, 355(6356):161–163, 1992.
- 550 V. Cabannes, B. Kiani, R. Balestrieri, Y. LeCun, and A. Bietti. The SSL interplay: Augmentations,  
551 inductive bias, and generalization. In *ICML*, 2023.
- 552 M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of  
553 visual features by contrasting cluster assignments. *NeurIPS*, 2020.
- 554 M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging  
555 properties in self-supervised vision transformers. In *ICCV*, 2021.
- 556 T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of  
557 visual representations. In *ICML*, 2020a.
- 558 X. Chen and K. He. Exploring simple Siamese representation learning. In *CVPR*, 2021.
- 559 X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning.  
560 *arXiv preprint arXiv:2003.04297*, 2020b.
- 561 Y. Chen, H. Zhang, M. Cameron, and T. Sejnowski. Predictive sequence learning in the hippocampal  
562 formation. *Neuron*, 112:2645–2658, 2024.
- 563 A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning.  
564 In *AISTATS*, 2011.
- 565 E. Fini, V. G. T. Da Costa, X. Alameda-Pineda, E. Ricci, K. Alahari, and J. Mairal. Self-supervised  
566 models are continual learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
567 and Pattern Recognition*, 2022.
- 568 K. Friston. A theory of cortical responses. *Philosophical Transactions of the Royal Society B:  
569 Biological Sciences*, 360(1456):815–836, 2005.
- 570 F. Fu, Y. Gao, Z. Lu, H. Wu, and S. Zhao. Unsupervised continual learning of image representation  
571 via rememory-based simsiam. In *ICASSP*, 2024.
- 572 J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires,  
573 Z. Guo, M. Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised  
574 learning. *NeurIPS*, 2020.
- 575 M. S. Halvagal, A. Laborieux, and F. Zenke. Implicit variance regularization in non-contrastive SSL.  
576 *NeurIPS*, 2023.
- 577 K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual  
578 representation learning. In *CVPR*, 2020.
- 579 O. Henaff. Data-efficient image recognition with contrastive predictive coding. In *ICML*, 2020.
- 580 M. W. Hirsch, S. Smale, and R. L. Devaney. *Differential Equations, Dynamical Systems, and an  
581 Introduction to Chaos*. Academic Press, 2012.
- 582 Y.-C. Hsu, Y.-C. Liu, A. Ramasamy, and Z. Kira. Re-evaluating continual learning scenarios: A  
583 categorization and case for strong baselines. *arXiv preprint arXiv:1810.12488*, 2018.
- 584 A. Krizhevsky. Learning multiple layers of features from tiny images. 2009.

- 594 A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural  
595 networks. In *NIPS*, 2012.
- 596
- 597 D. Kumaran, D. Hassabis, and J. L. McClelland. What learning systems do intelligent agents need?  
598 complementary learning systems theory updated. *Trends in Cognitive Sciences*, 20:512–534, 2016.
- 599 Z. Lin, Y. Wang, and H. Lin. Continual contrastive learning for image classification. In *ICME*, 2022.
- 600
- 601 Z. Liu, E. S. Lubana, M. Ueda, and H. Tanaka. What shapes the loss landscape of self supervised  
602 learning? In *ICLR*, 2023.
- 603
- 604 D. Madaan, J. Yoon, Y. Li, Y. Liu, and S. J. Hwang. Representational continuity for unsupervised  
605 continual learning. In *ICLR*, 2022.
- 606
- 607 J. L. McClelland, B. L. McNaughton, and R. C. O'Reilly. Why there are complementary learning  
608 systems in the hippocampus and neocortex: insights from the successes and failures of connectionist  
609 models of learning and memory. *Psychological Review*, 102(3):419, 1995.
- 610
- 611 D. Mumford. On the computational architecture of the neocortex: II the role of cortico-cortical loops.  
*Biological Cybernetics*, 66(3):241–251, 1992.
- 612
- 613 P. Nakkiran, B. Neyshabur, and H. Sedghi. The deep bootstrap framework: Good online learners are  
614 good offline generalizers. In *ICLR*, 2021.
- 615
- 616 A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding.  
*arXiv preprint arXiv:1807.03748*, 2018.
- 617
- 618 Q. Pham, C. Liu, and S. Hoi. DualNet: Continual learning, fast and slow. *NeurIPS*, 2021.
- 619
- 620 Q. Pham, C. Liu, and S. C. Hoi. Continual learning, fast and slow. *IEEE Transactions on Pattern  
621 Analysis and Machine Intelligence*, 2023.
- 622
- 623 R. P. Rao and D. H. Ballard. Predictive coding in the visual cortex: a functional interpretation of  
624 some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1):79–87, 1999.
- 625
- 626 O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla,  
627 M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of  
628 Computer Vision*, 115:211–252, 2015.
- 629
- 630 F. Sarnthein, G. Bachmann, S. Anagnostidis, and T. Hofmann. Random teachers are good teachers.  
631 In *ICML*. PMLR, 2023.
- 632
- 633 J. Smith, C. Taylor, S. Baer, and C. Dovrolis. Unsupervised progressive learning and the STAM  
634 architecture. In *IJCAI*, 2021.
- 635
- 636 M. V. Srinivasan, S. B. Laughlin, and A. Dubs. Predictive coding: a fresh view of inhibition in the  
637 retina. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, 216(1205):  
638 427–459, 1982.
- 639
- 640 Y. Tang, Z. D. Guo, P. H. Richemond, B. A. Pires, Y. Chandak, R. Munos, M. Rowland, M. G. Azar,  
641 C. Le Lan, C. Lyle, et al. Understanding self-predictive learning for reinforcement learning. In  
642 *ICML*, 2023.
- 643
- 644 Y. Tian, X. Chen, and S. Ganguli. Understanding self-supervised learning dynamics without con-  
645 trastive pairs. In *ICML*, 2021.
- 646
- 647 G. M. Van de Ven, H. T. Siegelmann, and A. S. Tolias. Brain-inspired replay for continual learning  
648 with artificial neural networks. *Nature communications*, 2020.
- 649
- 650 N. Vyas, A. Atanasov, B. Bordelon, D. Morwani, S. Sainathan, and C. Pehlevan. Feature-learning  
651 networks are consistent across widths at realistic scales. In *NeurIPS*, 2023.
- 652
- 653 J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via  
654 redundancy reduction. In *ICML*, 2021.

## 648 A LIMITATIONS AND FUTURE WORK (EXTENDED VERSION)

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

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

## 667 B DETAILS OF LEARNING DYNAMICS ANALYSIS

668 In this section, we complement the missing details of learning dynamics analysis provided in Section 4.  
 669 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  
 670 continuous limit of gradient descent. This corresponds to considering the following gradient descent  
 671 in discrete updates and taking the limit as  $\eta \rightarrow 0$ .

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

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

### 679 B.1 DERIVATION OF MATRIX DYNAMICS

680 Recall the PhiNet loss function:

$$682 \bar{L} := \frac{1}{2} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}^{(1)}, \mathbf{x}^{(2)} | \mathbf{x}} \left[ \|\mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x}^{(2)})\|^2 + \|\mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x})\|^2 \right].$$

685 Let us derive its matrix gradient.

$$\begin{aligned} 686 \nabla_{\mathbf{W}_f} \bar{L} &= \frac{1}{2} \nabla_{\mathbf{W}_f} \mathbb{E} \left[ (\mathbf{x}^{(1)\top} \mathbf{W}_f^\top \mathbf{W}_h^\top - \text{SG}(\mathbf{x}^{(2)\top} \mathbf{W}_f^\top)) (\mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x}^{(2)})) \right. \\ 687 &\quad \left. + (\mathbf{x}^{(1)} \mathbf{W}_f^\top \mathbf{W}_h^\top \mathbf{W}_g^\top - \text{SG}(\mathbf{x}^\top \mathbf{W}_f^\top)) (\mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbf{x}^{(1)} - \text{SG}(\mathbf{W}_f \mathbf{x})) \right] \\ 688 &= \left\{ \mathbf{W}_h^\top \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] - \mathbf{W}_h^\top \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(2)} \mathbf{x}^{(1)\top}] \right\} \\ 689 &\quad + \left\{ \mathbf{W}_h^\top \mathbf{W}_g^\top \mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] - \mathbf{W}_h^\top \mathbf{W}_g^\top \mathbf{W}_f \mathbb{E}[\mathbf{x} \mathbf{x}^{(1)\top}] \right\} \\ 690 &= \mathbf{W}_h^\top \left\{ (\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] - \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(2)} \mathbf{x}^{(1)\top}] - \mathbf{W}_g^\top \mathbf{W}_f \mathbb{E}[\mathbf{x} \mathbf{x}^{(1)\top}] \right\} \\ 691 &= \mathbf{W}_h^\top \left\{ (1 + \sigma^2) (\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) \mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top) \right\} \mathbf{W}_f, \end{aligned}$$

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

$$698 \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}^{(1)} | \mathbf{x}} [\mathbf{x}^{(1)} \mathbf{x}^{(1)\top}] = \mathbb{E}_{\mathbf{x}} [\mathbf{x} \mathbf{x}^\top] + \sigma^2 \mathbf{I} = (1 + \sigma^2) \mathbf{I},$$

$$699 \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}^{(1)}, \mathbf{x}^{(2)} | \mathbf{x}} [\mathbf{x}^{(2)} \mathbf{x}^{(1)\top}] = \mathbb{E}_{\mathbf{x}} [\mathbf{x} \mathbf{x}^\top] = \mathbf{I},$$

$$700 \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}^{(1)} | \mathbf{x}} [\mathbf{x} \mathbf{x}^{(1)\top}] = \mathbb{E}_{\mathbf{x}} [\mathbf{x} \mathbf{x}^\top] = \mathbf{I}.$$

702 Similarly, we derive  $\nabla_{\mathbf{W}_g} \bar{L}$  and  $\nabla_{\mathbf{W}_h} \bar{L}$ .  
 703

$$\begin{aligned} 704 \nabla_{\mathbf{W}_g} \bar{L} &= \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)\top} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top \mathbf{W}_h^\top - \mathbf{W}_f \mathbb{E}[\mathbf{x} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top \mathbf{W}_h^\top \\ 705 &= \{(1 + \sigma^2) \mathbf{W}_h - \mathbf{I}\} \mathbf{W}_f \mathbf{W}_f^\top \mathbf{W}_h^\top. \\ 706 \end{aligned}$$

$$\begin{aligned} 708 \nabla_{\mathbf{W}_h} \bar{L} &= \left\{ \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)\top} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top - \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(2)\top} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top \right\} \\ 709 &\quad + \left\{ \mathbf{W}_g^\top \mathbf{W}_g \mathbf{W}_h \mathbf{W}_f \mathbb{E}[\mathbf{x}^{(1)\top} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top - \mathbf{W}_g^\top \mathbf{W}_f \mathbb{E}[\mathbf{x} \mathbf{x}^{(1)\top}] \mathbf{W}_f^\top \right\} \\ 710 &= \{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) \mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\} \mathbf{W}_f \mathbf{W}_f^\top. \\ 711 \\ 712 \\ 713 \end{aligned}$$

714 From these, we obtain the matrix dynamics.  
 715

## 716 B.2 EIGENSPACE ALIGNMENT

717 Our aim is to show that the three matrices  $\Phi$ ,  $\mathbf{W}_g$ , and  $\mathbf{W}_h$  share a common eigenspace, i.e.,  
 718 simultaneously diagonalizable, asymptotically in time  $t$ . Let  
 719

$$721 \mathbf{C}_1 := [\Phi, \mathbf{W}_g], \quad \mathbf{C}_2 := [\Phi, \mathbf{W}_h], \quad \text{and} \quad \mathbf{C}_3 := [\mathbf{W}_g, \mathbf{W}_h],$$

722 where  $[\mathbf{A}, \mathbf{B}] := \mathbf{AB} - \mathbf{BA}$  is the commutator (matrix). By noting that commutative matrices are  
 723 simultaneously diagonalizable, we show that the time-dependent commutators  $\mathbf{C}_1(t)$ ,  $\mathbf{C}_2(t)$ , and  
 724  $\mathbf{C}_3(t)$  asymptotically converges to  $\mathbf{O}$  as  $t \rightarrow \infty$ .  
 725

726 Hereafter, we assume the symmetry assumption (A1) on  $\mathbf{W}_g$  and  $\mathbf{W}_h$ , and heavily use the following  
 727 formulas on commutators implicitly:  
 728

- 729 •  $[\mathbf{A}, \mathbf{A}] = \mathbf{O}$ .
- 730 •  $[\mathbf{A}, \mathbf{B}] = -[\mathbf{B}, \mathbf{A}]$ .
- 731 •  $[\mathbf{A}, \mathbf{BC}] = [\mathbf{A}, \mathbf{B}]\mathbf{C} + \mathbf{B}[\mathbf{A}, \mathbf{C}]$ .
- 732 •  $[\mathbf{AB}, \mathbf{C}] = \mathbf{A}[\mathbf{B}, \mathbf{C}] + [\mathbf{A}, \mathbf{C}]\mathbf{B}$ .
- 733

734 First, compute  $\dot{\mathbf{C}}_1$  based on the matrix dynamics of  $\Phi$  (can be found in Appendix B.3) and  $\mathbf{W}_g$ :  
 735

$$\begin{aligned} 736 \dot{\mathbf{C}}_1 &= \dot{\Phi} \mathbf{W}_g + \Phi \dot{\mathbf{W}}_g - \dot{\mathbf{W}}_g \Phi - \mathbf{W}_g \dot{\Phi} \\ 737 &= \{-(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) \\ 738 &\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\} \\ 739 &\quad + \{(1 + \sigma^2) \mathbf{W}_h - \mathbf{I}\} \Phi \mathbf{W}_h + \rho \mathbf{W}_g\} \Phi \\ 740 &\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) \\ 741 &\quad - (\mathbf{W}_h \mathbf{W}_g \Phi + \Phi \mathbf{W}_g \mathbf{W}_h) + 2\rho \Phi\} \\ 742 &= -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] \\ 743 &\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] \\ 744 &\quad + \{(1 + \sigma^2) \mathbf{W}_h - \mathbf{I}\} \Phi \mathbf{W}_h, \Phi\} \\ 745 &= -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) \\ 746 &\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) \\ 747 &\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)\} \\ 748 &\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) \\ 749 &\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. \\ 750 \\ 751 \\ 752 \\ 753 \\ 754 \\ 755 \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  $(\xi, \eta, \zeta)$ :

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

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).

810  
811 **Lemma 1** (Tian et al. (2021, Lemma 2)). Let  $\mathbf{H}(t)$  be time-varying positive semidefinite matrices  
812 whose minimal eigenvalues are bounded away from zero:

$$813 \quad \inf_{t \geq 0} \lambda_{\min}(\mathbf{H}(t)) \geq \lambda_0 > 0.$$

814 Then, the following dynamics

$$816 \quad \frac{d\mathbf{w}(t)}{dt} = -\mathbf{H}(t)\mathbf{w}(t)$$

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

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

### 823 B.3 DECOUPLING INTO EIGENVALUE DYNAMICS

825 We have obtained the following matrix dynamics:

$$\begin{aligned} 827 \quad \dot{\mathbf{W}}_f &= -\mathbf{W}_h^\top \{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g)\mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\}\mathbf{W}_f - \rho\mathbf{W}_f, \\ 828 \quad \dot{\mathbf{W}}_g &= -\{(1 + \sigma^2)\mathbf{W}_h - \mathbf{I}\}\mathbf{W}_f \mathbf{W}_f^\top \mathbf{W}_h^\top - \rho\mathbf{W}_g, \\ 830 \quad \dot{\mathbf{W}}_h &= -\{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g)\mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\}\mathbf{W}_f \mathbf{W}_f^\top - \rho\mathbf{W}_h. \end{aligned}$$

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

$$\begin{aligned} 834 \quad \dot{\Phi} &= \dot{\mathbf{W}}_f \mathbf{W}_f^\top + \mathbf{W}_f \dot{\mathbf{W}}_f^\top \\ 835 &= -\mathbf{W}_h^\top \{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g)\mathbf{W}_h - (\mathbf{I} + \mathbf{W}_g^\top)\}\Phi \\ 836 &\quad - \Phi \{(1 + \sigma^2)\mathbf{W}_h^\top (\mathbf{I} + \mathbf{W}_g^\top \mathbf{W}_g) - (\mathbf{I} + \mathbf{W}_g)\} \mathbf{W}_h - 2\rho\Phi. \end{aligned}$$

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

$$\begin{aligned} 840 \quad \dot{\Phi} &= -(1 + \sigma^2)\{\mathbf{W}_h(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h, \Phi\} + \{\mathbf{W}_h, \Phi\} + (\mathbf{W}_h \mathbf{W}_g \Phi + \Phi \mathbf{W}_g \mathbf{W}_h) - 2\rho\Phi, \quad (5) \\ 841 \quad \dot{\mathbf{W}}_g &= -\{(1 + \sigma^2)\mathbf{W}_h - \mathbf{I}\}\Phi \mathbf{W}_h - \rho\mathbf{W}_g, \\ 843 \quad \dot{\mathbf{W}}_h &= -\{(1 + \sigma^2)(\mathbf{I} + \mathbf{W}_g^2)\mathbf{W}_h + (\mathbf{I} + \mathbf{W}_g)\}\Phi - \rho\mathbf{W}_h, \end{aligned}$$

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

846 Next, we decouple them into the corresponding eigenvalues. The parameter matrices are simultaneously  
847 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  
848 symmetry assumption (A1) and common eigenspace assumption (A2). Here, we can easily show  
849 that the eigenspace is time-independent, namely,  $\dot{\mathbf{U}} = \mathbf{O}$ , using the same argument of Tian et al.  
850 (2021, Appendix B.1). By multiplying  $\mathbf{U}^\top$  and  $\mathbf{U}$  from left and right, respectively,  $\Phi$ -dynamics can  
851 be written as follows:

$$853 \quad \dot{\Lambda}_\Phi = -2(1 + \sigma^2)(\mathbf{I} + \Lambda_g^2)\Lambda_h^2 \Lambda_\Phi + 2\Lambda_h \Lambda_\Phi + 2\Lambda_h \Lambda_g \Lambda_\Phi - 2\rho\Lambda_\Phi,$$

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

$$857 \quad \dot{\phi} = -2(1 + \sigma^2)(1 + \gamma^2)\psi^2\phi + 2\psi\phi + 2\psi\phi\gamma - 2\rho\phi.$$

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

$$\begin{aligned} 860 \quad \dot{\gamma} &= -(1 + \sigma^2)(\psi - 1)\phi\psi - \rho\gamma, \\ 861 \quad \dot{\psi} &= -\{(1 + \sigma^2)(1 + \gamma^2)\psi + (1 + \gamma)\}\phi - \rho\psi. \end{aligned}$$

863 Note that  $\psi$  and  $\gamma$  correspond to (one of) eigenvalues of the linear networks  $h$  and  $g$ , respectively.  
Intuitively, we can regard  $\psi$  and  $\gamma$  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  $(\phi, \psi, \gamma)$ :

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

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

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

#### B.4 ADIABATIC ELIMINATION

The eigenvalue dynamics obtained in Appendix B.3 is jointly with respect to  $(\phi, \psi, \gamma)$ . Here, we eliminate  $\phi$  by confirming that  $\phi$  and  $\psi$  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  $\phi$ :

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

Note that the convergence to the invariant parabola is faster when weight decay  $\rho$  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()
```

918 Table 3: **X-PhiNet performs robustly well for different weight decays on CIFAR-5m.**  
919

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

931 Table 4: **X-PhiNet performs robustly well for one epoch training.**  
932

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

```

943   11 Lcos = D(p1, z2)/2 + D(p2, z1)/2 # loss
944   12 Lcor = mse_loss(y1, z0)/2 + mse_loss(y2, z0)/2
945   13 L = Lcos + Lcor
946   14 L.backward() # back-propagate
947   15 update(f, h) # SGD update
948   16 f_long = beta * f_long + (1-beta) * f # EMA for projection
949
950   18 def D(p, z): # negative cosine similarity
951   19     z = z.detach() # stop gradient
952   20     p = normalize(p, dim=1) # l2-normalize
953   21     z = normalize(z, dim=1) # l2-normalize
954   22     return -(p*z).sum(dim=1).mean()
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2189
2190
2191
2192
2193
2194
2195
2196
2197
2197
2198
2199
2199
2200
2201
2202
2203
2204
2205
2206
2206
2207
2208
2209
2210
2211
2212
2213
2213
2214
2215
2216
2217
2218
2219
2219
2220
2221
2222
2223
2224
2225
2226
2226
2227
2228
2229
2229
2230
2231
2232
2233
2234
2235
2235
2236
2237
2238
2239
2239
2240
2241
2242
2243
2244
2245
2245
2246
2247
2248
2249
2249
2250
2251
2252
2253
2254
2255
2255
2256
2257
2258
2259
2259
2260
2261
2262
2263
2264
2264
2265
2266
2267
2268
2268
2269
2270
2271
2272
2272
2273
2274
2275
2275
2276
2277
2278
2278
2279
2280
2281
2282
2282
2283
2284
2285
2285
2286
2287
2288
2288
2289
2290
2291
2292
2292
2293
2294
2295
2295
2296
2297
2298
2298
2299
2299
2300
2301
2302
2302
2303
2304
2305
2305
2306
2307
2308
2308
2309
2310
2311
2311
2312
2313
2314
2314
2315
2316
2317
2317
2318
2319
2319
2320
2321
2322
2322
2323
2324
2325
2325
2326
2327
2328
2328
2329
2330
2331
2331
2332
2333
2334
2334
2335
2336
2337
2337
2338
2339
2339
2340
2341
2342
2342
2343
2344
2345
2345
2346
2347
2348
2348
2349
2350
2351
2351
2352
2353
2354
2354
2355
2356
2357
2357
2358
2359
2359
2360
2361
2362
2362
2363
2364
2364
2365
2366
2366
2367
2368
2368
2369
2370
2370
2371
2372
2372
2373
2374
2374
2375
2376
2376
2377
2378
2378
2379
2380
2380
2381
2382
2382
2383
2384
2384
2385
2386
2386
2387
2388
2388
2389
2390
2390
2391
2392
2392
2393
2394
2394
2395
2396
2396
2397
2398
2398
2399
2399
2400
2401
2401
2402
2403
2403
2404
2405
2405
2406
2407
2407
2408
2409
2409
2410
2411
2411
2412
2413
2413
2414
2415
2415
2416
2417
2417
2418
2419
2419
2420
2421
2421
2422
2423
2423
2424
2425
2425
2426
2427
2427
2428
2429
2429
2430
2431
2431
2432
2433
2433
2434
2435
2435
2436
2437
2437
2438
2439
2439
2440
2441
2441
2442
2443
2443
2444
2445
2445
2446
2447
2447
2448
2449
2449
2450
2451
2451
2452
2453
2453
2454
2455
2455
2456
2457
2457
2458
2459
2459
2460
2461
2461
2462
2463
2463
2464
2465
2465
2466
2467
2467
2468
2469
2469
2470
2471
2471
2472
2473
2473
2474
2475
2475
2476
2477
2477
2478
2479
2479
2480
2481
2481
2482
2483
2483
2484
2485
2485
2486
2487
2487
2488
2489
2489
2490
2491
2491
2492
2493
2493
2494
2495
2495
2496
2497
2497
2498
2499
2499
2500
2501
2501
2502
2503
2503
2504
2505
2505
2506
2507
2507
2508
2509
2509
2510
2511
2511
2512
2513
2513
2514
2515
2515
2516
2517
2517
2518
2519
2519
2520
2521
2521
2522
2523
2523
2524
2525
2525
2526
2527
2527
2528
2529
2529
2530
2531
2531
2532
2533
2533
2534
2535
2535
2536
2537
2537
2538
2539
2539
2540
2541
2541
2542
2543
2543
2544
2545
2545
2546
2547
2547
2548
2549
2549
2550
2551
2551
2552
2553
2553
2554
2555
2555
2556
2557
2557
2558

```

972 D.2 CONTINUAL LEARNING  
 973

974 **Epochs per task.** In Table 2, we trained on each task for one epoch. However, in Madaan et al.  
 975 (2022), 200 epochs are trained for each task on Split CIFAR10, and the number of iterations differs  
 976 from this case. The effect of early stopping may be apparent when the number of iterations is different.  
 977 Thus, we trained on each task for two epochs to match the number of iterations. The result is shown  
 978 in Table 5. The performance of X-PhiNet is still high even when the number of epochs per task is set  
 979 to 2 epochs.

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

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

983  
 984 **Replay with Mixup.** Replay is one of the most promising methods for improving the performance  
 985 of continual learning while additional memory costs are required (Hsu et al., 2018; Van de Ven et al.,  
 986 2020; Madaan et al., 2022; Lin et al., 2022). We thus examined the performance of our method in  
 987 combination with the mixup-based replay method proposed in (Madaan et al., 2022). Table 6 shows  
 988 that when only one epoch is trained for each task, X-PhiNet shows considerably higher accuracy  
 989 than the other methods. On the other hand, when we train two epochs for each task, the accuracy of  
 990 other methods such as BYOL, BarlowTwins and RM-SimSiam also increases, showing an accuracy  
 991 comparable to that of X-PhiNet.

992 **Table 6: X-PhiNet performs higher or comparable results for split-CIFAR5m even with Mixup.**  
 993 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	Acc	90.18 <sub>0.65</sub>	91.51 <sub>0.42</sub>	90.74 <sub>0.63</sub>	91.66 <sub>0.21</sub>	91.92 <sub>0.17</sub>	91.78 <sub>0.29</sub>	<b>92.43</b> <sub>0.14</sub>
(1epoch)	Fg	0.36 <sub>3.07</sub>	-1.70 <sub>0.09</sub>	-0.97 <sub>3.05</sub>	-2.21 <sub>0.54</sub>	-1.14 <sub>0.27</sub>	-0.73 <sub>0.07</sub>	-0.69 <sub>0.65</sub>

		BYOL	SimSiam	Barlow Twins	PhiNet	RM-SimSiam	X-PhiNet (MSE)	X-PhiNet (Cos)
Split C-5m	Acc	<b>92.36</b> <sub>0.03</sub>	91.77 <sub>0.01</sub>	<b>92.36</b> <sub>0.70</sub>	91.00 <sub>1.78</sub>	<b>92.48</b> <sub>0.12</sub>	<b>92.12</b> <sub>0.28</sub>	<b>92.26</b> <sub>0.35</sub>
(2epoch)	Fg	0.64 <sub>0.01</sub>	2.24 <sub>0.01</sub>	-0.07 <sub>0.64</sub>	1.97 <sub>1.41</sub>	1.05 <sub>0.31</sub>	2.02 <sub>0.71</sub>	1.82 <sub>0.71</sub>

1000  
 1001 **Split CIFAR10 and Split CIFAR100.** Up to this point, we have experimented with continual  
 1002 learning using the CIFAR-5m-based dataset. Now, we test on the standard benchmarks, Split  
 1003 CIFAR10 and Split CIFAR100. Table 7 shows that in both Split CIFAR10 and Split CIFAR100,  
 1004 X-PhiNet outperforms SimSiam. However, PhiNet sometimes shows higher accuracy than X-PhiNet.  
 1005 Note that PhiNet is a special case of X-PhiNet, and we have set the momentum of X-PhiNet to 0.99  
 1006 in this study. If we carefully select the momentum value, X-PhiNet’s performance might improve,  
 1007 surpassing PhiNet. When using mixup for replay, X-PhiNet shows significantly higher accuracy  
 1008 compared to other methods.

1009  
 1010 **Effect of exponential moving average** X-PhiNet has an additional hyperparameter, the exponential  
 1011 moving average. We set  $\beta = 0.99$  in all the experiments in this paper. As shown in table.8, in  
 1012 tasks such as continual learning, where it is important to apply a strong exponential moving average,  
 1013 accuracy increases as  $\beta$  increases and then decreases again from a certain point.

1014 E ADDITIONAL EXPERIMENTS ON THE ROBUSTNESS OF PHINET

1015 E.1 ADDITIONAL ABLATION STUDY WITH CIFAR10

1016 **Usage of original input:** We first investigate that what is the best way to input the original signal to  
 1017 the model. To this end, we first replace one of augmented signals in SimSiam as an original input.

Table 7: **X-PhiNet performs good results when memorization is important.** We trained X-PhiNet on Split CIFAR10. In Split CIFAR-5m, 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 <sub>0.29</sub>	89.35 <sub>0.08</sub>	<b>91.25<sub>0.09</sub></b>	90.65 <sub>0.43</sub>	90.90 <sub>0.50</sub>	<b>90.97<sub>0.49</sub></b>
	Fg	5.31 <sub>0.65</sub>	3.70 <sub>0.16</sub>	4.86 <sub>0.50</sub>	1.02 <sub>0.31</sub>	5.72 <sub>0.82</sub>	3.95 <sub>0.37</sub>
Split C100 (FineTune)	Acc	77.93 <sub>0.64</sub>	78.19 <sub>0.41</sub>	<b>78.50<sub>0.25</sub></b>	<b>78.31<sub>0.16</sub></b>	77.50 <sub>0.04</sub>	77.44 <sub>0.28</sub>
	Fg	7.06 <sub>1.00</sub>	-0.57 <sub>0.98</sub>	6.51 <sub>0.31</sub>	5.49 <sub>1.27</sub>	8.46 <sub>0.21</sub>	4.46 <sub>0.34</sub>
Split C10 (Mixup)	Acc	90.68 <sub>0.89</sub>	91.14 <sub>0.84</sub>	89.89 <sub>0.69</sub>	90.69 <sub>0.21</sub>	90.49 <sub>0.32</sub>	<b>91.56<sub>0.12</sub></b>
	Fg	0.85 <sub>0.16</sub>	1.08 <sub>0.60</sub>	1.04 <sub>0.22</sub>	-2.11 <sub>1.61</sub>	1.80 <sub>0.09</sub>	1.36 <sub>0.15</sub>
Split C100 (Mixup)	Acc	81.77 <sub>0.14</sub>	82.47 <sub>0.70</sub>	80.76 <sub>0.18</sub>	82.16 <sub>0.76</sub>	83.32 <sub>0.04</sub>	<b>83.88<sub>0.26</sub></b>
	Fg	1.23 <sub>0.81</sub>	-1.35 <sub>0.05</sub>	1.18 <sub>1.27</sub>	-1.23 <sub>2.11</sub>	1.28 <sub>0.34</sub>	-0.07 <sub>0.30</sub>

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

		EMA $\beta$					
		0.999	0.997	0.99	0.97	0.9	0.7
Split C100 (Mixup)	Acc	83.19 <sub>0.22</sub>	82.62 <sub>0.06</sub>	83.88 <sub>0.26</sub>	82.72 <sub>0.98</sub>	81.76 <sub>0.64</sub>	81.73 <sub>0.06</sub>
	Fg	-0.22 <sub>0.34</sub>	0.33 <sub>0.06</sub>	-0.07 <sub>0.30</sub>	0.73 <sub>0.64</sub>	2.15 <sub>0.11</sub>	1.82 <sub>0.69</sub>

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, StopGradient, and 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 an role of a slow learner through experiments. In Figure 9 of the additional pdf, 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 learn 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

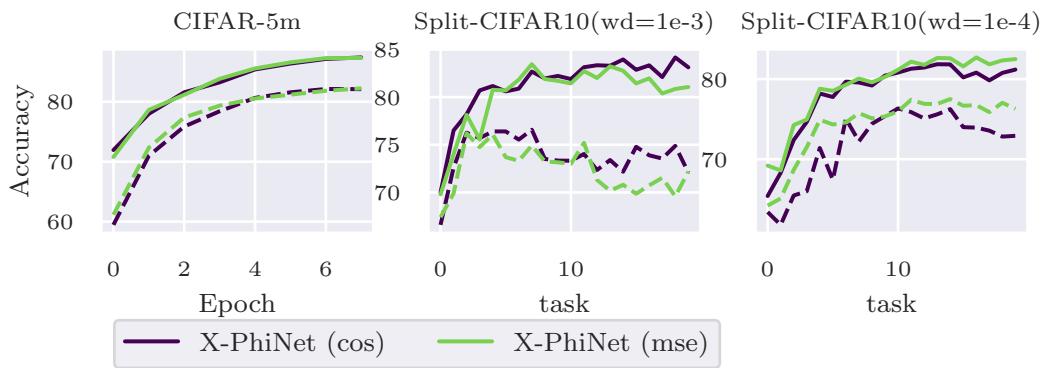
Table 10 shows the ablation of additional predictors in ImageNet. In this case, we used a higher weight decay of 1e-3.  $g = h$  has a lower accuracy than other methods, which is consistent with the results in CIFAR10.

## 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 size. 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.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
Table 9: Ablation study for PhiNet using CIFAR10 data. We use SGD with momentum as a 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 <sub>0.39</sub>	90.39 <sub>0.10</sub>	90.98 <sub>0.02</sub>
SimSiam (Orig-In)	–	–	–	72.82 <sub>0.18</sub>	76.67 <sub>1.13</sub>	69.03 <sub>11.37</sub>
PhiNet	MSE	✓	$g$	77.77 <sub>1.13</sub>	90.77 <sub>0.22</sub>	91.38 <sub>0.19</sub>
	MSE	✓	$g = \mathbf{I}$	77.63 <sub>0.11</sub>	91.01 <sub>0.12</sub>	91.50 <sub>0.07</sub>
	MSE		$g$	62.80 <sub>0.47</sub>	91.40 <sub>0.23</sub>	89.01 <sub>0.55</sub>
	MSE	✓	$h$	74.87 <sub>0.58</sub>	91.23 <sub>0.12</sub>	91.18 <sub>0.34</sub>
	Cos	✓	$g$	80.06 <sub>0.47</sub>	87.73 <sub>0.26</sub>	88.27 <sub>0.24</sub>
	Cos	✓	$g = \mathbf{I}$	75.34 <sub>3.27</sub>	87.38 <sub>0.17</sub>	87.90 <sub>0.10</sub>
	Cos		$g$	27.57 <sub>4.41</sub>	9.98 <sub>0.00</sub>	9.98 <sub>0.00</sub>
	Cos	✓	$h$	75.99 <sub>0.28</sub>	85.97 <sub>0.23</sub>	85.04 <sub>0.11</sub>



1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
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.

1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
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 the computational costs.

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

1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
Table 13 shows the results for STL10, where PhiNet performs comparably to SimSiam. However, as 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

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

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

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 <sub>2.46</sub>	87.26 <sub>0.26</sub>	87.53 <sub>0.20</sub>	87.17 <sub>0.05</sub>
	PhiNet	84.28 <sub>0.41</sub>	87.32 <sub>0.16</sub>	87.22 <sub>0.12</sub>	87.01 <sub>0.28</sub>
0.0005	SimSiam	88.89 <sub>0.34</sub>	89.23 <sub>0.11</sub>	89.39 <sub>0.12</sub>	88.57 <sub>0.40</sub>
	PhiNet	89.33 <sub>0.13</sub>	89.26 <sub>0.02</sub>	89.36 <sub>0.27</sub>	88.62 <sub>0.36</sub>
0.001	SimSiam	89.54 <sub>0.05</sub>	89.61 <sub>0.11</sub>	89.37 <sub>0.03</sub>	nan <sub>nan</sub>
	PhiNet	89.52 <sub>0.06</sub>	89.71 <sub>0.07</sub>	89.28 <sub>0.23</sub>	10.00 <sub>0.01</sub>

	Pretrained	VOC 07 detection			VOC 07+12 detection		
		AP <sub>50</sub>	AP	AP <sub>75</sub>	AP <sub>50</sub>	AP	AP <sub>75</sub>
MoCo v2		73.2	<b>46.6</b>	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	<b>46.6</b>	49.8	<b>82.7</b>	<b>57.3</b>	64.6
PhiNet (lr=0.05, wd=1e-4)		72.7	46.35	<b>50.4</b>	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	<b>63.9</b>
X-PhiNet (lr=0.5, wd=1e-5)		<b>74.9</b>	45.9	50.1	<b>82.7</b>	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.

performance to MoCo across various tasks. Furthermore, as shown in Table 15, we can find that 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

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

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.

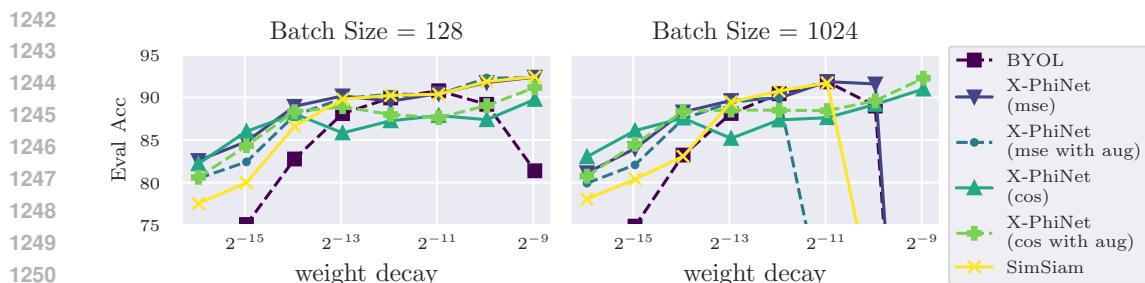


Figure 10: Even when  $x$  is augmented, PhiNet is more robust to weight decay than SimSiam (CIFAR10). We used the same setting as in Figure.6 in the original paper. “with aug” performs data augmentation for  $x$ , which is not data augmented in Figure.6.

on batch size, are dominant in this setting. Additionally, Table 18 includes a comparison of training times, demonstrating that PhiNet can be trained in a time comparable to SimSiam.

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)

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)

## 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  $\text{srank}(M) = \|M\|_F^2 / \|M\|^2$ , 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 widerer 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 available<sup>1</sup>.

### F.2 SETTINGS FOR TRAINING ON IMAGENET

In our ImageNet (Russakovsky et al., 2015) experiments, we follow the formal implementation of SimSiam by Pytorch<sup>2</sup> (Chen and He, 2021). Table 20 shows the model and experimental setup for

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

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

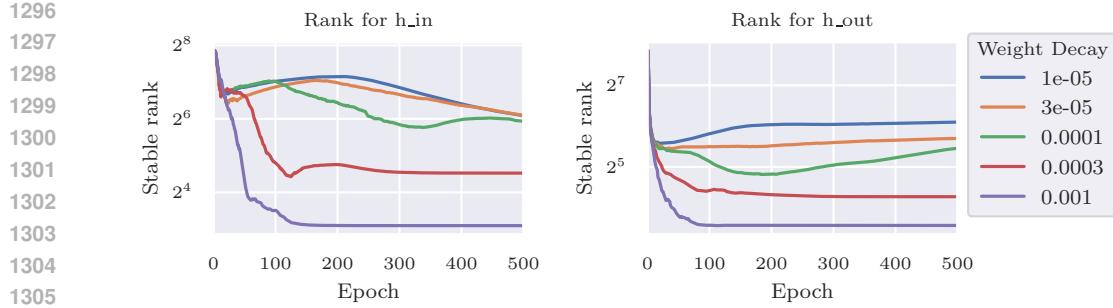


Figure 11: **The smaller the weight decay, the larger the rank of the additional predictor.** We trained PhiNet on CIFAR10 with SGD and evaluate the rank for layers in additional predictor blocks.

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 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.

### 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).

Table 20: The experimental setups of Table 1.

1352		Optimiser	SGD
1353		Momentum	0.9
1354	Learning	Learning Rate	0.05
1355		Epochs	100
1356		Encoder	ResNet50
1357		Backbone	
1358		Projector output dimension	2048
1359		Predictor $h$	
1360		Latent dimension $m$	2048
1361		Hidden dimension $h$	512
1362		Activation function	ReLU
1363		Batch normalization	Yes
1364		Predictor $g$	
1365		Latent dimension $m$	2048
1366		Hidden dimension $g$	512
1367		Activation function (Hidden)	ReLU
1368	Computational resource	Activation function (Output)	Tanh
1369		Batch normalization	Yes
1370		Computational resource	GPUs
1371			4×V100
1372			

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:

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