

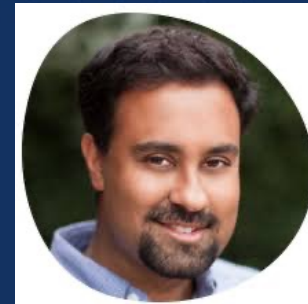
# Understanding self-supervised Learning Dynamics without Contrastive Pairs



**Yuandong Tian<sup>1</sup>**



**Xinlei Chen<sup>1</sup>**



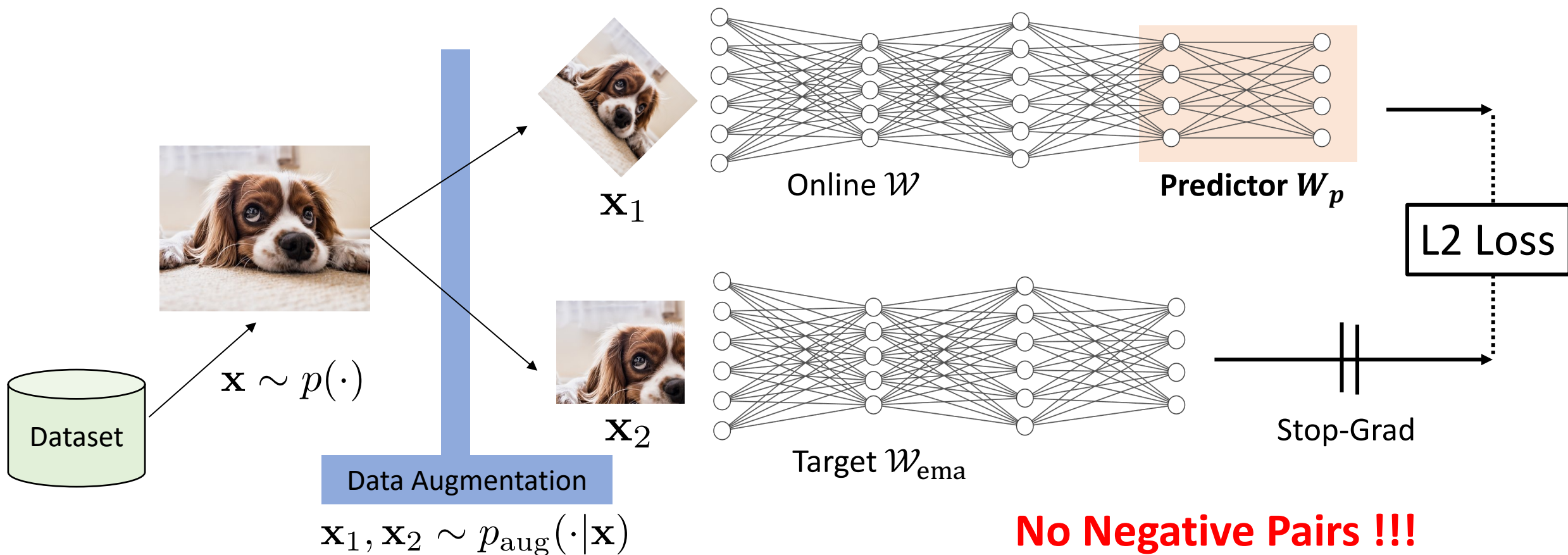
**Surya Ganguli<sup>1,2</sup>**

<sup>1</sup> Facebook AI Research

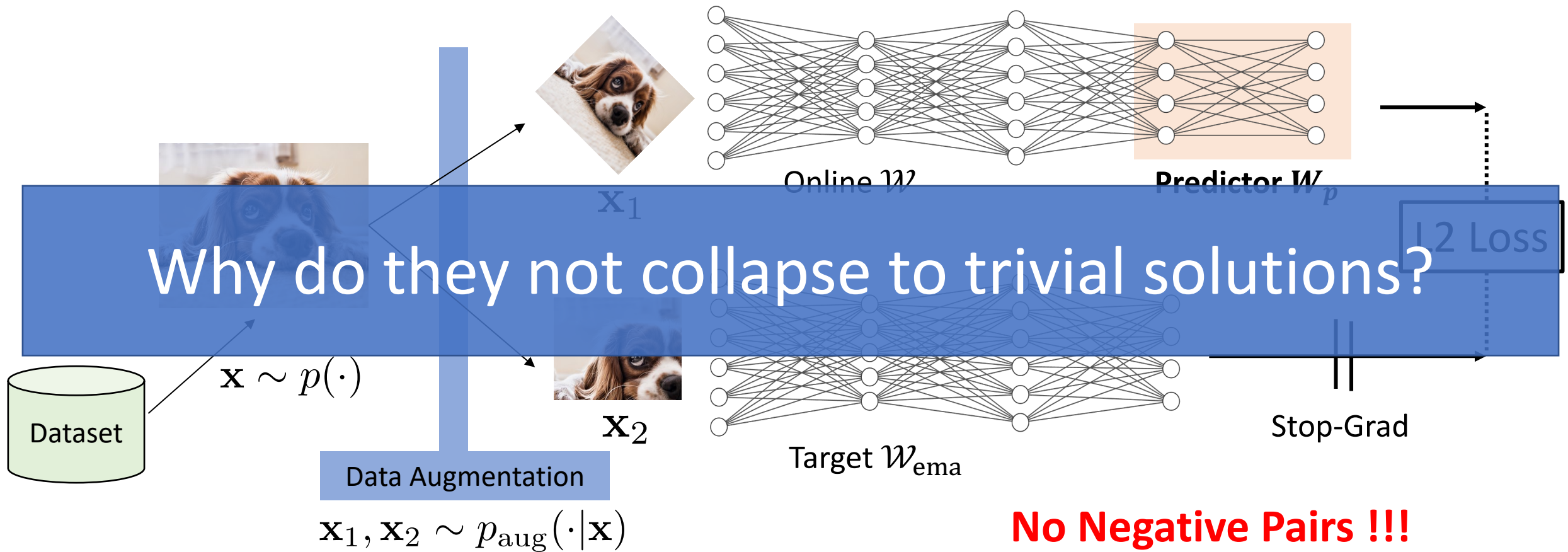
<sup>2</sup> Stanford University

**ICML 2021**

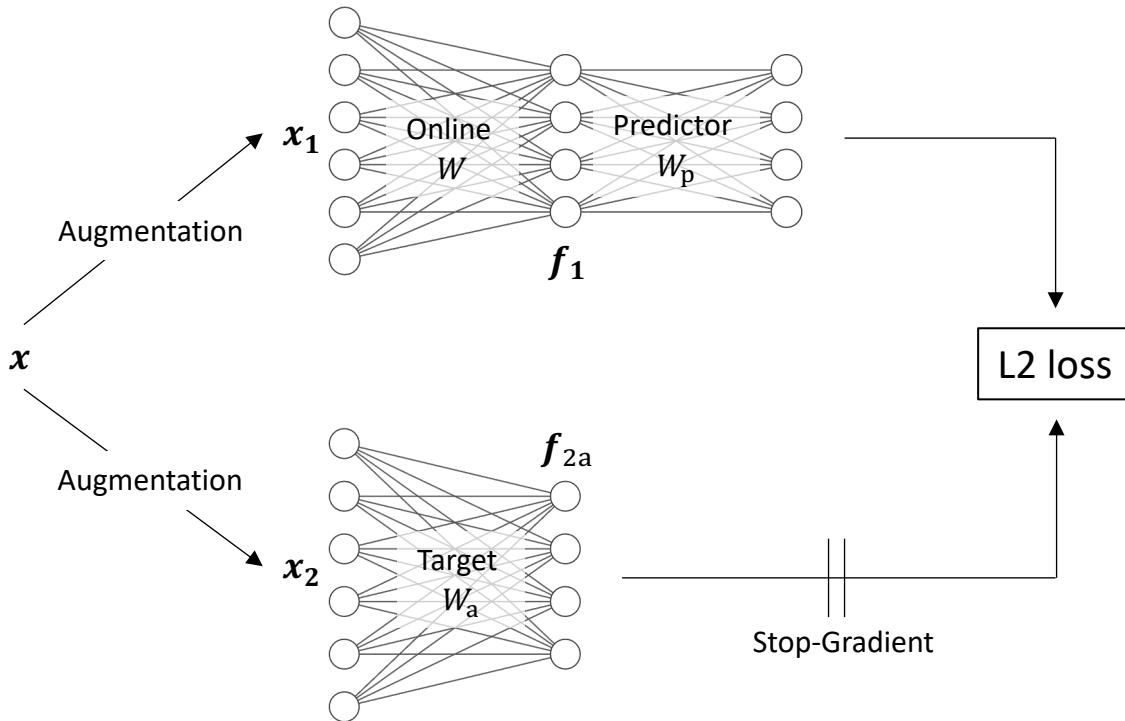
# Non-contrastive SSL (BYOL/SimSiam)



# Non-contrastive SSL (BYOL/SimSiam)?



# A simple model



Objective:

$$J(W, W_p) := \frac{1}{2} \mathbb{E}_{x_1, x_2} [\|W_p \mathbf{f}_1 - \text{StopGrad}(\mathbf{f}_{2a})\|_2^2]$$

Linear online network  $W$

Linear target network  $W_a$

Linear predictor  $W_p$

# The Dynamics of Training Procedure

**Lemma 1.** *BYOL learning dynamics following Eqn. 1:*

$$\dot{W}_p = \alpha_p (-W_p W (X + X') + W_a X) W^\top - \eta W_p$$

$$\dot{W} = W_p^\top (-W_p W (X + X') + W_a X) - \eta W$$

$$\dot{W}_a = \beta (-W_a + W)$$

$$\bar{\mathbf{x}}(\mathbf{x}) := \mathbb{E}_{\mathbf{x}' \sim p_{\text{aug}}(\cdot | \mathbf{x})} [\mathbf{x}']$$

$$X = \mathbb{E} [\bar{\mathbf{x}} \bar{\mathbf{x}}^\top] \quad \text{Covariance of the data}$$

$$X' = \mathbb{E}_{\mathbf{x}} [\mathbb{V}_{\mathbf{x}' | \mathbf{x}} [\mathbf{x}']] \quad \text{Covariance of the augmentation}$$

**Part I** Why we need (1) an **extra predictor** and (2) **stop-gradient**?

**Part II** Why the system doesn't **collapse** to trivial solutions?

**Part III** The role played by different hyperparameters

Hyperparameter	Description
$\alpha_p$	Relative learning rate of the predictor
$\eta$	Weight decay
$\beta$	The rate of Exponential Moving Average (EMA)

**Part IV** Novel non-contrastive SSL algorithm **DirectPred**

# Part I No Predictor / No Stop-Gradient do not work

If there is no EMA ( $W = W_a$ ), then the dynamics changes:

No Predictor

$$\dot{W} = -\underbrace{(X' + \eta I)}_{\text{PSD matrix}} W$$

No Stop-Gradient (Here  $\tilde{W}_p := W_p - I$ )

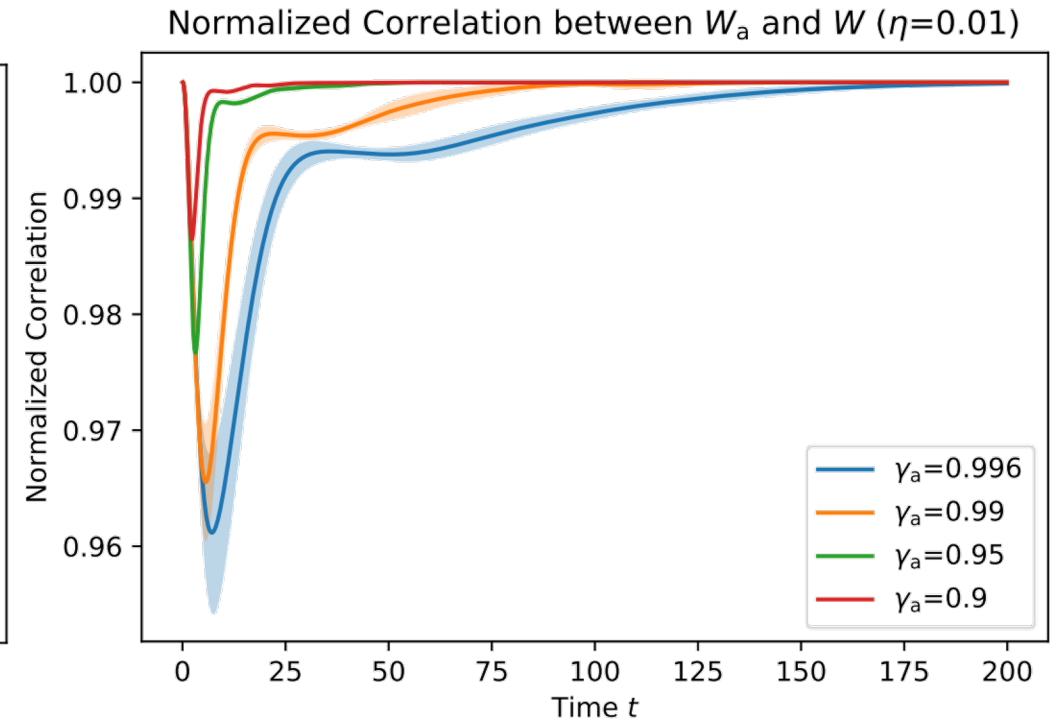
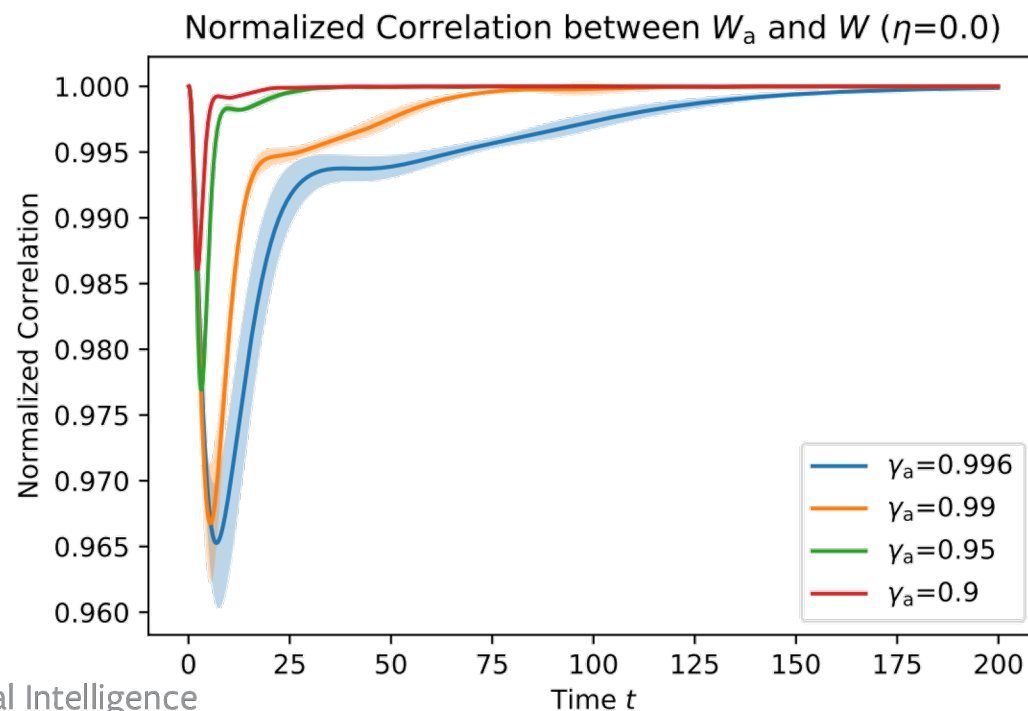
$$\frac{d}{dt} \text{vec}(W) = -\underbrace{\left[ X' \otimes (W_p^\top W_p + I) + X \otimes \tilde{W}_p^\top \tilde{W}_p + \eta I_{n_1 n_2} \right]}_{\text{PSD matrix}} \text{vec}(W)$$

In both cases,  $W \rightarrow 0$

# Part II Assumptions

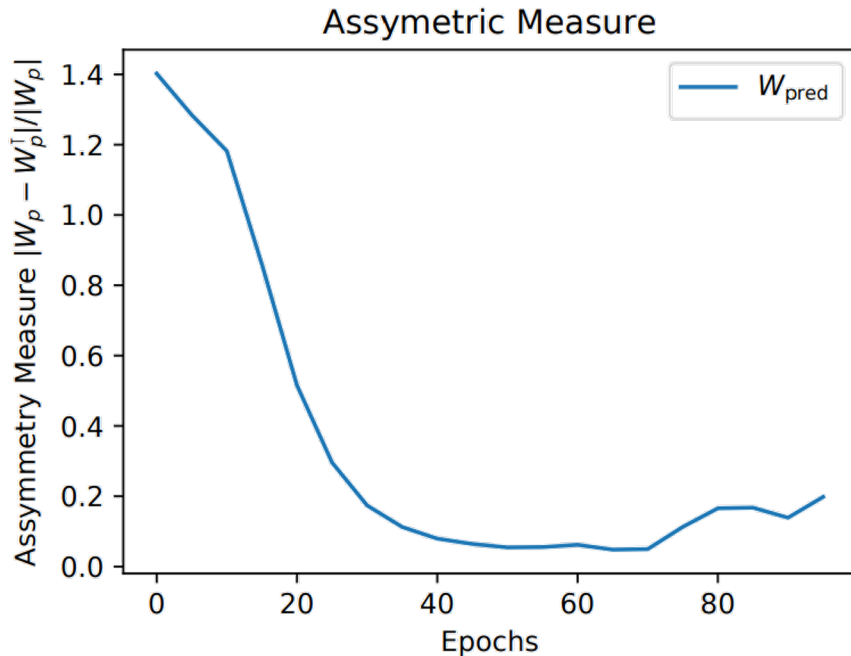
Assumption 1 (Isotropic Data and Augmentation):  $X = I$  and  $X' = \sigma^2 I$

Assumption 2: the EMA weight  $W_a(t) = \tau(t)W(t)$  is a linear function of  $W(t)$



# Symmetrization of the dynamics

Assumption 3 (Symmetric predictor  $W_p$ ):  $W_p(t) = W_p^T(t)$



$W_p$  becomes increasingly symmetric over training

	No predictor bias		With predictor bias	
	sym $W_p$	regular $W_p$	sym $W_p$	regular $W_p$
<i>One-layer linear predictor</i>				
EMA	75.09 ± 0.48	74.51 ± 0.47	74.52 ± 0.29	74.16 ± 0.33
no EMA	<b>36.62 ± 1.85</b>	72.85 ± 0.16	<b>36.04 ± 2.74</b>	72.13 ± 0.53
<i>Two-layer predictor with BatchNorm and ReLU</i>				
EMA	71.58 ± 6.46	78.85 ± 0.25	77.64 ± 0.41	78.53 ± 0.34
no EMA	<b>35.59 ± 2.10</b>	65.98 ± 0.71	<b>41.92 ± 4.25</b>	65.59 ± 0.66

Perfect symmetric  $W_p$  might hurt training



# Symmetrized Dynamics

Under the three assumptions, the dynamics becomes:

$$\begin{aligned}\dot{W}_p &= -\frac{\alpha_p}{2}(1 + \sigma^2)\{W_p, F\} + \alpha_p\tau F - \eta W_p \\ \dot{F} &= -(1 + \sigma^2)\{W_p^2, F\} + \tau\{W_p, F\} - 2\eta F\end{aligned}$$

$\{A, B\} := AB + BA$  is the anti-commutator.

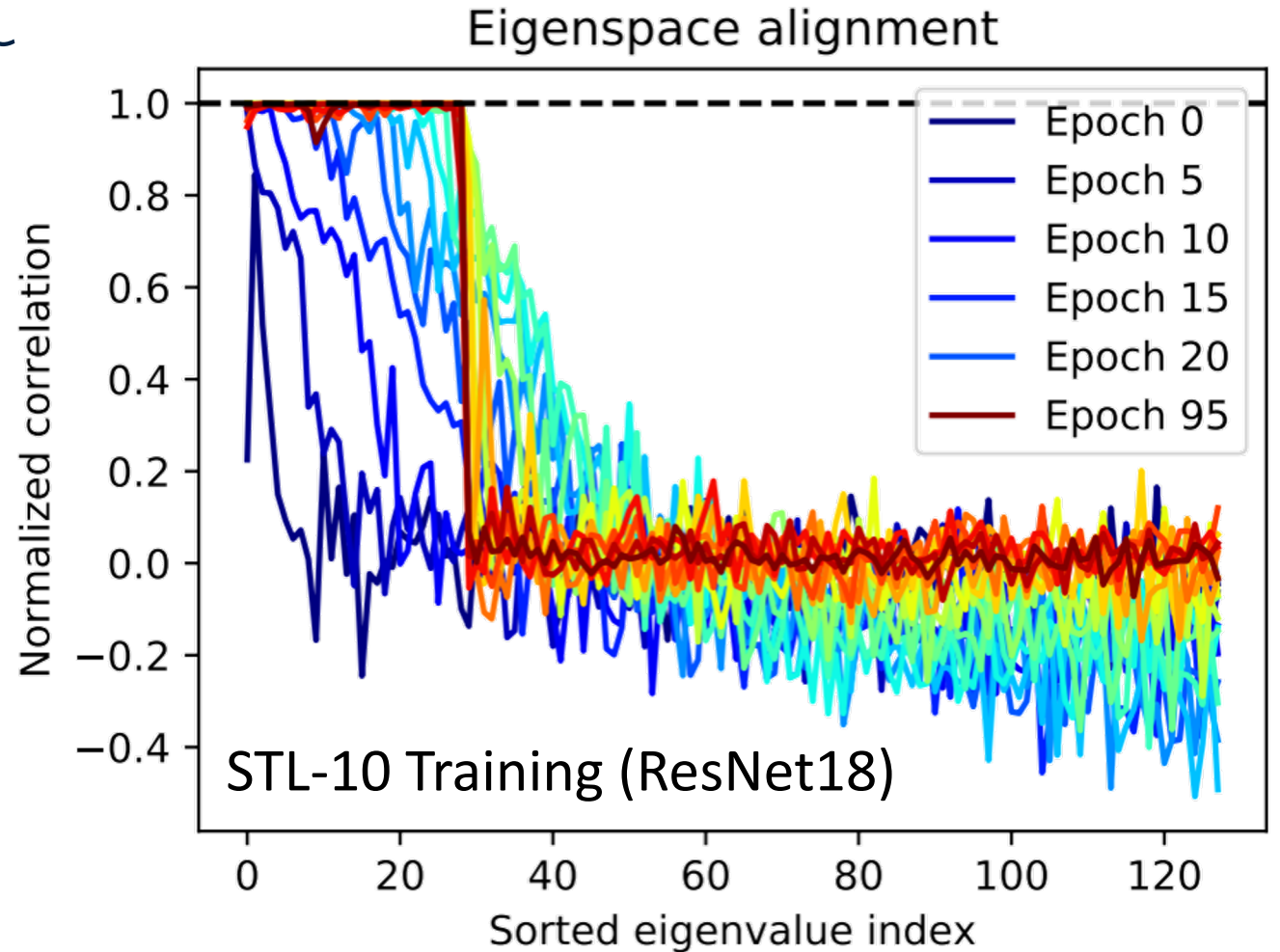
Here  $F := E[ff^T] = WXW^T$  is the correlation matrix of the input of the predictor  $W_p$ .  $F$  is well-defined even with nonlinear network.

# Eigenspace Alignment

Theorem 3: Under certain conditions,

$$FW_p - W_pF \rightarrow 0$$

and the eigenspace of  $W_p$  and  $F$  gradually **aligns**.



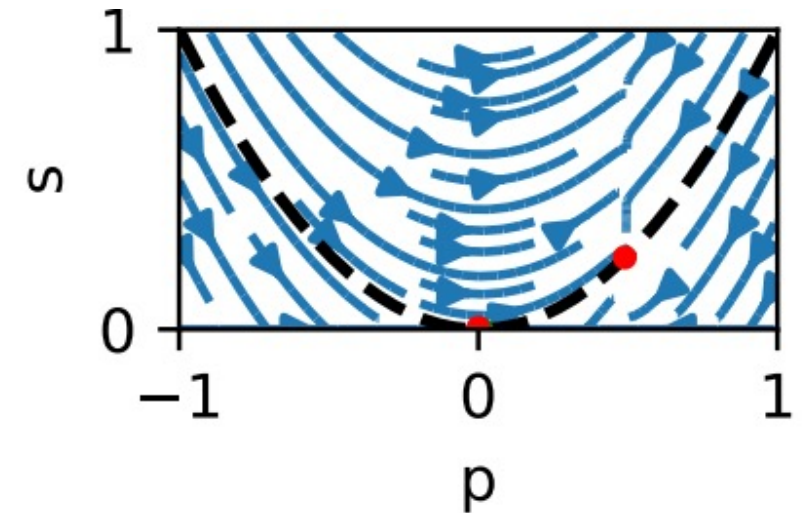
# Why non-contrastive SSL doesn't collapse?

When eigenspace aligns, the dynamics becomes decoupled:

$$\begin{aligned}\dot{p}_j &= \alpha_p s_j [\tau - (1 + \sigma^2)p_j] - \eta p_j \\ \dot{s}_j &= 2p_j s_j [\tau - (1 + \sigma^2)p_j] - 2\eta s_j \\ s_j \dot{\tau} &= \beta(1 - \tau)s_j - \tau \dot{s}_j / 2.\end{aligned}$$

Where  $p_j$  and  $s_j$  are eigenvalues of  $W_p$  and  $F$

Invariance holds:  $s_j(t) = \alpha_p^{-1} p_j^2(t) + e^{-2\eta t} c_j$



# Why non-contrastive SSL doesn't collapse?

1D dynamics of the eigenvalue  $p_j$  of  $W_p$ :

$$\dot{p}_j = p_j^2 \left[ \tau(t) - (1 + \sigma^2)p_j \right] - \eta p_j$$

EMA

Variance due to data augmentation

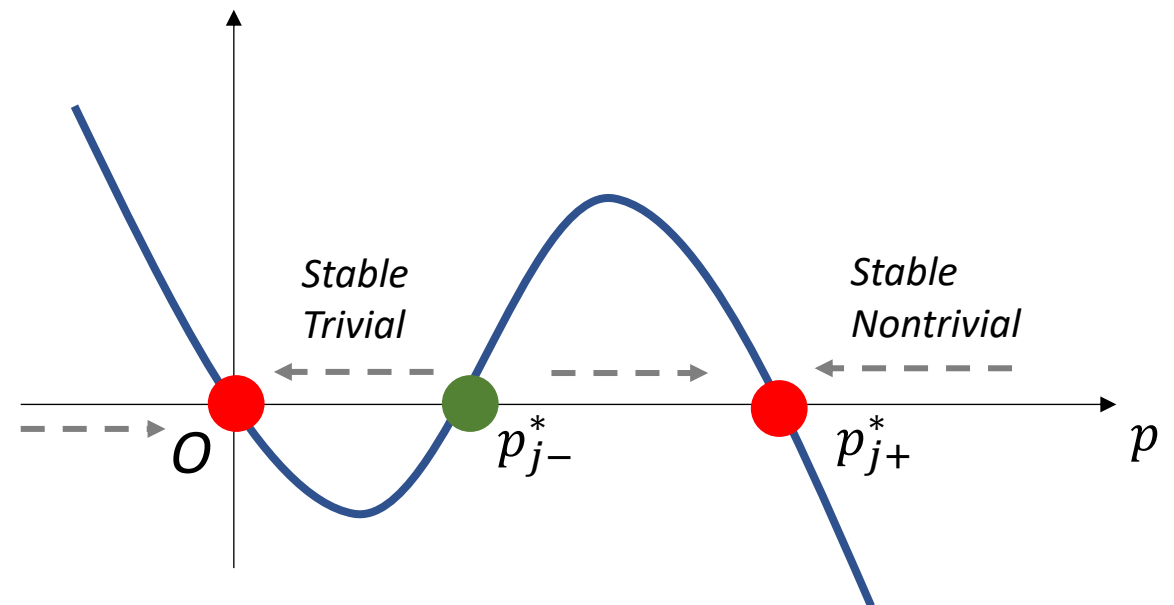
Weight Decay

# Why non-contrastive SSL doesn't collapse?

1D dynamics of the eigenvalue  $p_j$  of  $W_p$ :

$$\dot{p}_j = p_j^2 \left[ \tau(t) - (1 + \sigma^2)p_j \right] - \eta p_j$$

EMA  
Variance due to data augmentation  
Weight Decay



● Stable stationary point      ● Unstable stationary point

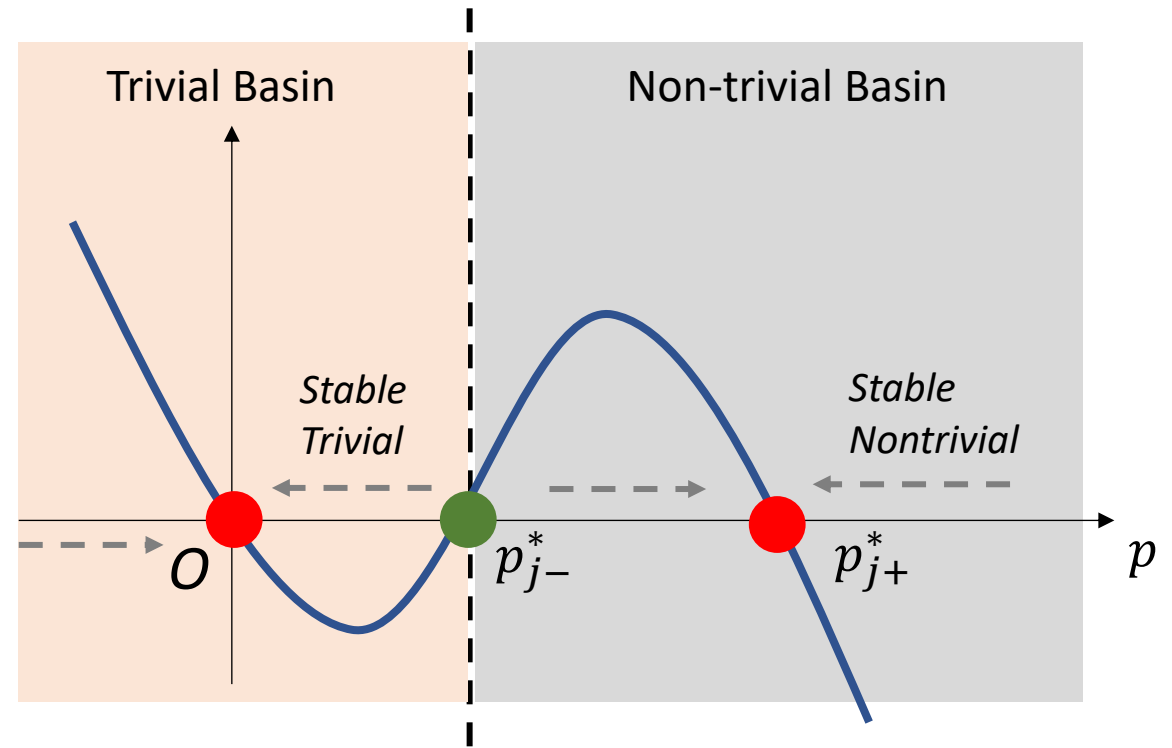
# Why non-contrastive SSL doesn't collapse?

1D dynamics of the eigenvalue  $p_j$  of  $W_p$ :

$$\dot{p}_j = p_j^2 \left[ \tau(t) - (1 + \sigma^2)p_j \right] - \eta p_j$$

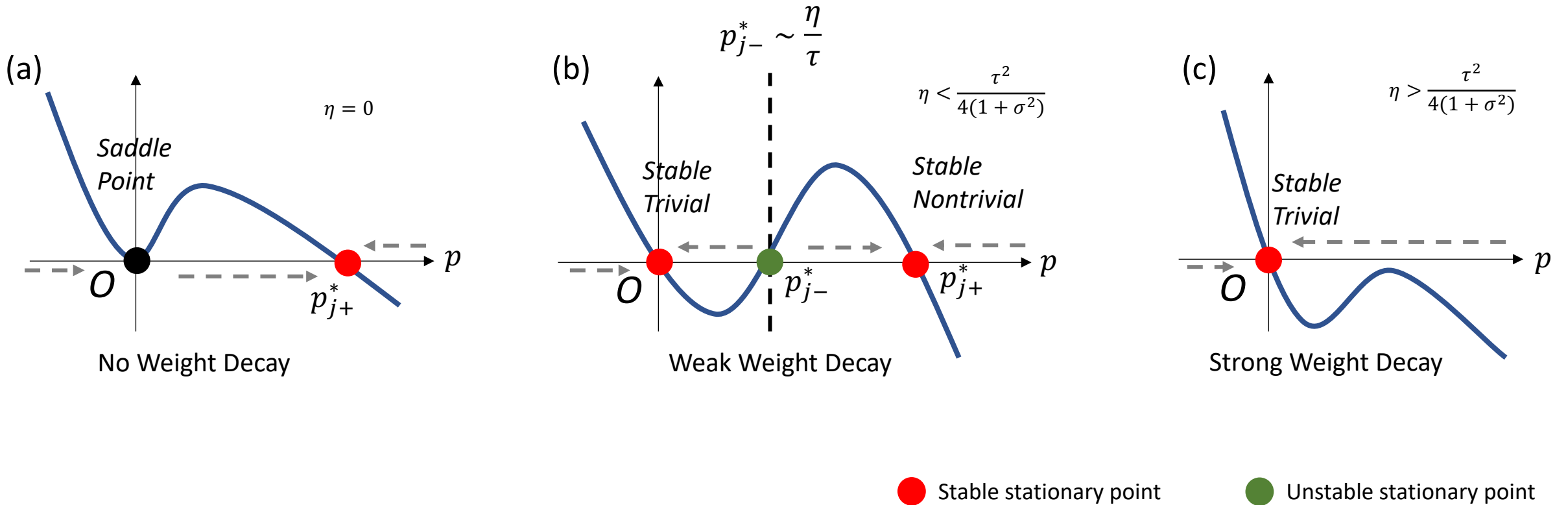
EMA (points to  $\tau(t)$ )  
 Variance due to data augmentation (points to  $\sigma^2$ )  
 Weight Decay (points to  $\eta$ )

$$p_{j-}^* = \frac{\tau - \sqrt{\tau^2 - 4\eta(1 + \sigma^2)}}{2(1 + \sigma^2)} \sim \frac{\eta}{\tau}$$



● Stable stationary point      ● Unstable stationary point

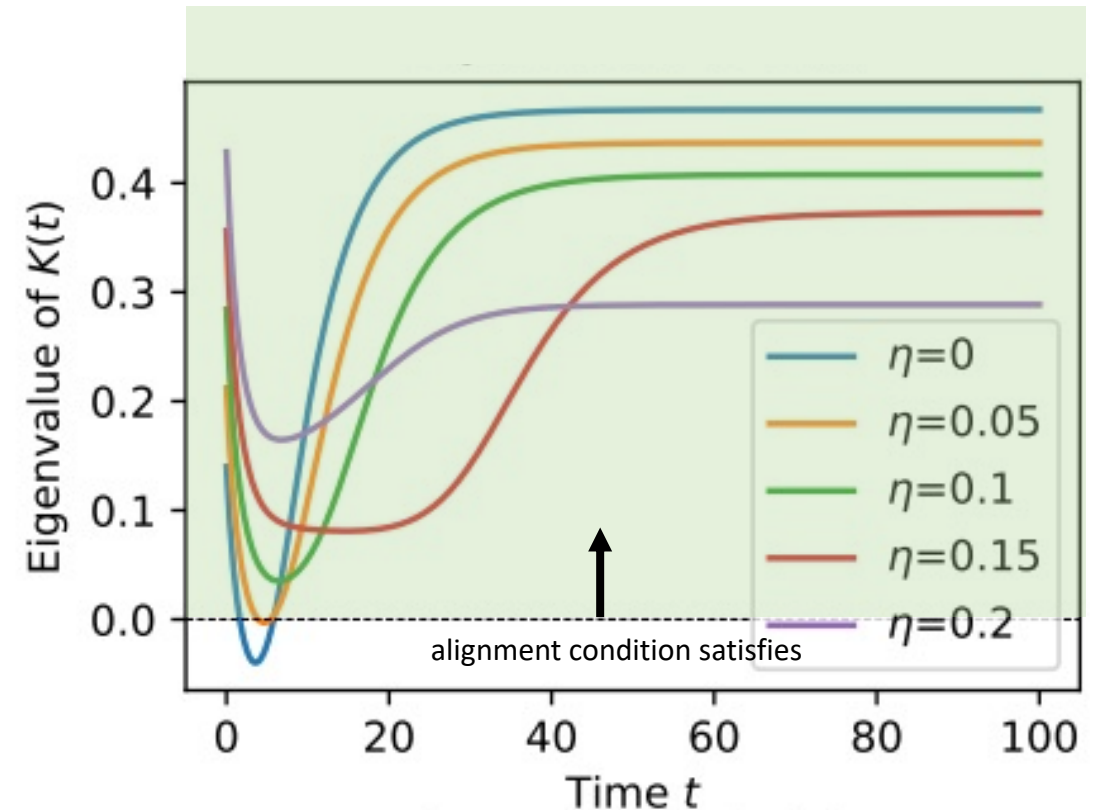
# Part III The Effect of Weight Decay $\eta$



# The Benefit of Weight Decay

Eigenspace alignment condition

$$p_j[\tau - (1 + \sigma^2)p_j] < \frac{1}{2}[\alpha_p(1 + \sigma^2)s_j + 3\eta]$$



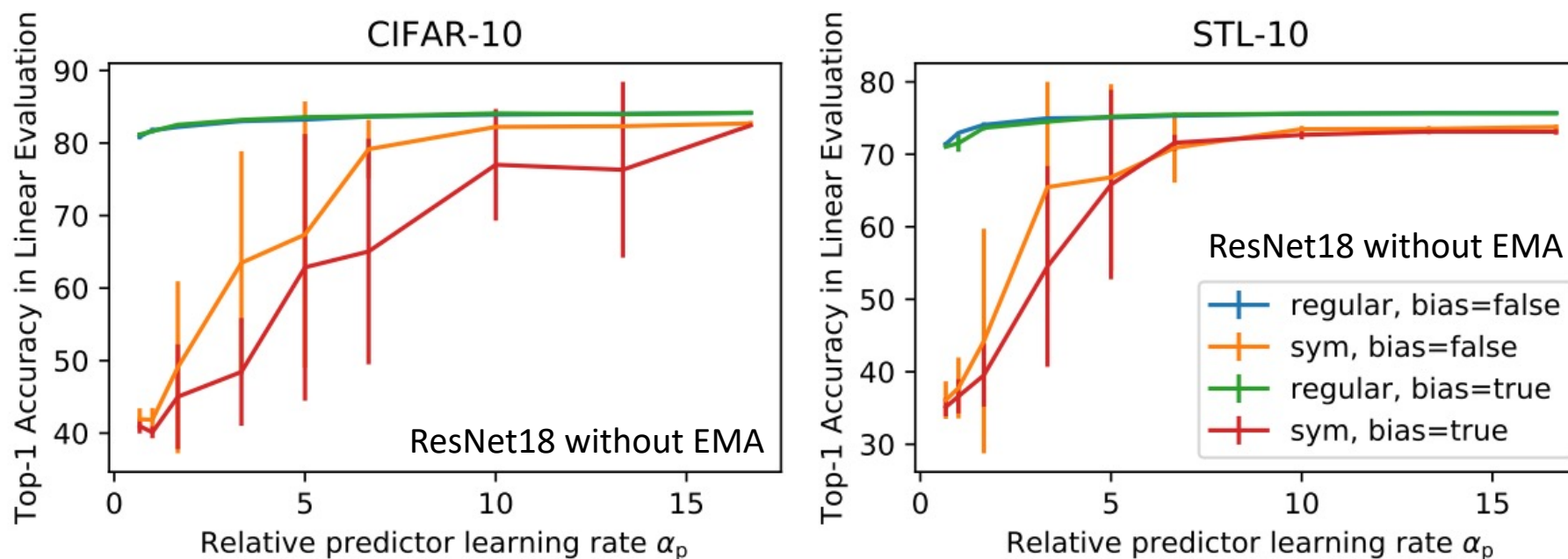
**Higher weight decay → alignment condition is more likely to satisfy!**



# Relative learning rate of the predictor $\alpha_p$

## Positive 😊

1. Large  $\alpha_p$  shrinks the size of trivial basin
2. Relax the condition of eigenspace alignment



# Exponential Moving Average rate $\beta$

$\beta$  large  $\rightarrow W_a(t)$  catches  $W(t)$  faster  $\rightarrow \tau$  grows faster to 1

**Positive** 😊: Slower rate (small  $\beta$ ) relaxes the condition of eigenspace alignment

**Negative** 😞: Slower rate makes the training slow and expands the size of trivial basin

## Part IV DirectPred

- Directly setting linear  $W_p$  rather than relying on gradient update.
  1. Estimate  $\hat{F} = \rho\hat{F} + (1 - \rho)E[\mathbf{f}\mathbf{f}^T]$
  2. Eigen-decompose  $\hat{F} = \hat{U}\Lambda_F\hat{U}^T$ ,  $\Lambda_F = \text{diag}[s_1, s_2, \dots, s_d]$
  3. Set  $W_p$  following the invariance:

$$p_j = \sqrt{s_j} + \epsilon \max_j s_j, \quad W_p = \hat{U} \text{diag}[p_j] \hat{U}^T$$

**Guaranteed Eigenspace Alignment 😊**

# Performance of DirectPred on STL-10/CIFAR-10

Downstream Classification Top-1	Number of epochs		
	100	300	500
<i>STL-10</i>			
<b>DirectPred</b>	<b>77.86 ± 0.16</b>	78.77 ± 0.97	78.86 ± 1.15
<b>DirectPred (freq=5)</b>	77.54 ± 0.11	<b>79.90 ± 0.66</b>	<b>80.28 ± 0.62</b>
SGD baseline	75.06 ± 0.52	75.25 ± 0.74	75.25 ± 0.74
<i>CIFAR-10</i>			
<b>DirectPred</b>	<b>85.21 ± 0.23</b>	<b>88.88 ± 0.15</b>	89.52 ± 0.04
<b>DirectPred (freq=5)</b>	84.93 ± 0.29	88.83 ± 0.10	<b>89.56 ± 0.13</b>
SGD baseline	84.49 ± 0.20	88.57 ± 0.15	89.33 ± 0.27

# Performance of DirectPred on ImageNet

Downstream classification (ImageNet):

BYOL variants	<i>Accuracy (60 ep)</i>		<i>Accuracy (300 ep)</i>	
	Top-1	Top-5	Top-1	Top-5
2-layer predictor <sup>*</sup>	<b>64.7</b>	<b>85.8</b>	<b>72.5</b>	90.8
linear predictor	59.4	82.3	69.9	89.6
<b>DirectPred</b>	64.4	<b>85.8</b>	72.4	<b>91.0</b>

<sup>\*</sup> 2-layer predictor is BYOL default setting.

DirectPred using linear predictor is better than SGD with linear predictor, and is comparable with 2-layer predictor.

# Conclusion

- A systematic analysis on the dynamics of non-contrastive self-supervised learning (SSL) methods
  - **Part I** Why we need (1) an **extra predictor** and (2) **stop-gradient**?
  - **Part II** Why training doesn't **collapse** to trivial solutions?
  - **Part III** The role played by different hyperparameters
- Propose **DirectPred**, a novel non-contrastive SSL method
  - Directly align the eigenspace of the predictor  $W_p$  with the correlation matrix  $F$
  - Comparable performance in downstream classification tasks, compared to vanilla BYOL
    - CIFAR-10/STL-10
    - ImageNet (60 epochs / 300 epochs)

Thanks!