# Self-supervised Learning: What we can learn from nonlinear dynamics and neuroscience

FIMI2025@Okinawa, Mar 1st 2025
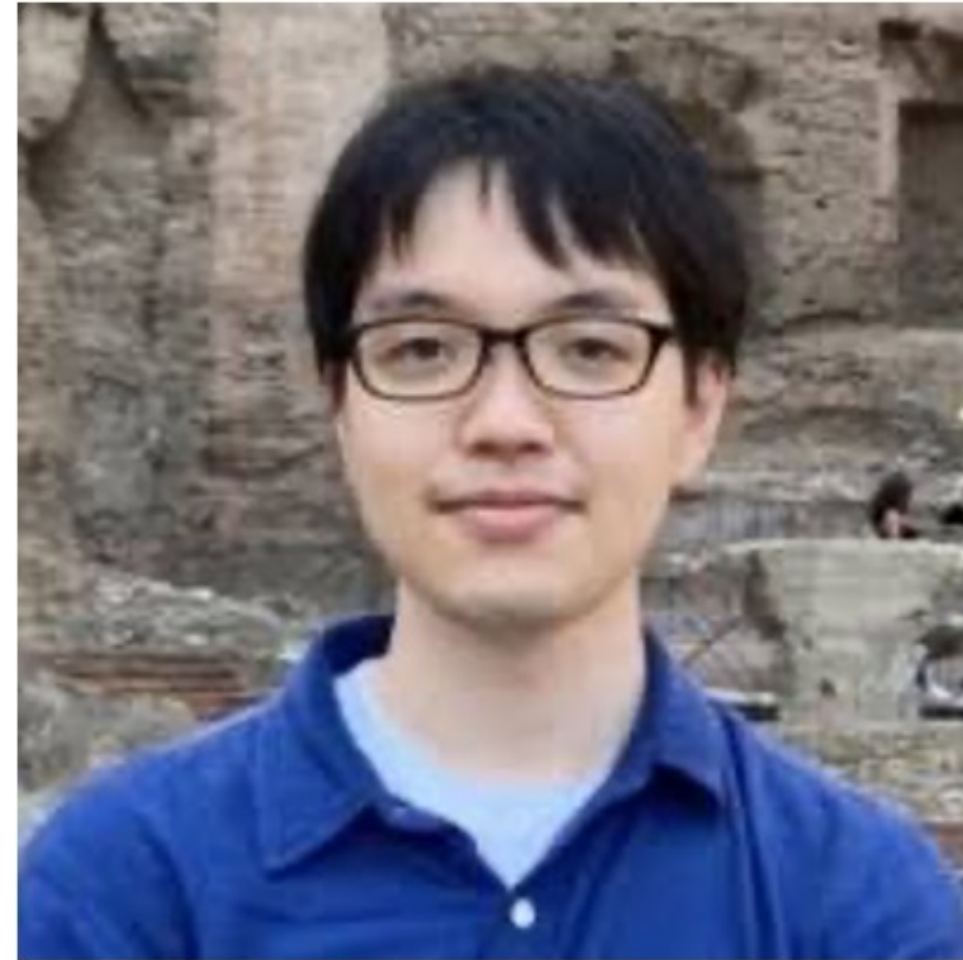
Han Bao

(Kyoto University → The Institute of Statistical Mathematics)

where most of the work has been done
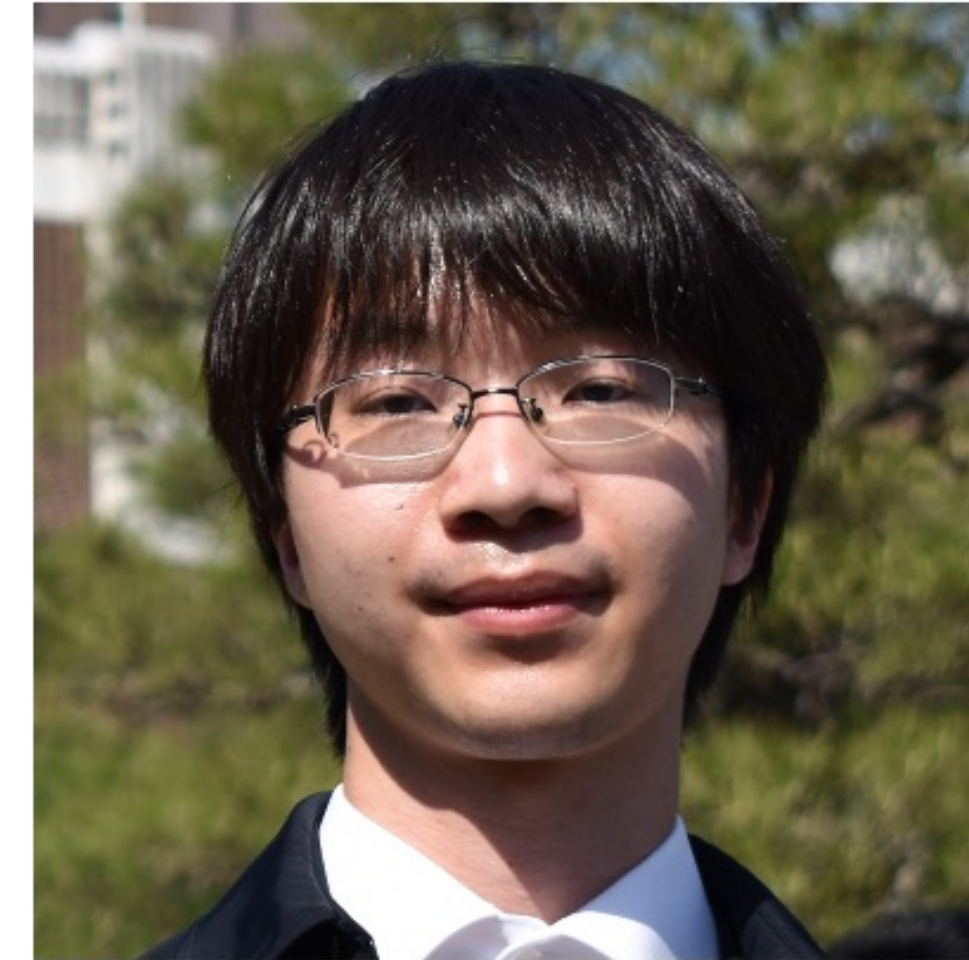
# This work was …

- Collaboration with great members!



Satoki Ishikawa
Science Tokyo



Makoto Yamada
OIST



Yuki Takezawa
Kyoto University

# Self-supervised Learning: What we can learn from **nonlinear dynamics** and **neuroscience**

## Part I
learning dynamics, stability, adaptivity, ⋯

## Part II
predictive coding, hippocampal model

Published in Advances in Neural Information Processing Systems 6 (NIPS **1993**)

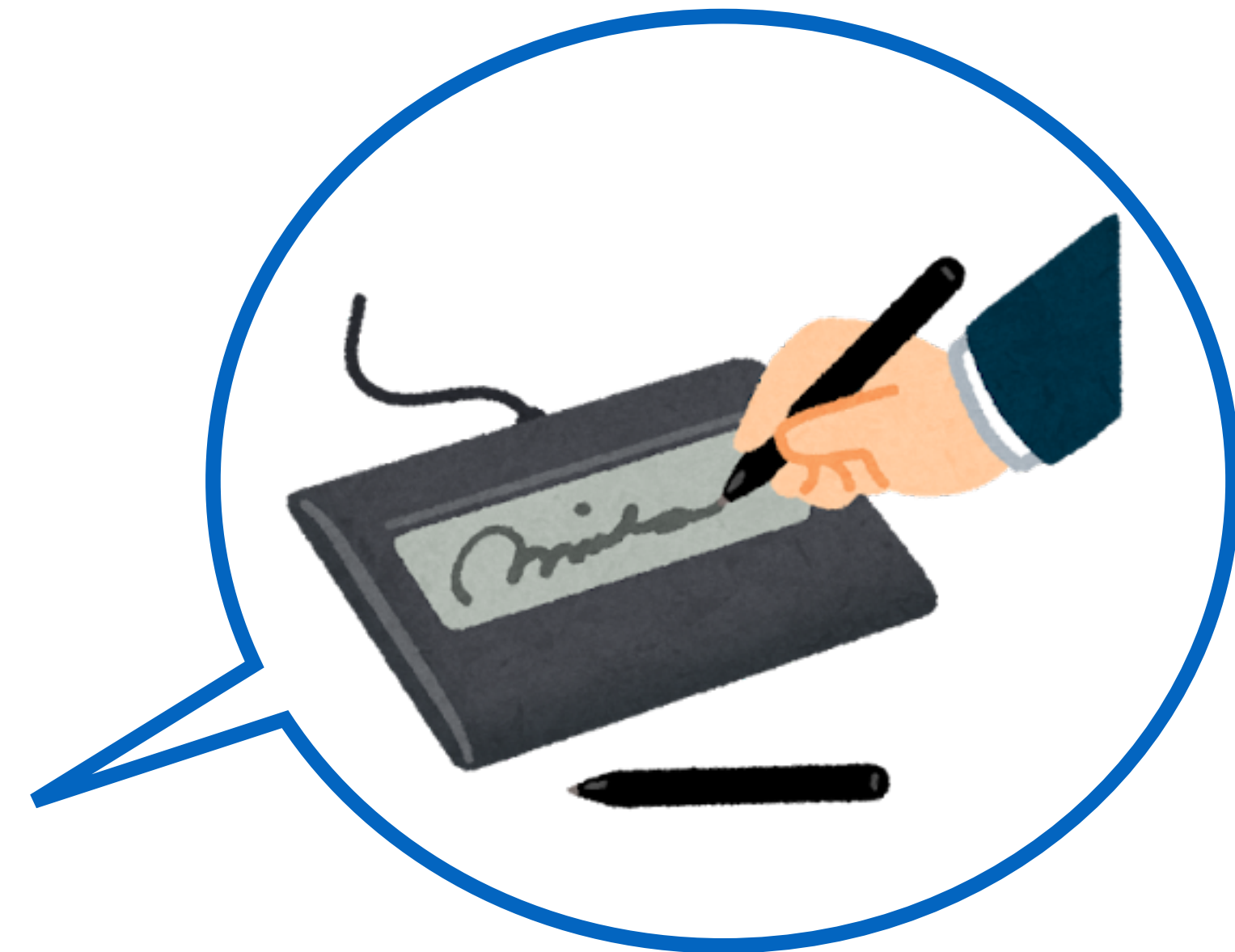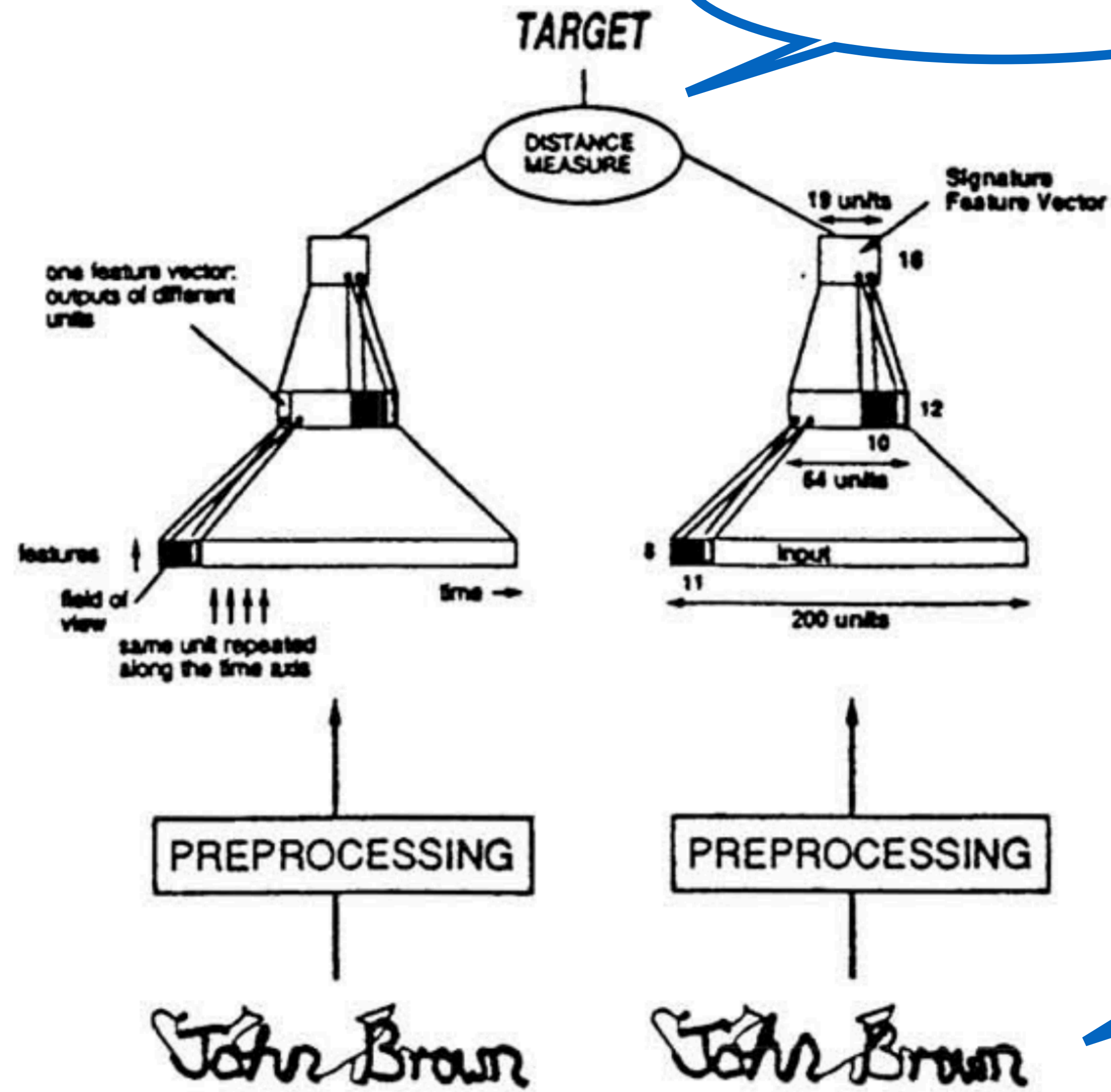## Signature Verification using a "Siamese" Time Delay Neural Network

Jane Bromley, Isabelle Guyon, Yann LeCun,
Eduard Säckinger and Roopak Shah
AT&T Bell Laboratories
Holmdel, NJ 07733
jbromley@big.att.com

# Once upon a time …



Cosine distance

**Learning a Similarity Metric Discriminatively, with Application to Face Verification**

Sumit Chopra          Raia Hadsell          Yann LeCun

Courant Institute of Mathematical Sciences
New York University
New York, NY, USA

2010: unsupervised, density estimation

**Noise-contrastive estimation: A new estimation principle for unnormalized statistical models**

**Michael Gutmann**
Dept of Computer Science
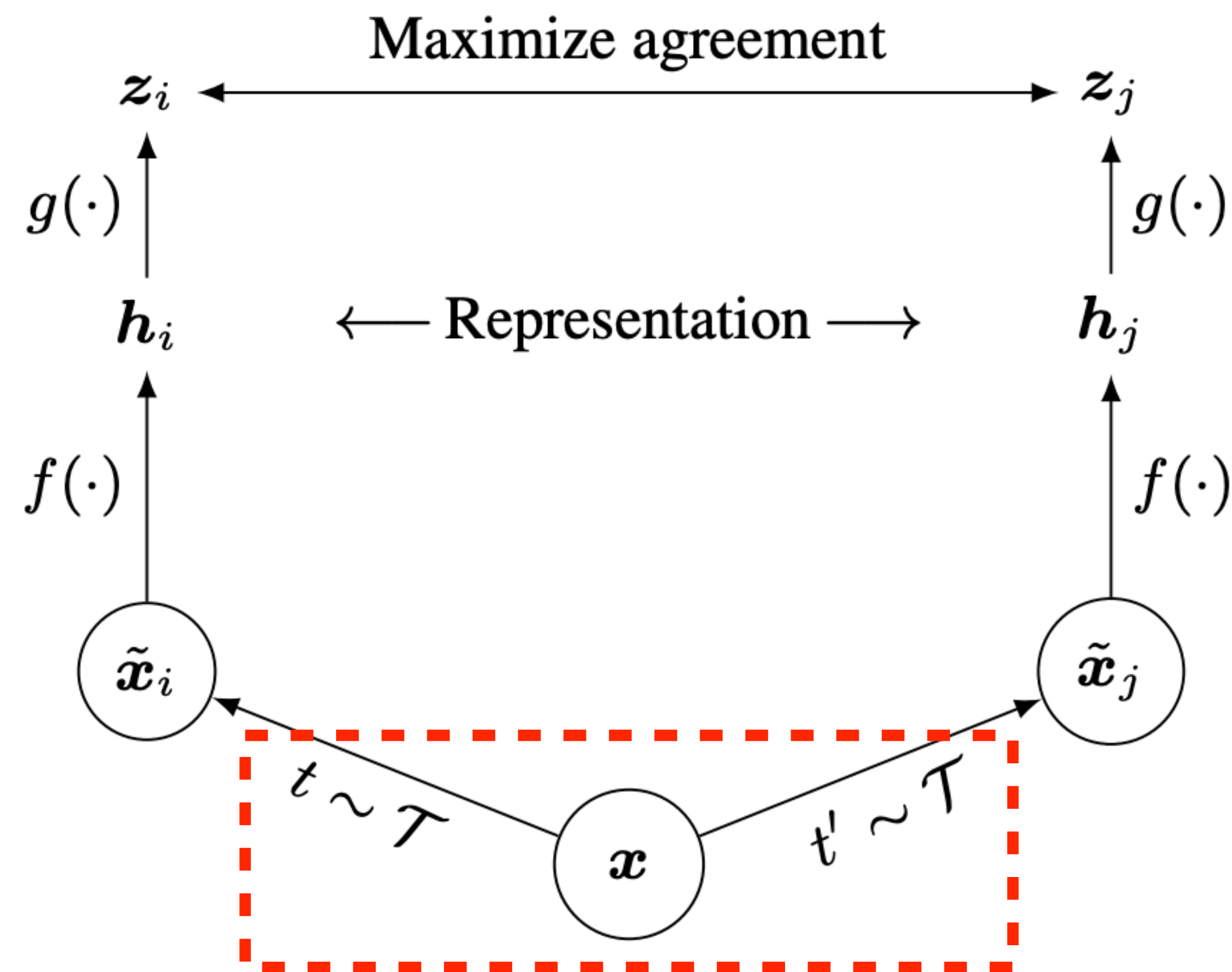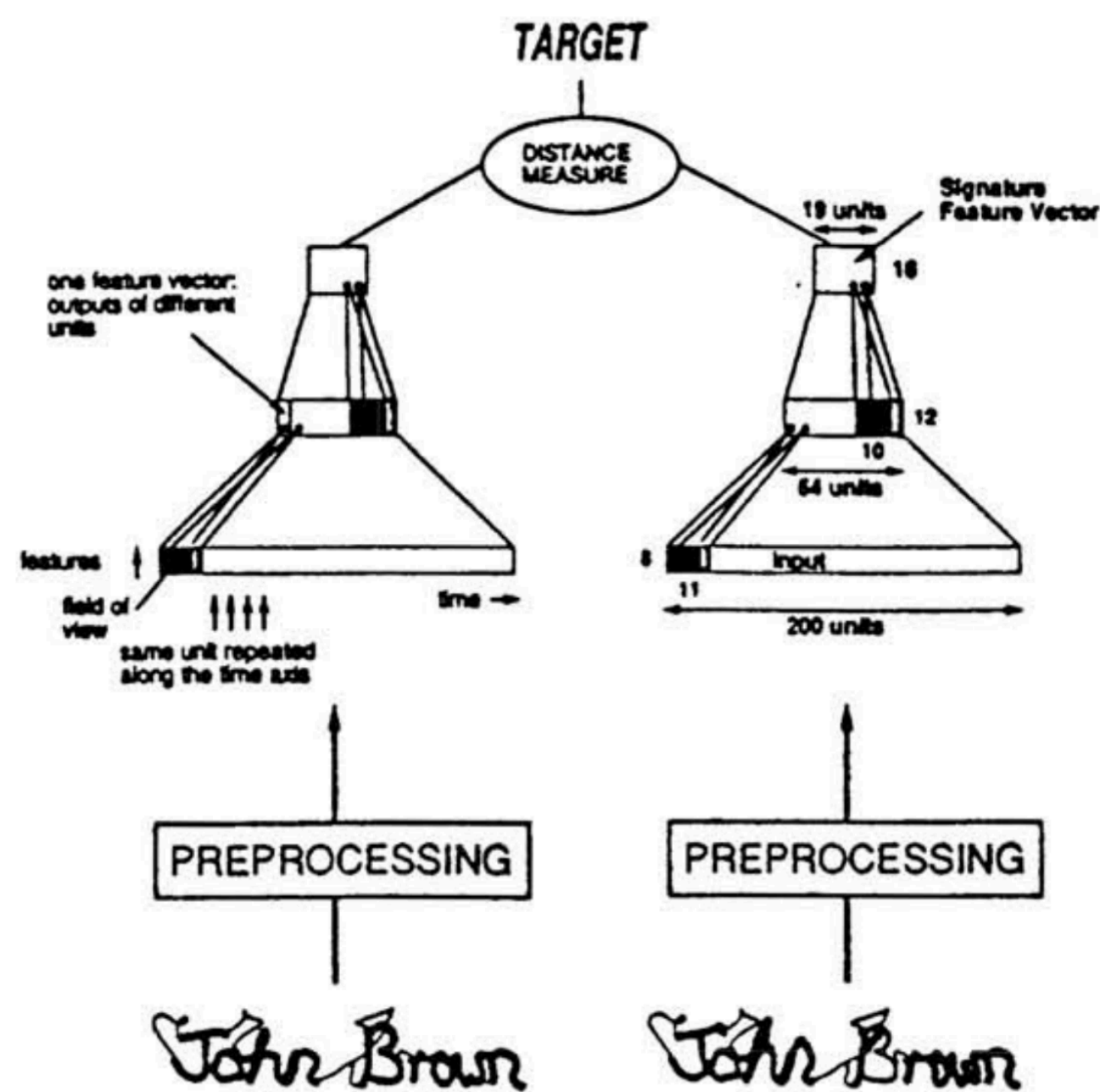and HIIT, University of Helsinki

**Aapo Hyvärinen**
Dept of Mathematics & Statistics, Dept of Computer
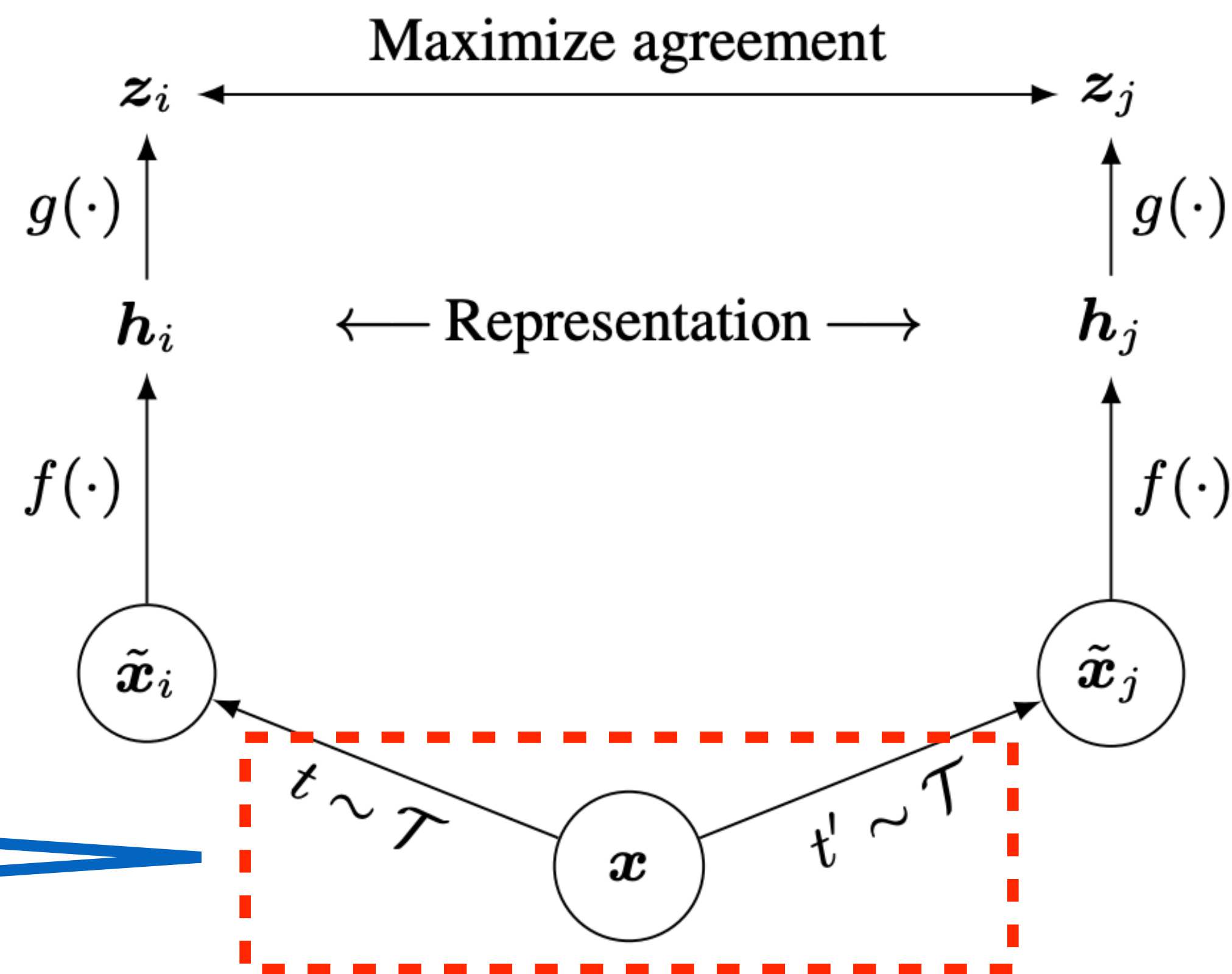Science and HIIT, University of Helsinki

# From supervised to unsupervised

1993: supervised

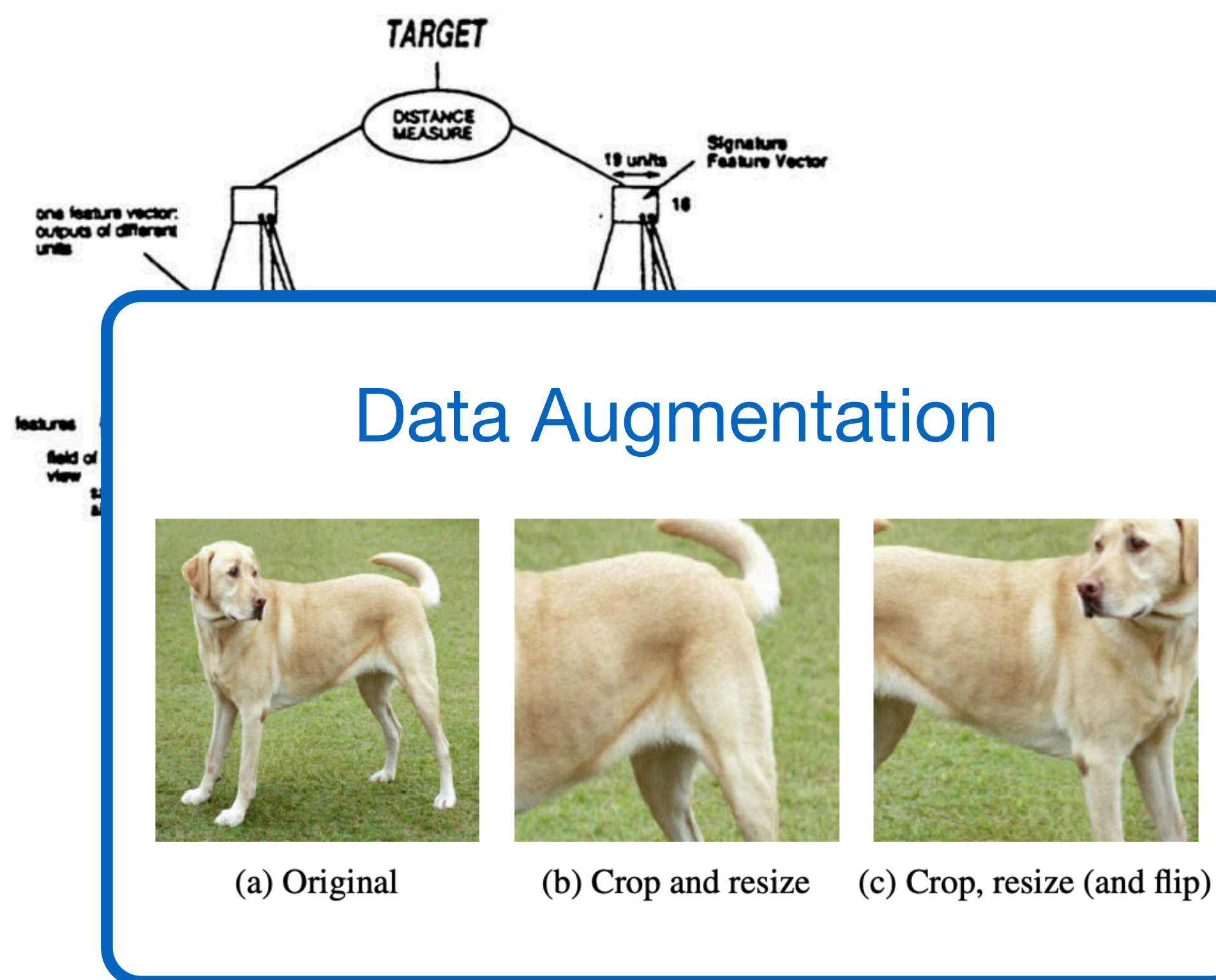2020: **un**supervised
SimCLR [Chen+ 20]

# From supervised to unsupervised
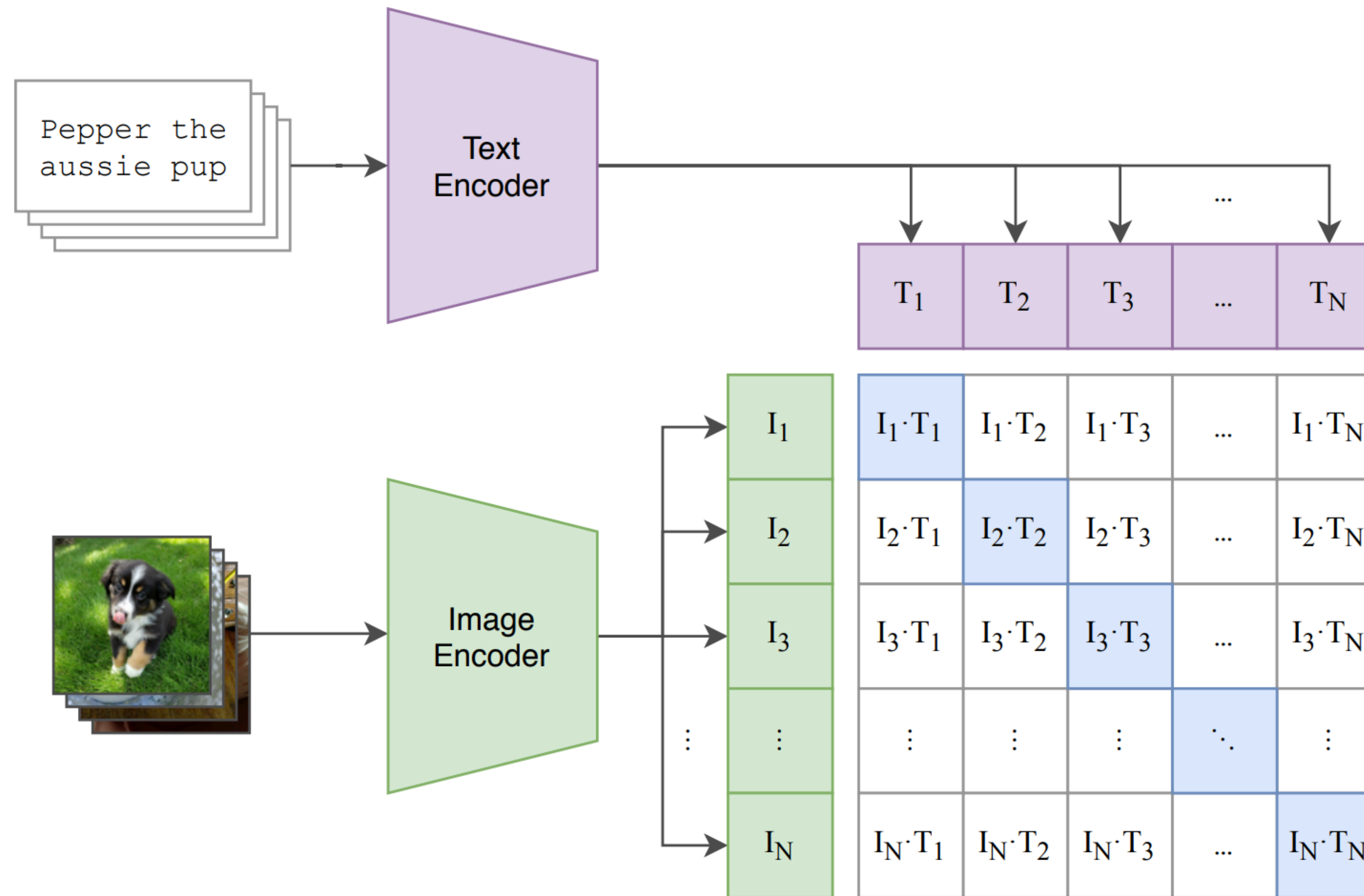
1993: supervised

2020: **un**supervised
SimCLR [Chen+ 20]



Data Augmentation

(a) Original    (b) Crop and resize    (c) Crop, resize (and flip)

Maximize agreement

$z_i \longleftrightarrow z_j$

$g(\cdot)$            $g(\cdot)$

$h_i \longleftarrow$ Representation $\longrightarrow h_j$

$f(\cdot)$            $f(\cdot)$

$\tilde{x}_i$            $\tilde{x}_j$

$t \sim \mathcal{T}$     $x$     $t' \sim \mathcal{T}$

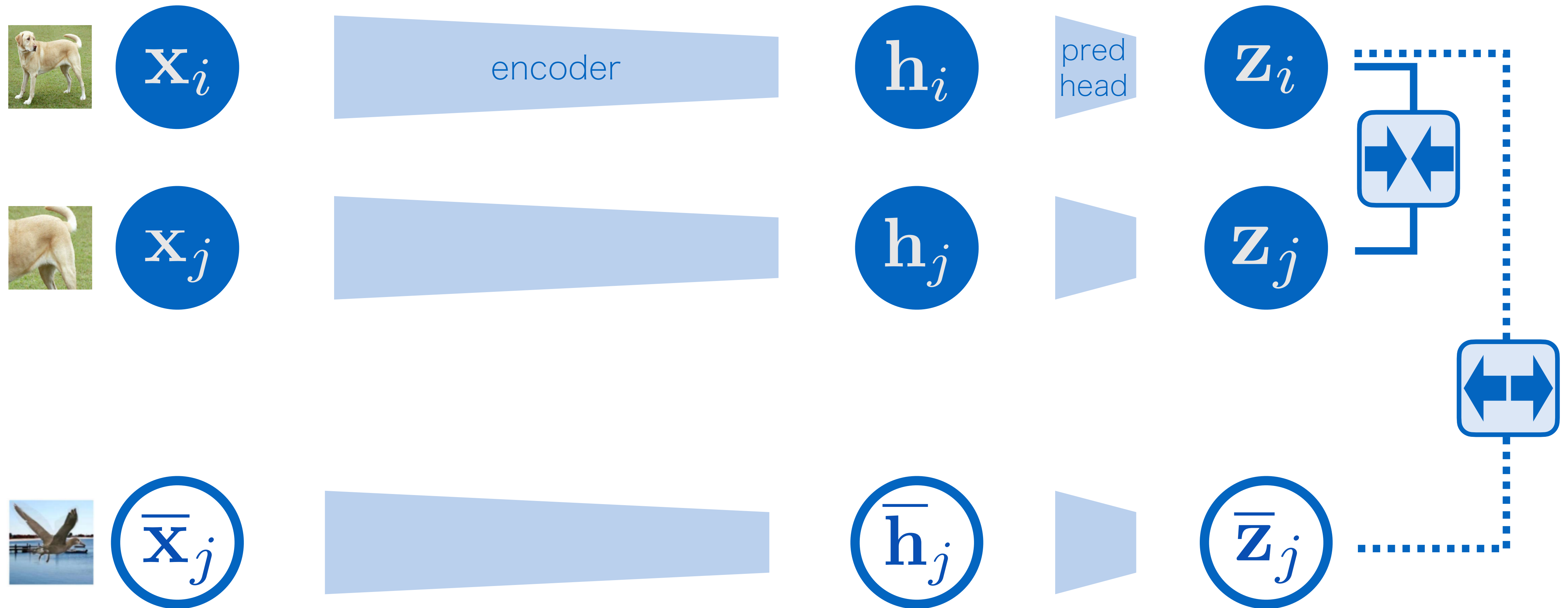Chen et al. (2020) A Simple Framework for Contrastive Learning of Visual Representations

# CLIP: multi-view representation learning



(1) Contrastive pre-training

# Massive negative sampling

SimCLR [Chen+ 20]



$\mathbf{x}_i$ — encoder — $\mathbf{h}_i$ — pred head — $\mathbf{z}_i$

$\mathbf{x}_j$ — encoder — $\mathbf{h}_j$ — $\mathbf{z}_j$

$\overline{\mathbf{x}}_j$ — encoder — $\overline{\mathbf{h}}_j$ — $\overline{\mathbf{z}}_j$

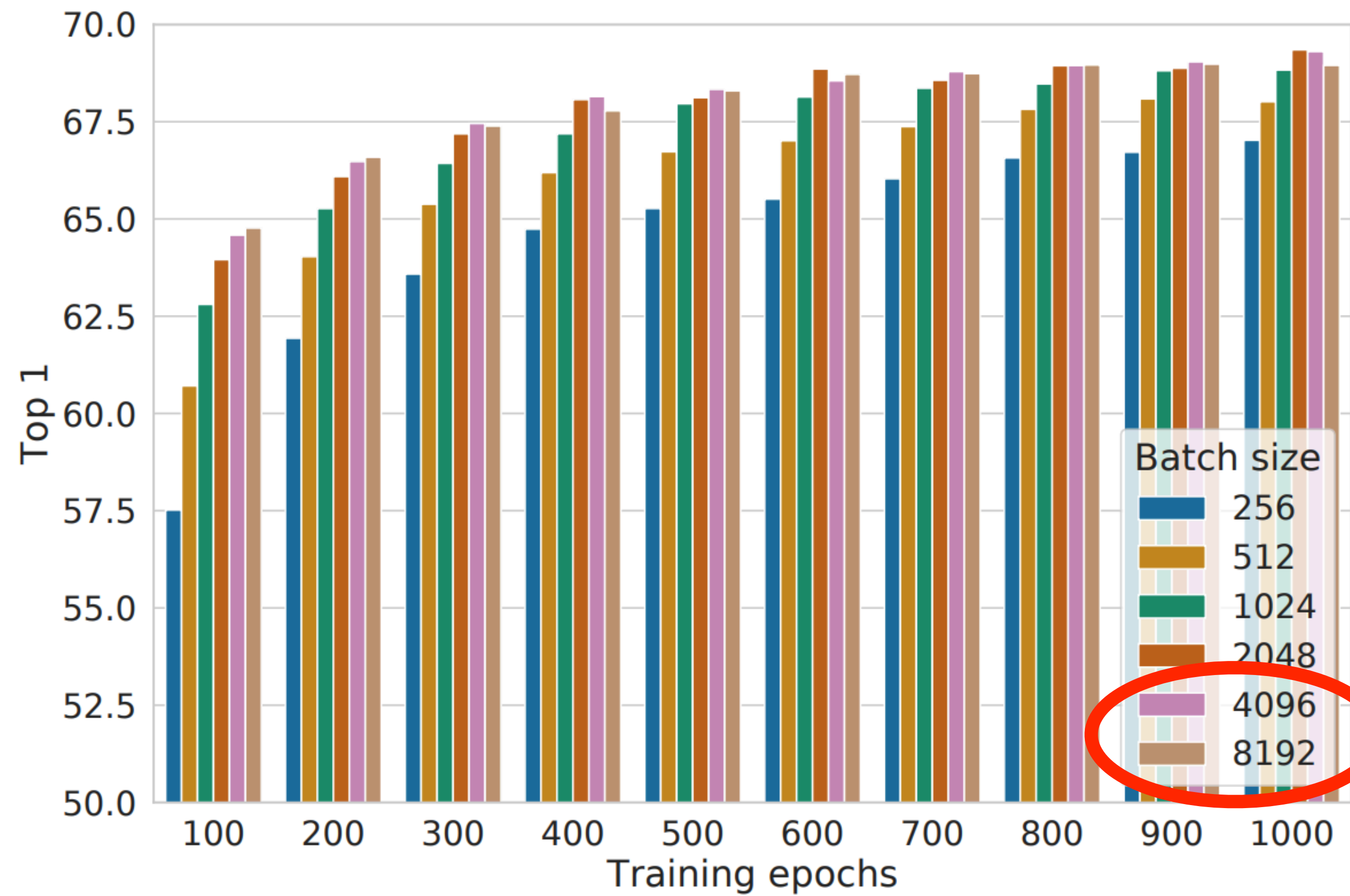Chen et al. (2020) A Simple Framework for Contrastive Learning of Visual Representations
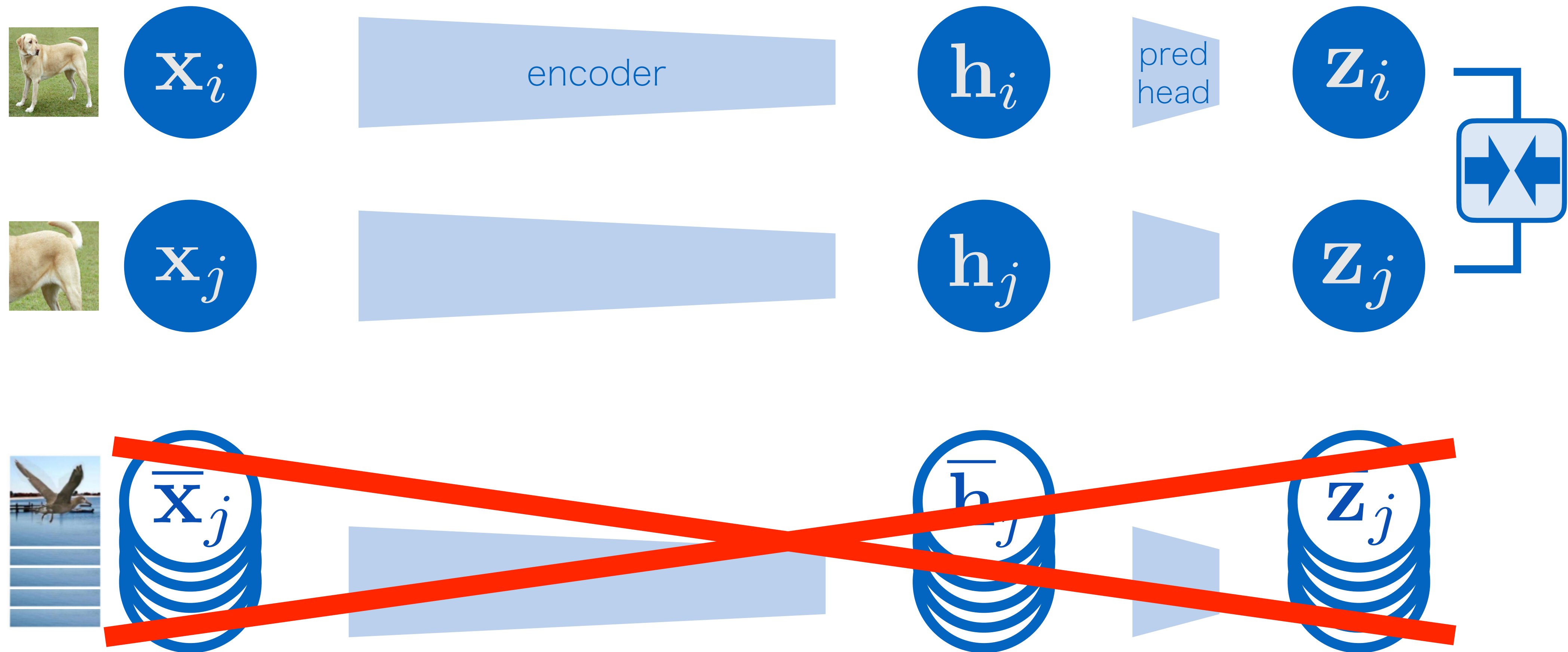
# Massive negative sampling
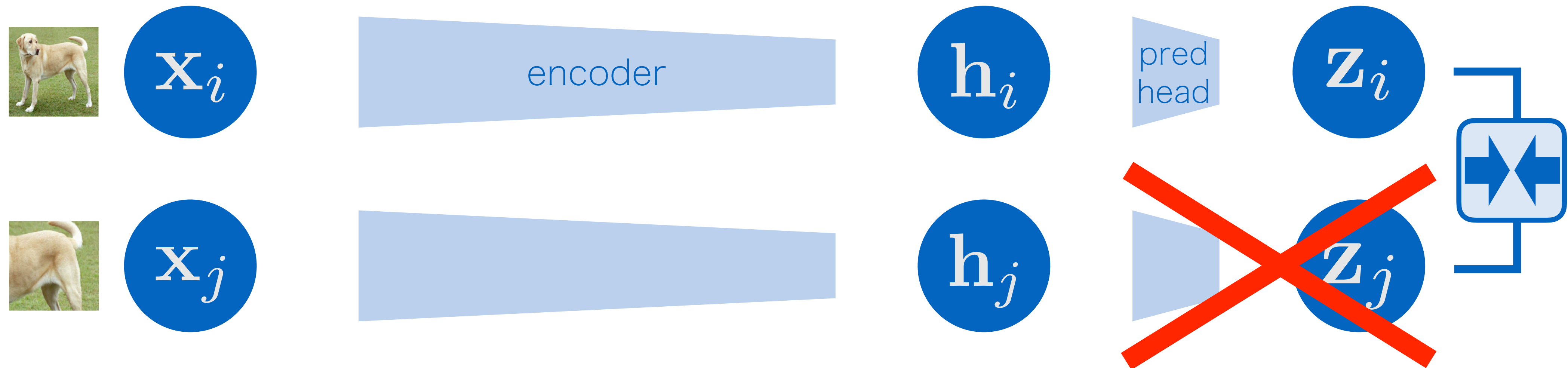
SimCLR [Chen+ 20]

# From contrastive to NON-contrastive
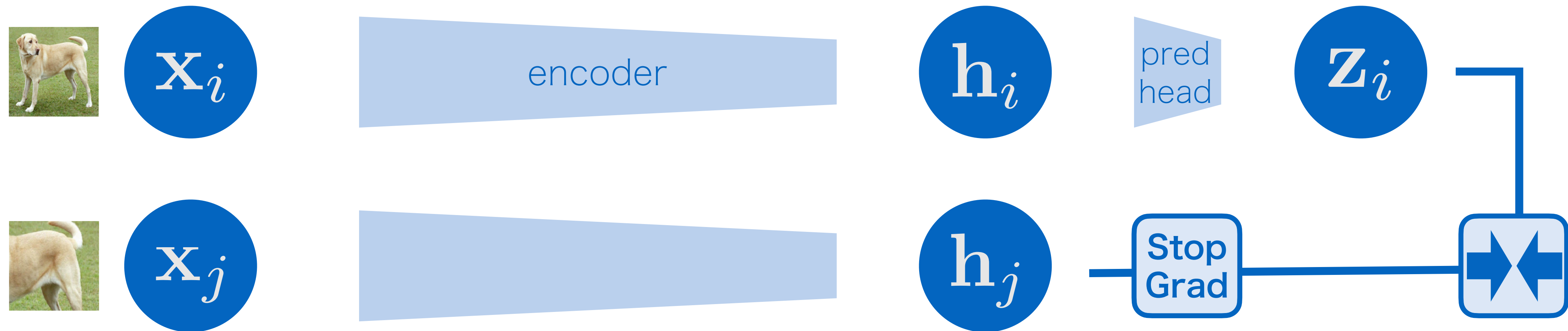
SimSiam [Chen-He 21]

# From contrastive to NON-contrastive

SimSiam [Chen-He 21]
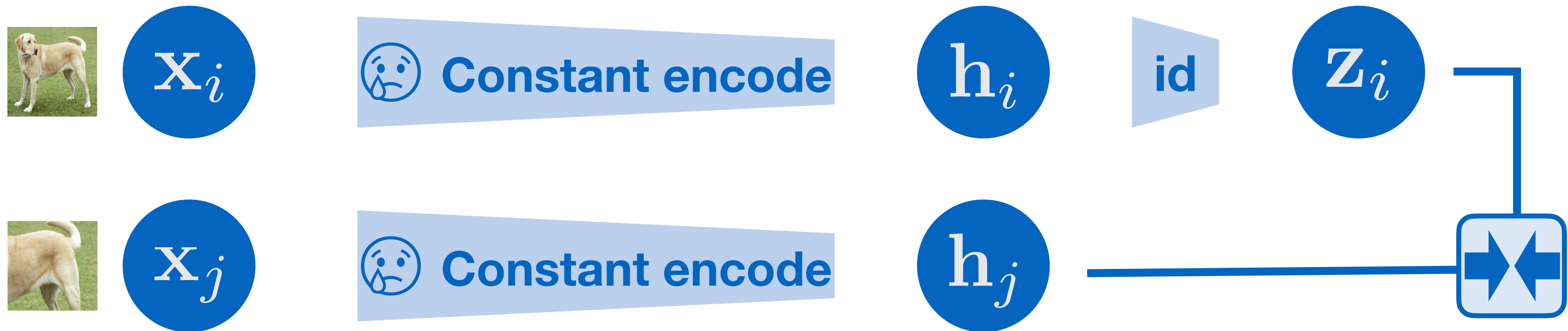
# From contrastive to NON-contrastive

SimSiam [Chen-He 21]



- Data augmentation
- Prediction head (but at anchor side only!)
- Stop gradient

Chen & He (2021) Exploring Simple Siamese Representation Learning.

SimSiam [Chen-He 21]

$\mathbf{x}_i$ — 😢 **Constant encode** — $\mathbf{h}_i$ — **id** — $\mathbf{z}_i$

$\mathbf{x}_j$ — 😢 **Constant encode** — $\mathbf{h}_j$

- Data augmentation

- Prediction head (but at anchor side only!)

- Stop gradient

🤔 **How to avoid *constant* encoder?**
trivial

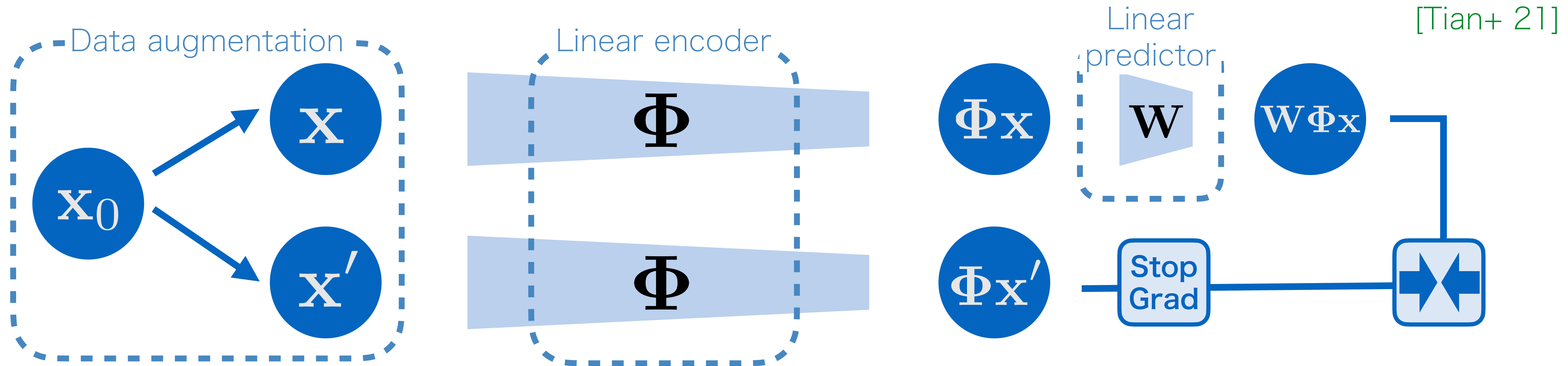Chen & He (2021) Exploring Simple Siamese Representation Learning.
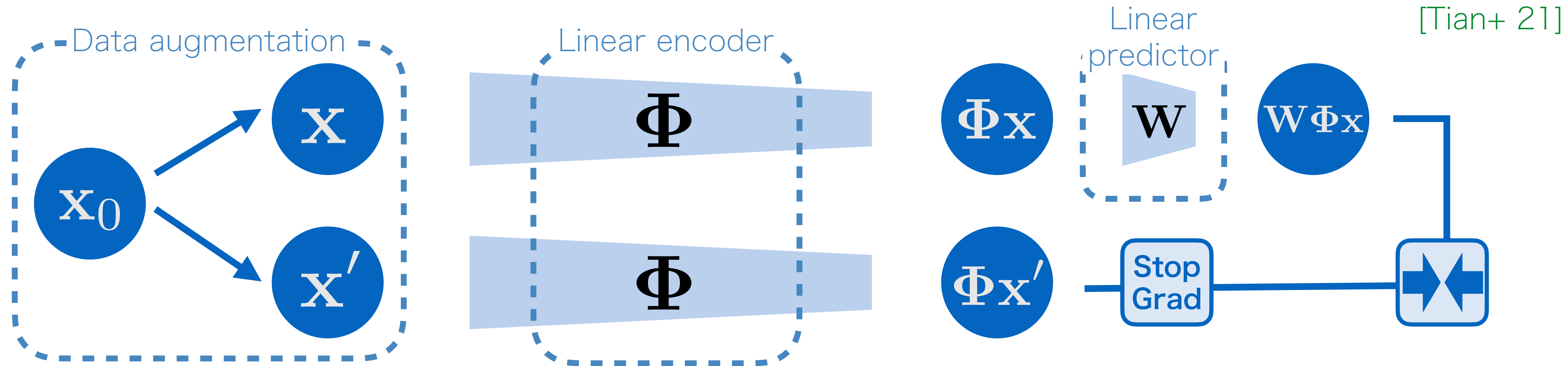
# What we can learn
# from **nonlinear dynamics** and neuroscience

**Bao, H.** (2023)
Feature Normalization Prevents Collapse of Non-contrastive Learning Dynamics.

# Theoretical model of SimSiam

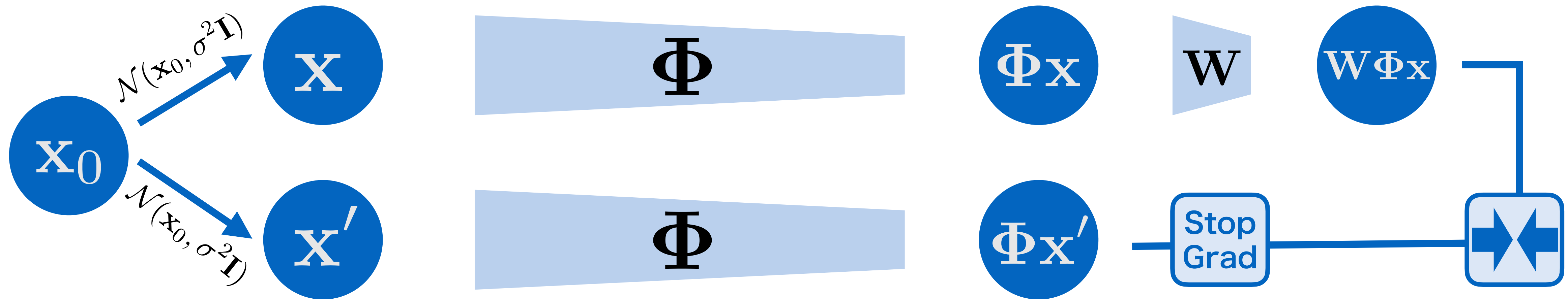Tian et al. (2021) Understanding self-supervised learning dynamics without contrastive pairs.

# Theoretical model of SimSiam

Data augmentation

Linear encoder

Linear predictor

[Tian+ 21]

$\mathbf{x}_0$   $\mathbf{x}$   $\mathbf{x}'$

$\Phi$

$\Phi$

$\Phi\mathbf{x}$   $\mathbf{W}$   $\mathbf{W}\Phi\mathbf{x}$

$\Phi\mathbf{x}'$   Stop Grad

$$\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I})$$

$$\mathbf{x}, \mathbf{x}' \sim \mathcal{N}(\mathbf{x}_0, \sigma^2 \mathbf{I})$$

Strength of data aug

Tian et al. (2021) Understanding self-supervised learning dynamics without contrastive pairs.

[Tian+ 21]



$$\mathcal{L}(\mathbf{\Phi}, \mathbf{W}) = \frac{1}{2}\mathbb{E}\|\mathbf{W}\mathbf{\Phi}\mathbf{x} - \mathrm{StopGrad}(\mathbf{\Phi}\mathbf{x}')\|_2^2$$

Tian et al. (2021) Understanding self-supervised learning dynamics without contrastive pairs.

# Learning dynamics

$$\mathcal{L}(\boldsymbol{\Phi}, \mathbf{W}) = \frac{1}{2}\mathbb{E}\|\mathbf{W}\boldsymbol{\Phi}\mathbf{x} - \mathrm{StopGrad}(\boldsymbol{\Phi}\mathbf{x}')\|_2^2$$

weight decay

**Discrete time**
gradient descent

$$\begin{aligned} \boldsymbol{\Phi}(t+1) &= \boldsymbol{\Phi}(t) - \eta(\nabla_{\boldsymbol{\Phi}}\mathcal{L} + \rho\boldsymbol{\Phi}(t)) \\ \mathbf{W}(t+1) &= \mathbf{W}(t) - \eta(\nabla_{\mathbf{W}}\mathcal{L} + \rho\mathbf{W}(t)) \end{aligned}$$

$\eta \to 0$

**Continuous time**
gradient flow

$$\begin{aligned} \dot{\boldsymbol{\Phi}} &= -\nabla_{\boldsymbol{\Phi}}\mathcal{L} - \rho\boldsymbol{\Phi} \\ \dot{\mathbf{W}} &= -\nabla_{\mathbf{W}}\mathcal{L} - \rho\mathbf{W} \end{aligned}$$

Tian et al. (2021) Understanding self-supervised learning dynamics without contrastive pairs.

[Tian+ 21]

**Matrix dynamics**: not easy to deal with 🥲

$$\dot{\mathbf{\Phi}} = -\nabla_{\mathbf{\Phi}}\mathcal{L} - \rho\mathbf{\Phi}$$
$$\dot{\mathbf{W}} = -\nabla_{\mathbf{W}}\mathcal{L} - \rho\mathbf{W}$$

$\xrightarrow{\mathbf{\Phi}\mathbf{\Phi}^{\top} \equiv \mathbf{F}}$

$$\dot{\mathbf{F}} = -2(1+\sigma^2)\mathbf{W}^2\mathbf{F} + 2\mathbf{W}\mathbf{F} - 2\rho\mathbf{F}$$
$$\dot{\mathbf{W}} = -(1+\sigma^2)\mathbf{W}\mathbf{F} + \mathbf{F} - \rho\mathbf{W}$$

$$\begin{cases} s: & j\text{-th eigval of } \mathbf{F} \\ p: & j\text{-th eigval of } \mathbf{W} \end{cases}$$

**Scalar dynamics**:

enabled by eigendecomposition 😁

$$\dot{s} = -2(1+\sigma^2)p^2s + 2ps - 2\rho s$$
$$\dot{p} = -(1+\sigma^2)ps + s - \rho p$$

[Tian+ 21]



Eigvals of matrices evolves as follows:

$$\dot{s} = -2(1+\sigma^2)p^2s + 2ps - 2\rho s$$

$$\dot{p} = -(1+\sigma^2)ps + s - \rho p$$

Simultaneous ODE

$$\dot{s} = -2(1+\sigma^2)p^2 s + 2ps - 2\rho s$$

$$\dot{p} = -(1+\sigma^2)ps + s - \rho p$$

Adiabatic elimination



Eigval ODE of predictor

$$\dot{p} = p^2\{1 - (1+\sigma^2)p\} - \rho p$$

- Two params: $\sigma^2$ (data aug) & $\rho$ (weight decay)
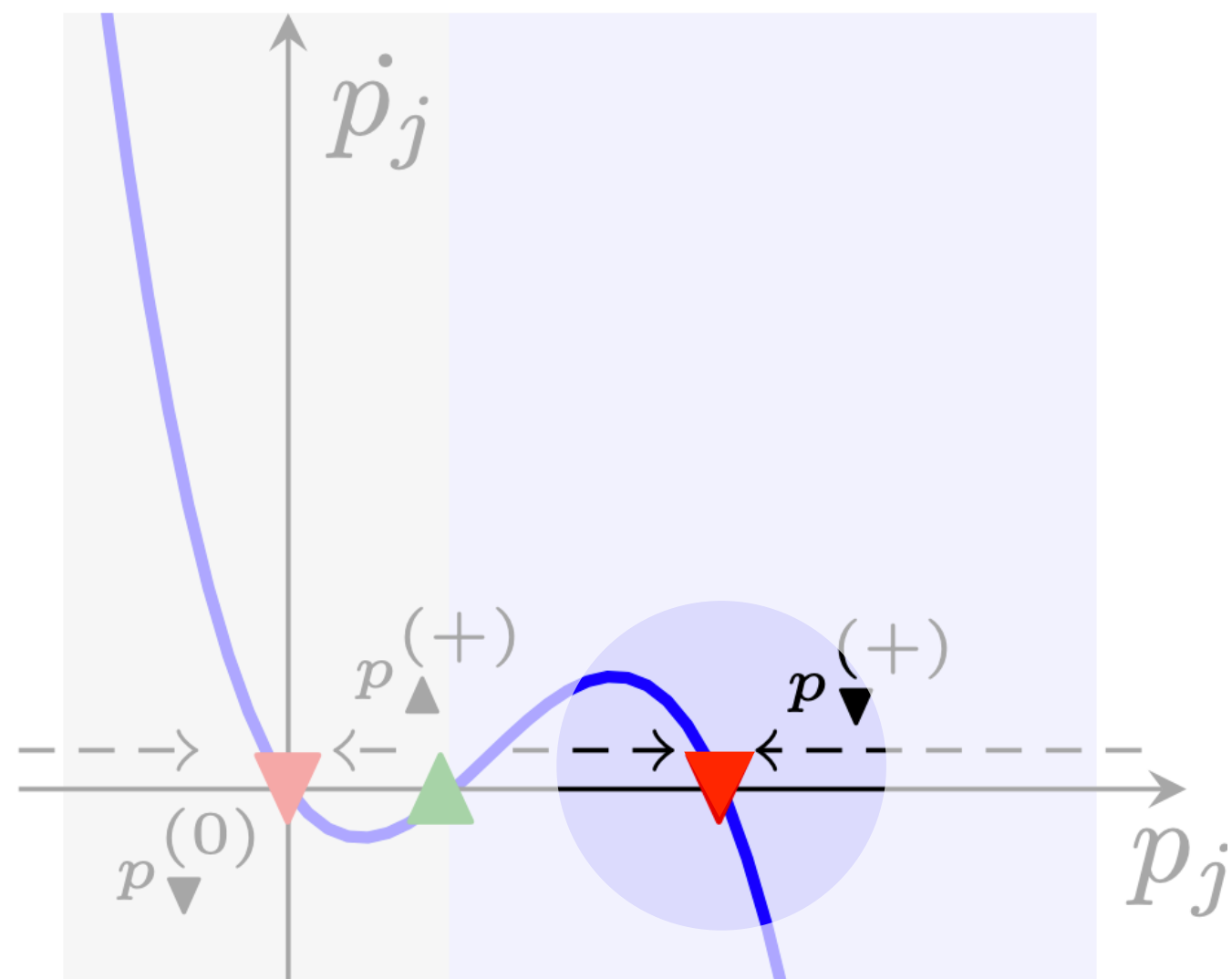
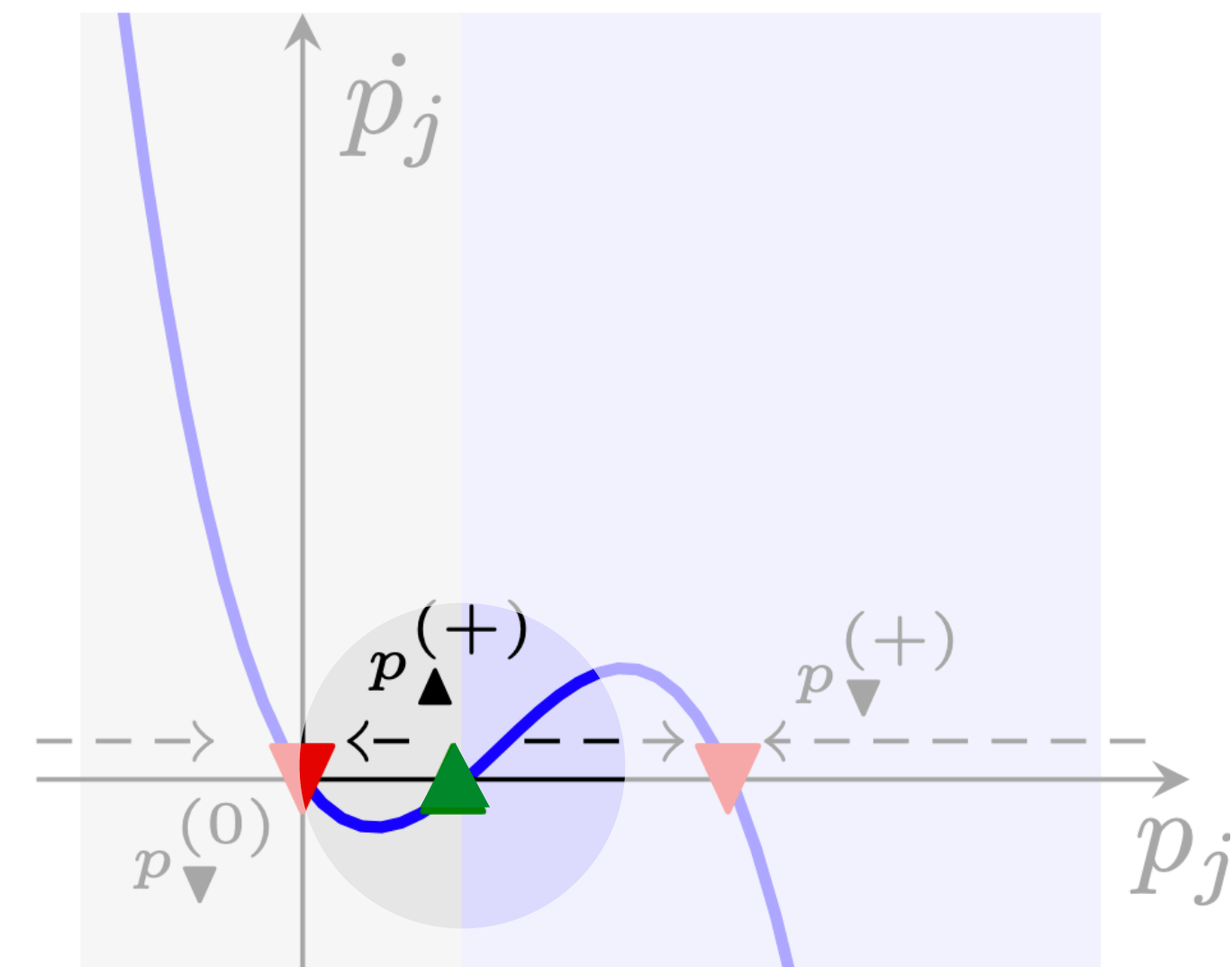- **Q.** How to avoid constant predictor?

- **Q.** How to avoid $p = 0$ ?

Tian et al. (2021) Understanding self-supervised learning dynamics without contrastive pairs.

# Quick pre-requisite

- Stability analysis of ODE $\dot{p} = f(p)$

- $\dot{p} = 0$ is equilibrium (but can be unstable)

- If $f(p) < 0$: **stable**

- If $f(p) > 0$: **unstable**
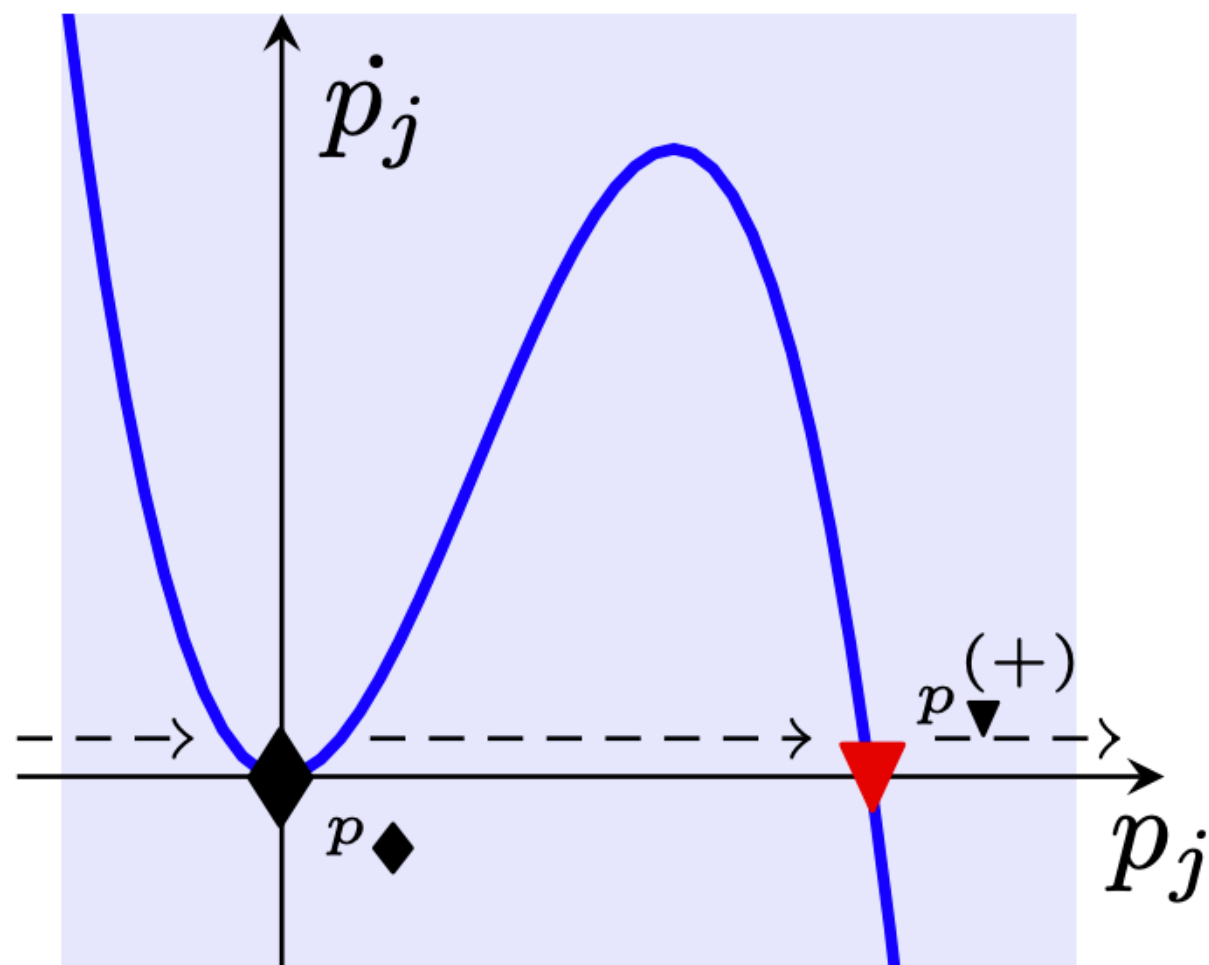


because $f(p)$ head for the equilibrium locally

# Bifurcation: too strong weight decay collapses 😢

Eigval ODE of projector

$$\dot{p} = p^2 \{1 - (1 + \sigma^2)p\} - \rho p$$

Case (a): $\rho = 0$
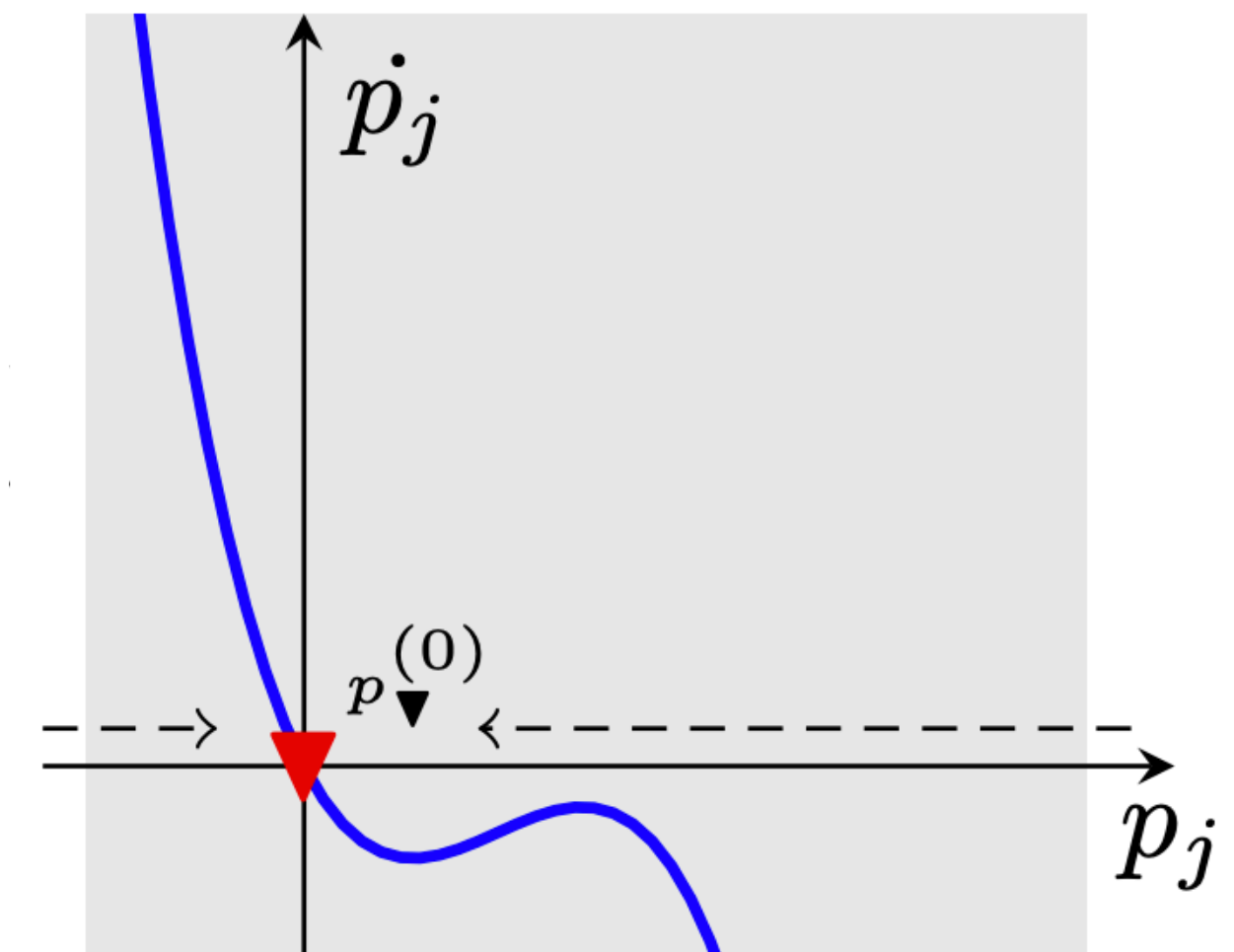
[Tian+ 21]

Eigval ODE of projector

$$\dot{p} = p^2\{1 - (1 + \sigma^2)p\} - \rho p$$



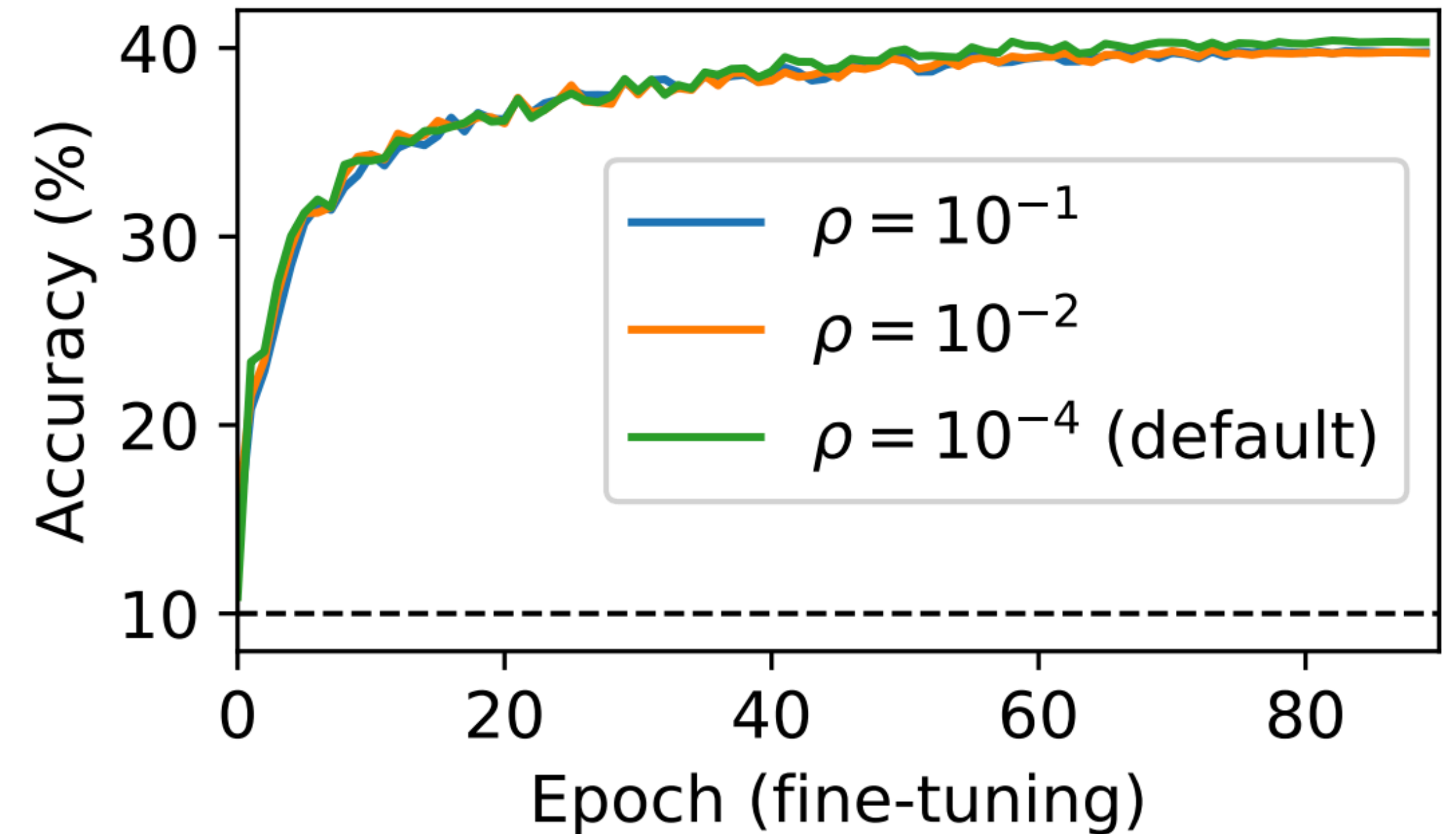Case (a): $\rho = 0$

Case (b): $\rho < \frac{1}{4(1+\sigma^2)}$

[Tian+ 21]

Eigval ODE of projector

$$\dot{p} = p^2\{1 - (1+\sigma^2)p\} - \rho p$$



Case (a): $\rho = 0$

Case (b): $\rho < \frac{1}{4(1+\sigma^2)}$

Case (c): $\rho > \frac{1}{4(1+\sigma^2)}$

strong weight decay:
trivial solution $p = 0$ only

Tian et al. (2021) Understanding self-supervised learning dynamics without contrastive pairs.

# But is this really happening?

- Pilot study: SimSiam on CIFAR-10
  - ❖ evaluation: linear probing accuracy
- [Chen-He 21] Let's use small enough WD!
- [Bao 23] Intensifying WD keeps working
- 🔑 small learning rate
  (to enter gradient flow regime)



larger WD $\rho$ still works

= accuracy does not breaks down

[Tian+ 21]

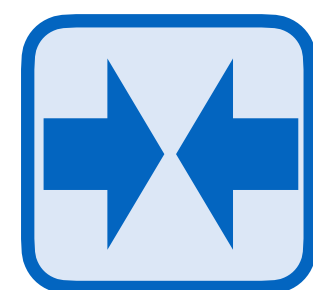$$\mathcal{L}(\mathbf{\Phi}, \mathbf{W}) = \frac{1}{2}\mathbb{E}\|\mathbf{W}\mathbf{\Phi}\mathbf{x} - \mathrm{StopGrad}(\mathbf{\Phi}\mathbf{x}')\|_2^2$$

**SimSiam impl**

$$\mathcal{L}(\mathbf{\Phi}, \mathbf{W}) = \mathbb{E}\left[-\frac{\langle\mathbf{W}\mathbf{\Phi}\mathbf{x}, \mathrm{StopGrad}(\mathbf{\Phi}\mathbf{x}')\rangle}{\|\mathbf{W}\mathbf{\Phi}\mathbf{x}\|\|\mathrm{StopGrad}(\mathbf{\Phi}\mathbf{x}')\|}\right]$$

# Cosine loss may prevent collapse

- If collapsing $(p = 0)$, predictor goes to zero $\mathbf{W} = \mathbf{O}$, blowing up cosine loss



$$\mathcal{L}(\boldsymbol{\Phi}, \mathbf{W}) = \mathbb{E}\left[-\frac{\langle \mathbf{W}\boldsymbol{\Phi}\mathbf{x}, \mathrm{StopGrad}(\boldsymbol{\Phi}\mathbf{x}')\rangle}{\|\mathbf{W}\boldsymbol{\Phi}\mathbf{x}\|\|\mathrm{StopGrad}(\boldsymbol{\Phi}\mathbf{x}')\|}\right]$$

🤔 **What does the cosine-loss dynamics look like?**

# Challenges of cosine loss: normalization

$$\mathcal{L}(\mathbf{\Phi}, \mathbf{W}) = \mathbb{E}\left[-\frac{\langle \mathbf{W}\mathbf{\Phi}\mathbf{x}, \mathrm{StopGrad}(\mathbf{\Phi}\mathbf{x}')\rangle}{\|\mathbf{W}\mathbf{\Phi}\mathbf{x}\|\|\mathrm{StopGrad}(\mathbf{\Phi}\mathbf{x}')\|}\right]$$

- Taking derivative wrt normalizer makes gradient complicated
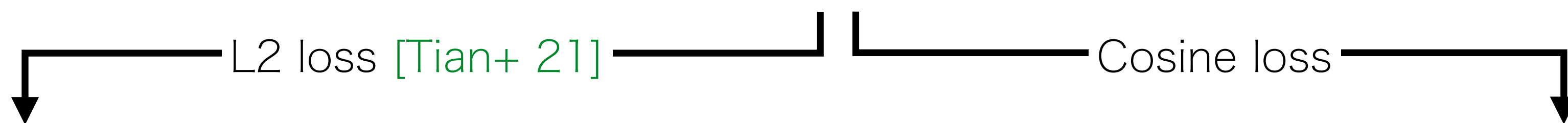
- Solution: **high-dimensional limit**



norm of random vector concentrates on a hypersphere

Wait, let me use LaTeX for the fraction.

$$\mathcal{L}(\mathbf{\Phi}, \mathbf{W}) = \mathbb{E}\left[-\frac{\langle \mathbf{W\Phi x}, \text{StopGrad}(\mathbf{\Phi x'})\rangle}{\|\mathbf{W\Phi x}\|\|\text{StopGrad}(\mathbf{\Phi x'})\|}\right]$$

- Taking derivative wrt normalizer makes gradient complicated

- Solution: **high-dimensional limit**

[Chen-He 21]

- *Prediction MLP.* The prediction MLP ($h$) has BN applied to its hidden fc layers. Its output fc does not have BN (ablation in Sec. 4.4) or ReLU. This MLP has 2 layers. The dimension of $h$'s input and output ($z$ and $p$) is $d = $ 2048, and $h$'s hidden layer's dimension is 512, making $h$ a bottleneck structure (ablation in supplement).



dim=1024

🔑 **Approx** $\|\mathbf{W\Phi x}\| = \text{const.}$

$$\dot{\mathbf{\Phi}} = -\nabla_{\mathbf{\Phi}}\mathcal{L} - \rho\mathbf{\Phi}$$

$$\dot{\mathbf{W}} = -\nabla_{\mathbf{W}}\mathcal{L} - \rho\mathbf{W}$$

L2 loss [Tian+ 21]

Cosine loss

$$\dot{s} = -2(1 + \sigma^2)p^2 s + 2ps - 2\rho s$$

$$\dot{p} = -(1 + \sigma^2)ps + s - \rho p$$

$$\dot{s}_j = -\frac{2}{(1+\sigma^2)N_\Phi N_\Psi}\left(\frac{2}{N_\Psi^2}s_j^2 p_j^3 + N_\times s_j p_j^2 - s_j p_j\right) - 2\rho s_j.$$
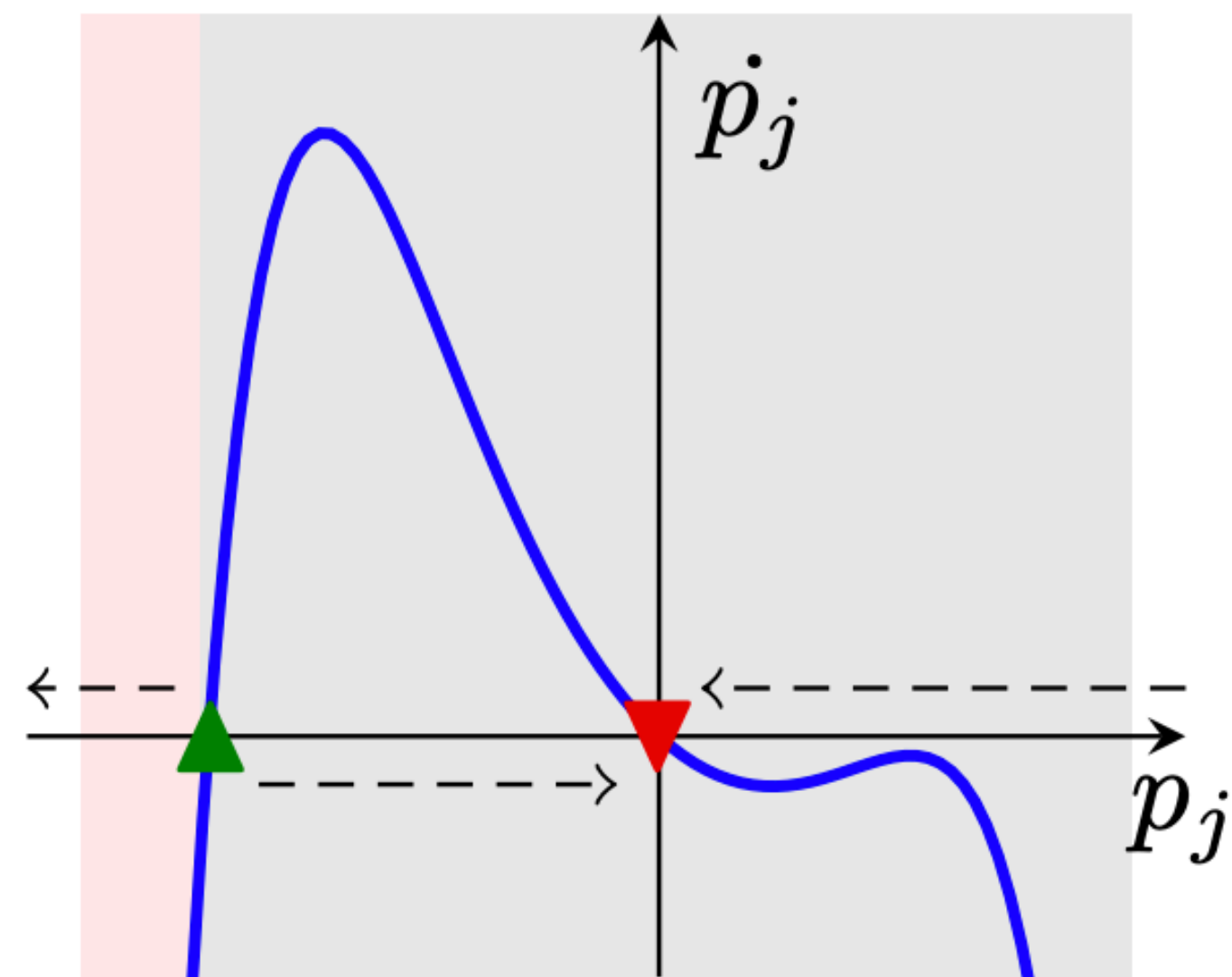
$$\dot{p}_j = -\frac{1}{(1+\sigma^2)N_\Phi N_\Psi}\left(\frac{2}{N_\Psi^2}s_j^2 p_j^2 + N_\times s_j p_j - s_j\right) - \rho p_j,$$

Adiabatic elimination



Adiabatic elimination



$$\dot{p} = p^2\{1 - (1 + \sigma^2)p\} - \rho p$$

$$\dot{p} = -\frac{2C_1 p^6 + C_2 p^3 - C_3 p^2}{1 + \sigma^2} - \rho p$$

# Bifurcation: collapsed solution is not stable 😁

Eigval ODE of projector  ( $C_i$ depends on $\|\mathbf{W}\mathbf{\Phi}\mathbf{x}\|$ )

$$\dot{p} = -\frac{2C_1p^6 + C_2p^3 - C_3p^2}{1 + \sigma^2} - \rho p$$

Eigval ODE of projector   ( $C_i$ depends on $\|\mathbf{W\Phi x}\|$ )

$$\dot{p} = -\frac{2C_1 p^6 + C_2 p^3 - C_3 p^2}{1 + \sigma^2} - \rho p$$

Eigval ODE of projector  ( $C_i$ depends on $\|\mathbf{W\Phi x}\|$ )

$$\dot{p} = -\frac{2C_1 p^6 + C_2 p^3 - C_3 p^2}{1 + \sigma^2} - \rho p$$

non-trivial solution exists



$\|\mathbf{W\Phi x}\|$
decay

$\|\mathbf{W\Phi x}\|$
decay

collapse $p = 0$ is **saddle**!

# Part 1: What we learn from nonlinear dynamics

- 😁 Dynamics analysis provides **stability analysis** beyond analyzing loss minimizer solely

  - ❖ Why StopGrad? Why encoder-predictor? etc.

- 😁 Difference loss function may yield more **adaptivity** during optimization

$$\dot{p} = p^2\{1 - (1 + \sigma^2)p\} - \rho p$$

$$\dot{p} = -\frac{2C_1 p^6 + C_2 p^3 - C_3 p^2}{1 + \sigma^2} - \rho p$$

L2 dynamics

cosine dynamics

- 🥲 What we don't answer: **feature learning**

  - ❖ since analytical solution to ODE typically requires strong Gaussianity assumptions

# What we can learn
# from nonlinear dynamics and **neuroscience**

Ishikawa, S.*, Yamada, M.*, **Bao, H.**, & Takezawa, Y. (ICLR2025)
PhiNets: Brain-inspired Non-contrastive Learning Based on Temporal Prediction Hypothesis.

# Predictive coding

- Brain predicts a future/neighboring input signal
  - ❖ dopamine is secreted if prediction makes a mistake

Contrastive predictive coding [van den Oord+ 18]



van den Oord et al. (2018) Representation Learning with Contrastive Predictive Coding.

# Predictive coding

- Brain predicts a future/neighboring input signal **at various level**

# Neocortex and hippocampus

Hippocampus

**Short-term** memory

external stimuli

Neocortex

**Long-term** memory

store memory

# Hippocampus as a self-supervised learning model

[Chen+ 24]

**prediction**

NeoCortex

Entorhinal Cortex

EC

I

II

III

IV

V

VI

NC

external stimuli

$+2\tau$

CA3

synaptic delay

$+\tau$

$+\tau$

CA1

store back

Transmission delay b/w CA3 and CA1 forms a self-supervised feedback
⇒ with prediction, neural activity is replicated more accurately

Chen et al. (2024) Predictive Sequence Learning in the Hippocampal Formation.

# Temporal prediction hypothesis

Temporal difference

additional predictor (CA1)

additional loss
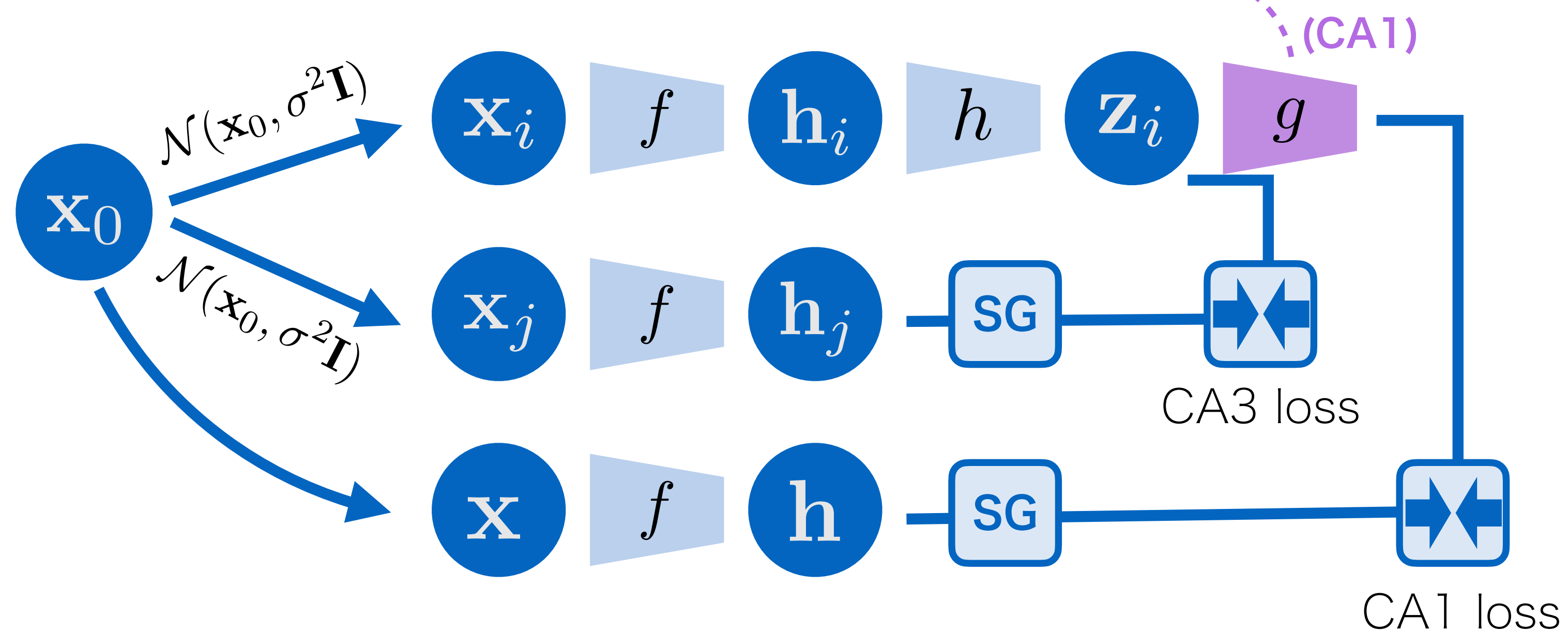
# Φ-Net

CA3 predictor · CA1 predictor

augmented signals

original signal

- Encoder $f$ is shared
- All layers are optimized by backprop simultaneously

# 🤔 But why [additional predictor]?

Analysis model



$$\mathcal{L}(\mathbf{W}_f, \mathbf{W}_g, \mathbf{W}_h) = \frac{1}{2}\mathbb{E}\left[\underbrace{\|\mathbf{W}_h\mathbf{W}_f\mathbf{x}_1 - \mathrm{SG}(\mathbf{W}_f\mathbf{x}_2)\|^2}_{\text{CA3 loss}} + \underbrace{\|\mathbf{W}_g\mathbf{W}_h\mathbf{W}_f\mathbf{x}_1 - \mathrm{SG}(\mathbf{W}_f\mathbf{x})\|^2}_{\text{CA1 loss}}\right]$$

Disclaimer: cosine loss is not considered for simplicity

# 🤔 But why <u>additional predictor</u>?

$$\mathcal{L}(\mathbf{W}_f, \mathbf{W}_g, \mathbf{W}_h) = \frac{1}{2}\mathbb{E}\left[\|\mathbf{W}_h\mathbf{W}_f\mathbf{x}_1 - \mathrm{SG}(\mathbf{W}_f\mathbf{x}_2)\|^2 + \|\mathbf{W}_g\mathbf{W}_h\mathbf{W}_f\mathbf{x}_1 - \mathrm{SG}(\mathbf{W}_f\mathbf{x})\|^2\right]$$

eigendecomposition
adiabatic elimination

$$\begin{cases} \dot{p}_h &= \{(1 + p_g) - (1 + \sigma^2)(1 + p_g^2)p_h\}p_h^2 - \rho p_h \qquad \text{(CA3 predictor)} \\ \dot{p}_g &= \{1 - (1 + \sigma^2)p_h\}p_h^3 - \rho p_g \qquad\qquad\qquad\quad \text{(CA1 predictor)} \end{cases}$$

cf. SimSiam dynamics

$$\dot{p} = p^2\{1 - (1 + \sigma^2)p\} - \rho p \qquad\qquad\qquad\qquad\qquad \text{(CA3 predictor)}$$

# PhiNet dynamics (2D)



MEDIUM ($\rho = 0.03$)

$\dot{p}_h = 0$

$\dot{p}_g = 0$

stable equilib

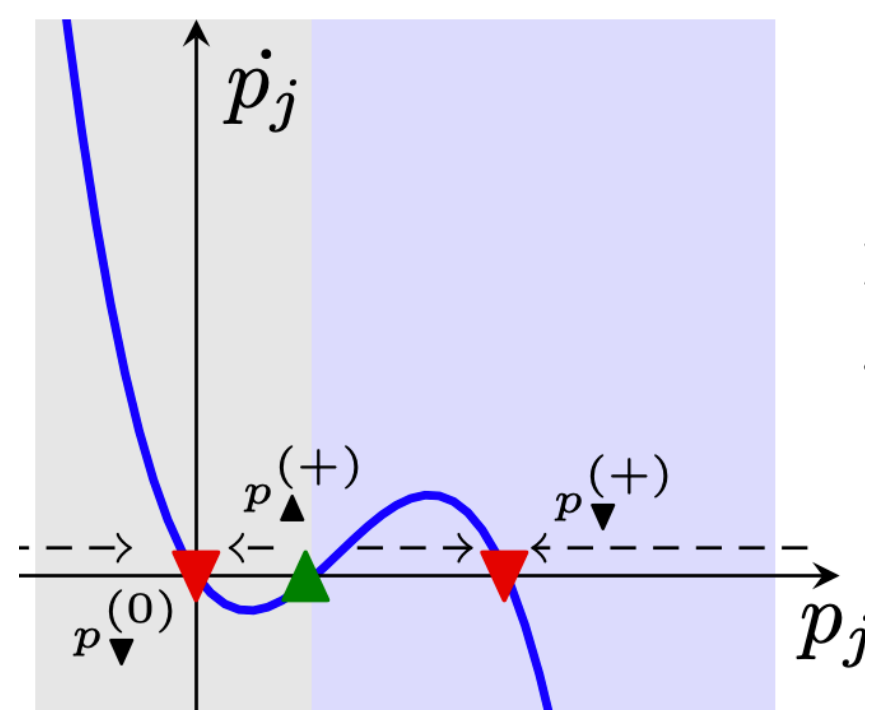# PhiNet (2D) vs SimSiam (1D)

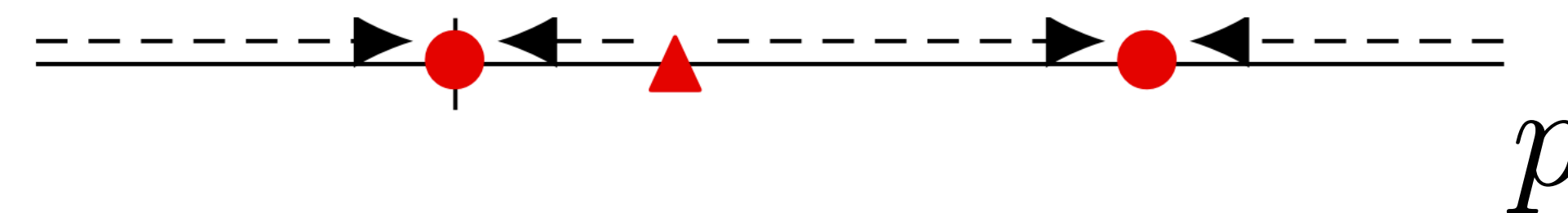PhiNet dynamics



extract nullclines

$p_g$

$p_h$

topologically
conjugate

SimSiam dynamics

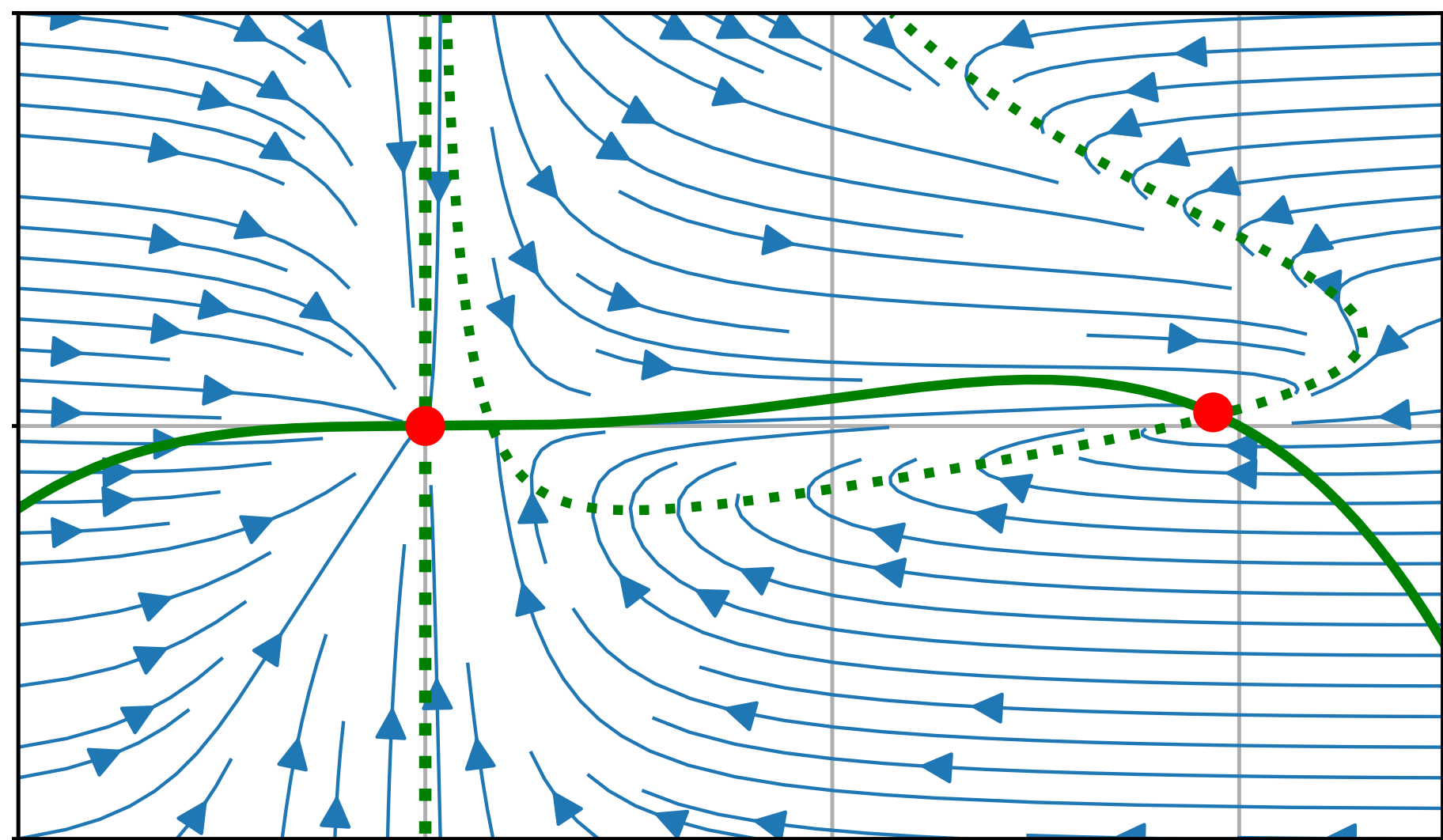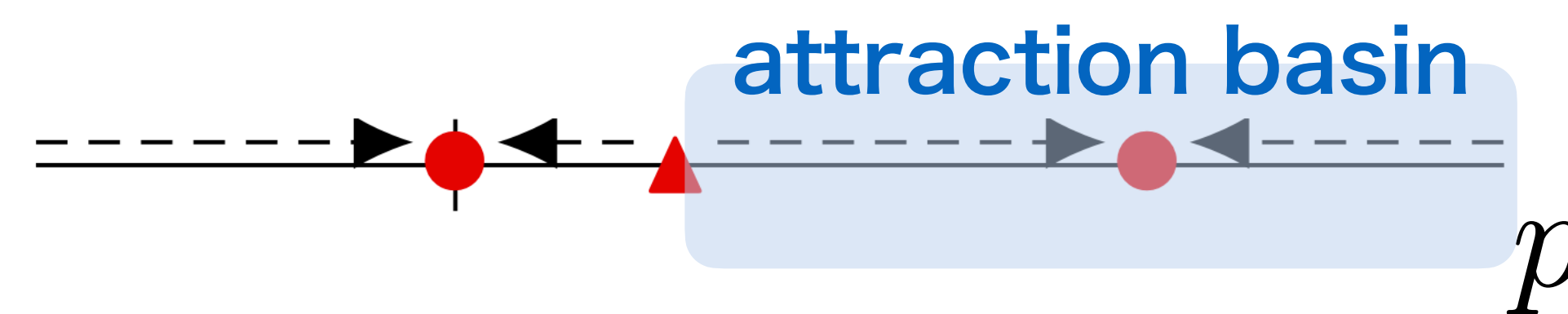extract nullclines

$p$

# PhiNet (2D) vs SimSiam (1D)

PhiNet dynamics



extract nullclines →

$p_g$

**attraction basin (wider 😁)**

$p_h$

↕ topologically conjugate

SimSiam dynamics



extract nullclines →

**attraction basin**

$p$

# PhiNet dynamics with different weight decay



STRONG ($\rho = 0.12$)

# PhiNet dynamics with different weight decay



MEDIUM $(\rho = 0.03)$

$p_g$ (CA1)

$p_h$ (CA3)

# PhiNet dynamics with different weight decay



LIGHT $(\rho = 0.003)$

# PhiNet dynamics with different weight decay



WEAK ($\rho = 0.0001$)

WEAK ($\rho = 0.0001$)

$p_g$ (CA1)

$p_h$ (CA3)

SimSiam cannot avoid collapse
if negatively initialized

# Enhanced stability wrt weight decay



(a1) Batch size=128 (a2) Batch size=1024

Legend: PhiNet(g=I, mse), PhiNet(mse), X-PhiNet(mse), X-PhiNet(cos), PhiNet(g=h, mse), SimSiam

More details from Makoto!

Dataset: CIFAR-10 / Evaluation: kNN accuracy

# Summary