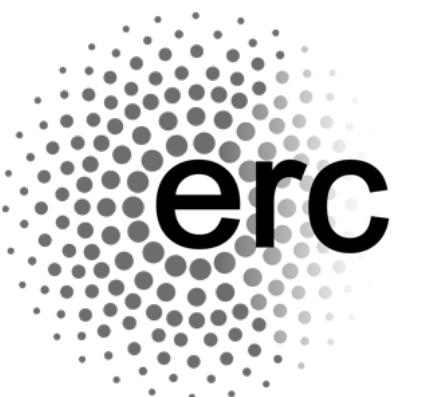


# MECHANICAL MODELS OF FEATURE LEARNING

By Ivan Dokmanić

August 8, 2025, Cargèse

---



PRACTICAL DEEP LEARNING IS  
DEEP  
NON-LINEAR  
FEATURE LEARNING

PRACTICAL DEEP LEARNING IS  
DEEP  
NON-LINEAR  
FEATURE LEARNING

A GREAT DEAL OF THEORY ADDRESSES  
SHALLOW LEARNING  
LINEAR ARCHITECTURES  
NON-FEATURE LEARNING

# AN ACCIDENT: THE MICROSCOPIC PICTURE IS 100% KNOWN

# AN ACCIDENT: THE MICROSCOPIC PICTURE IS 100% KNOWN

$$\frac{pV}{T} = \text{const.}$$

BOYLE 1662  
CHARLES 1787,  
GAY-LUSSAC 1802,  
...

# AN ACCIDENT: THE MICROSCOPIC PICTURE IS 100% KNOWN

$$\frac{pV}{T} = \text{const.}$$

BOYLE 1662  
CHARLES 1787,  
GAY-LUSSAC 1802,  
...

$$p = \frac{1}{3} \frac{N}{V} \overline{mv^2} \quad f(v) \propto v^2 e^{-\frac{mv^2}{2k_B T}}$$

MAXWELL & BOLTZMANN 1860-1872

# THE BLACK BOX SYNDROME

# THE BLACK BOX SYNDROME



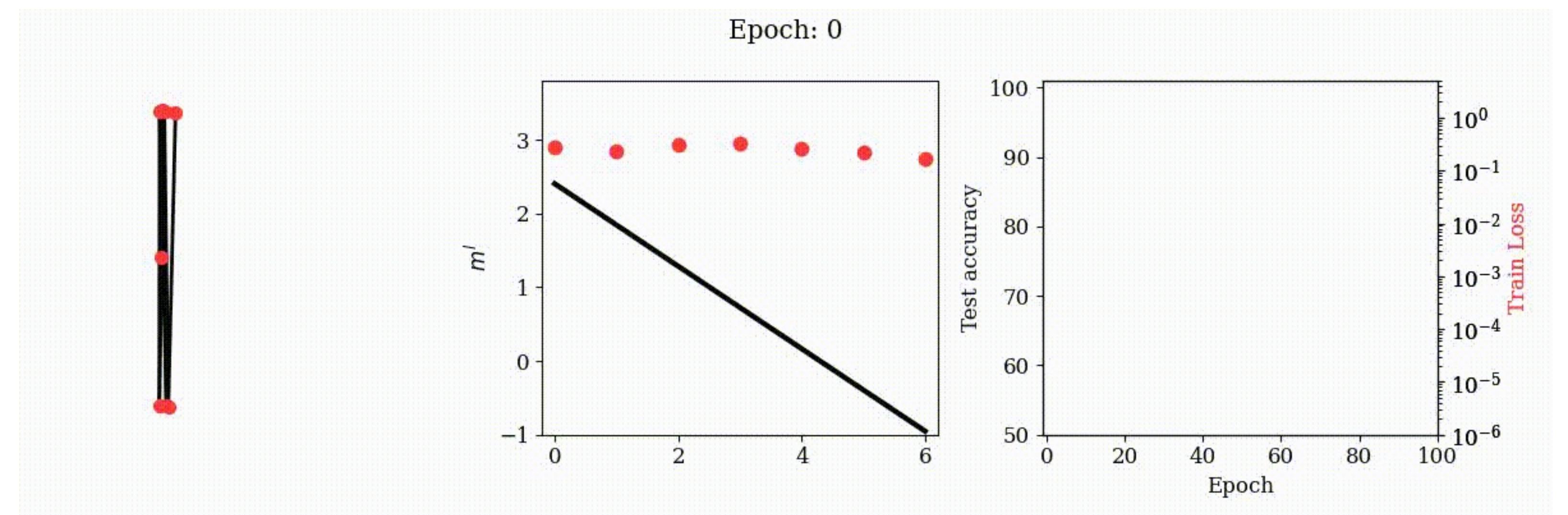
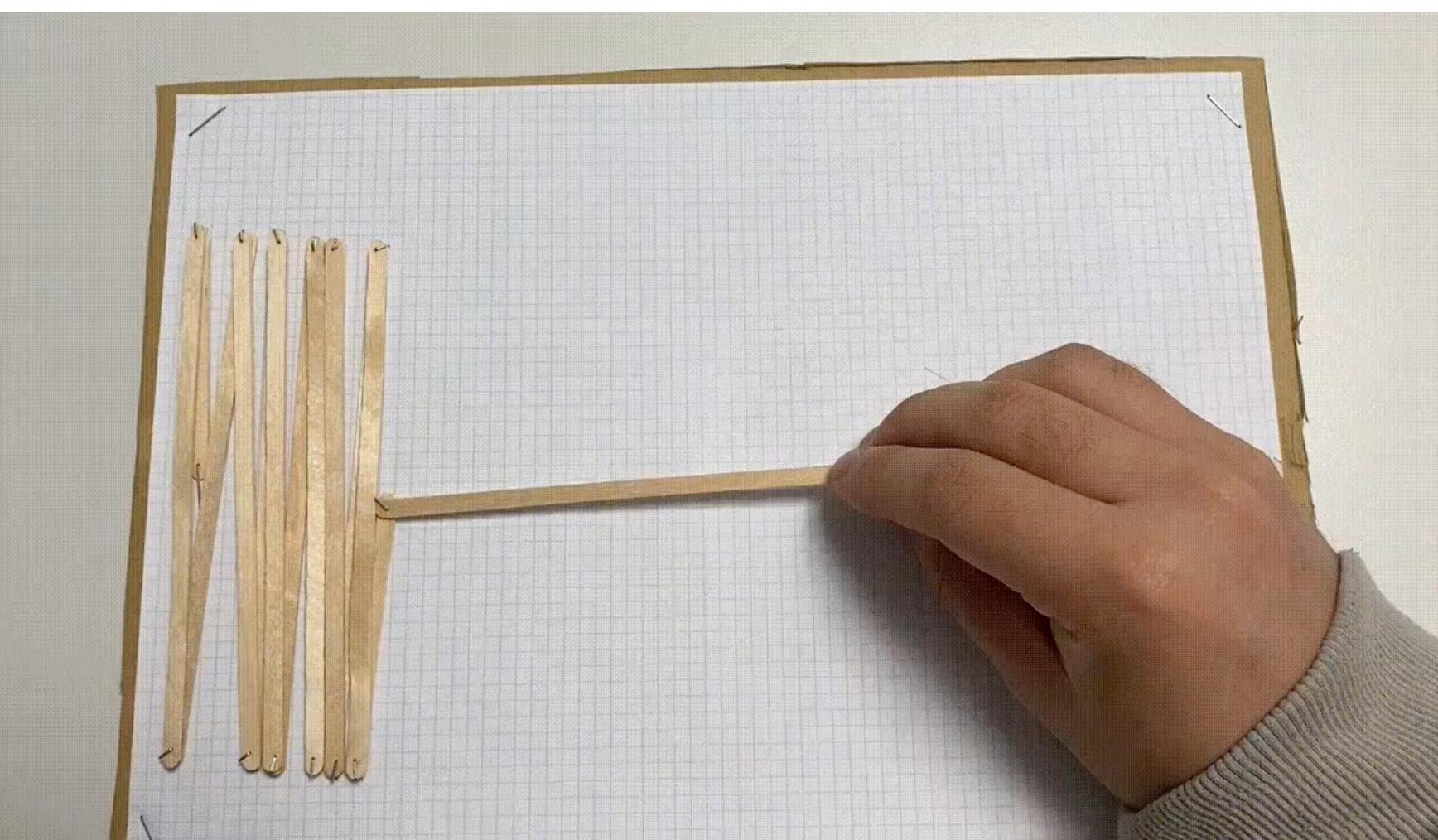
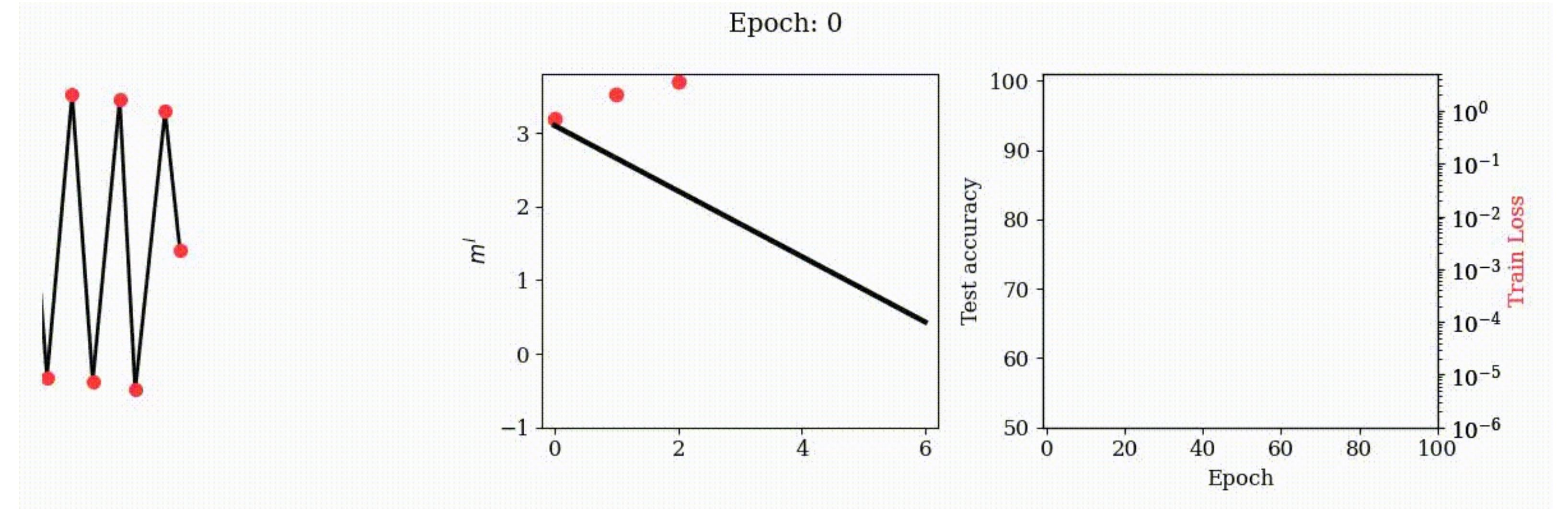
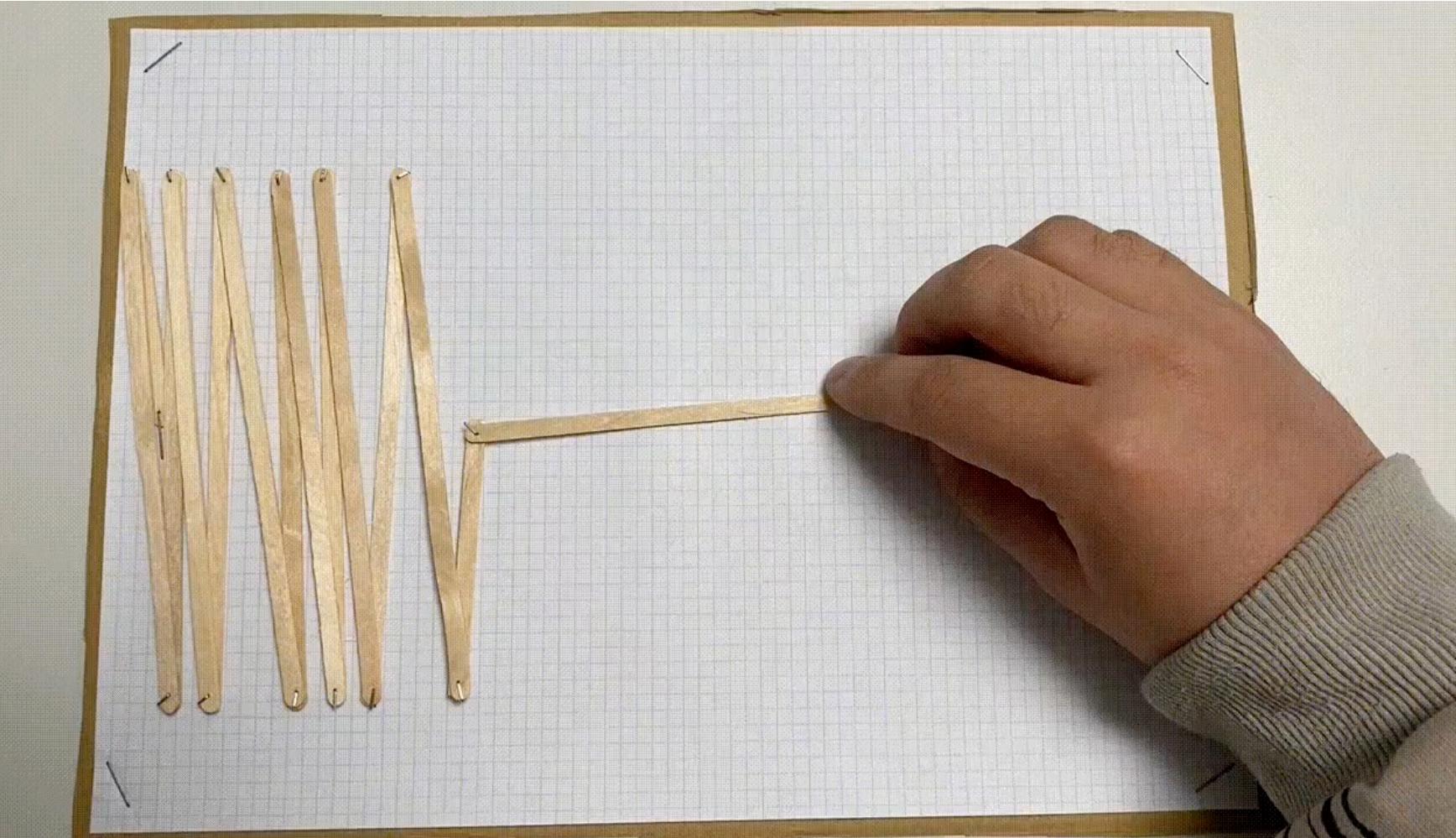
# THE BLACK BOX SYNDROME

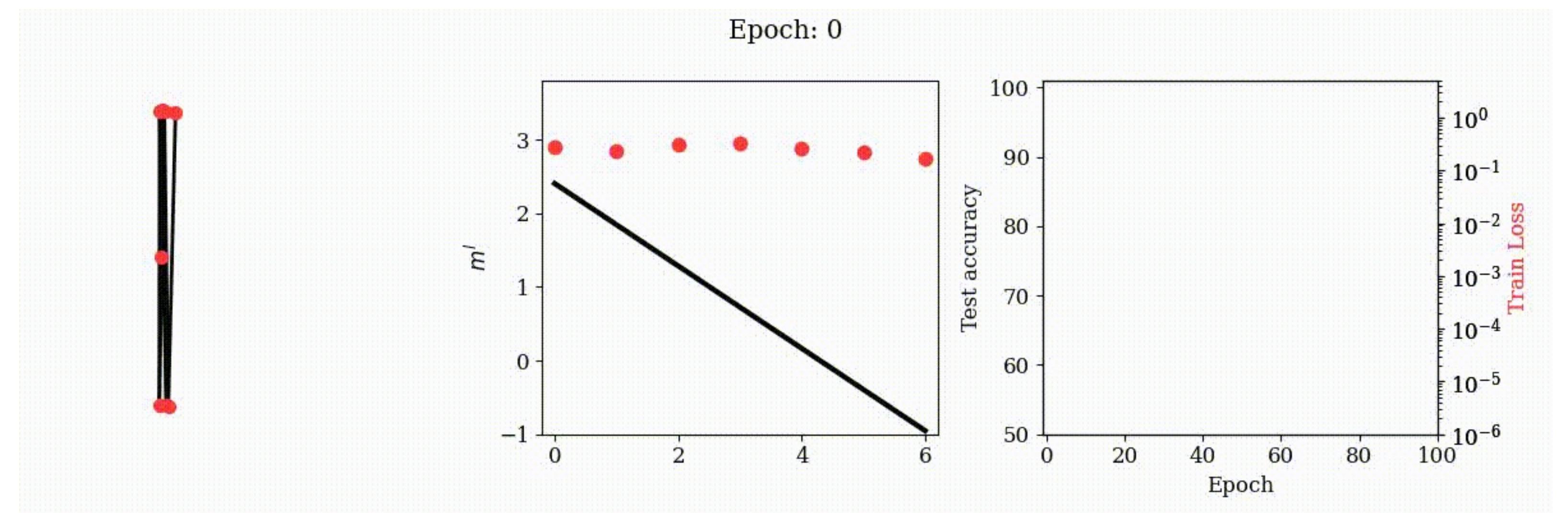
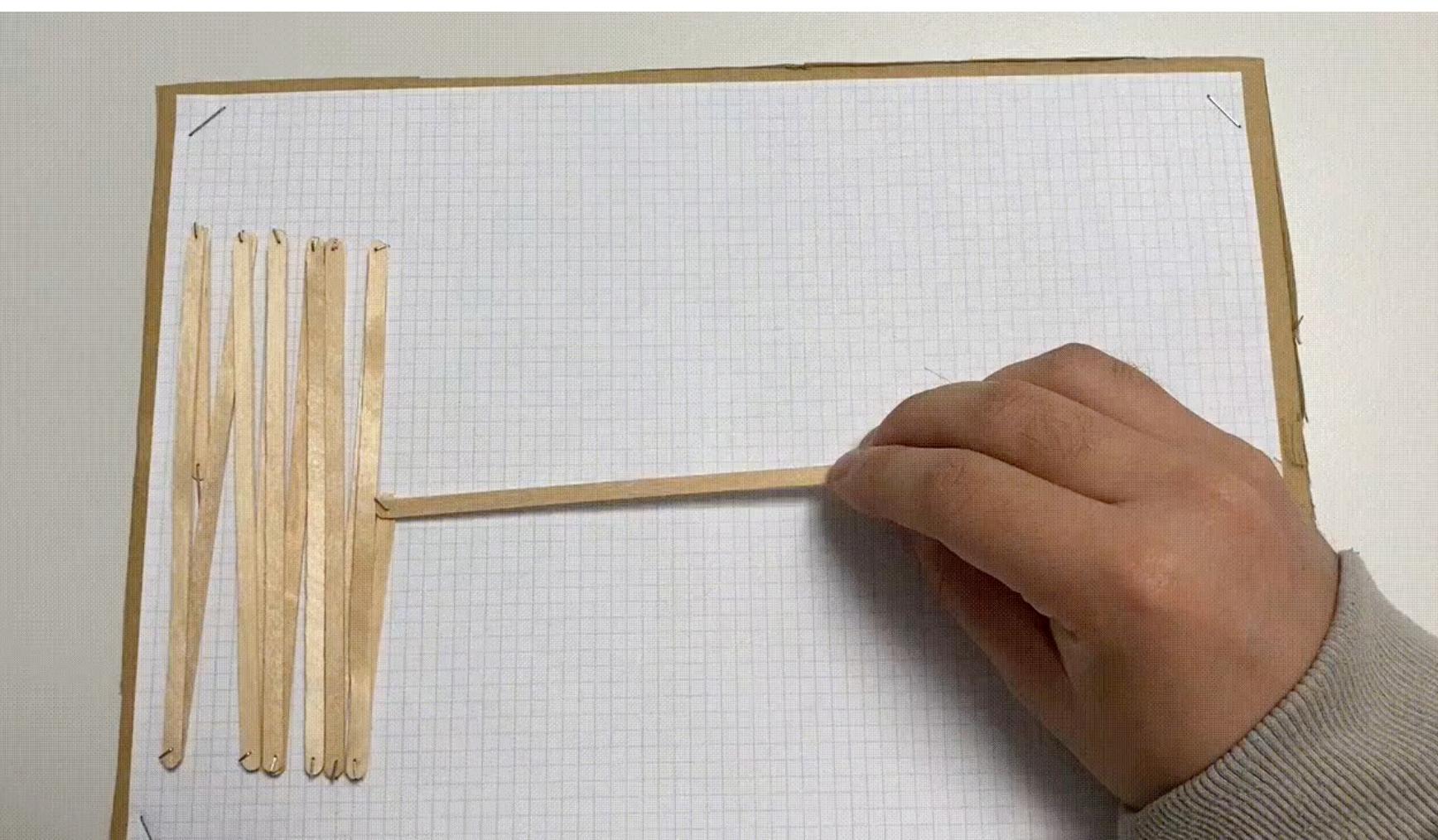
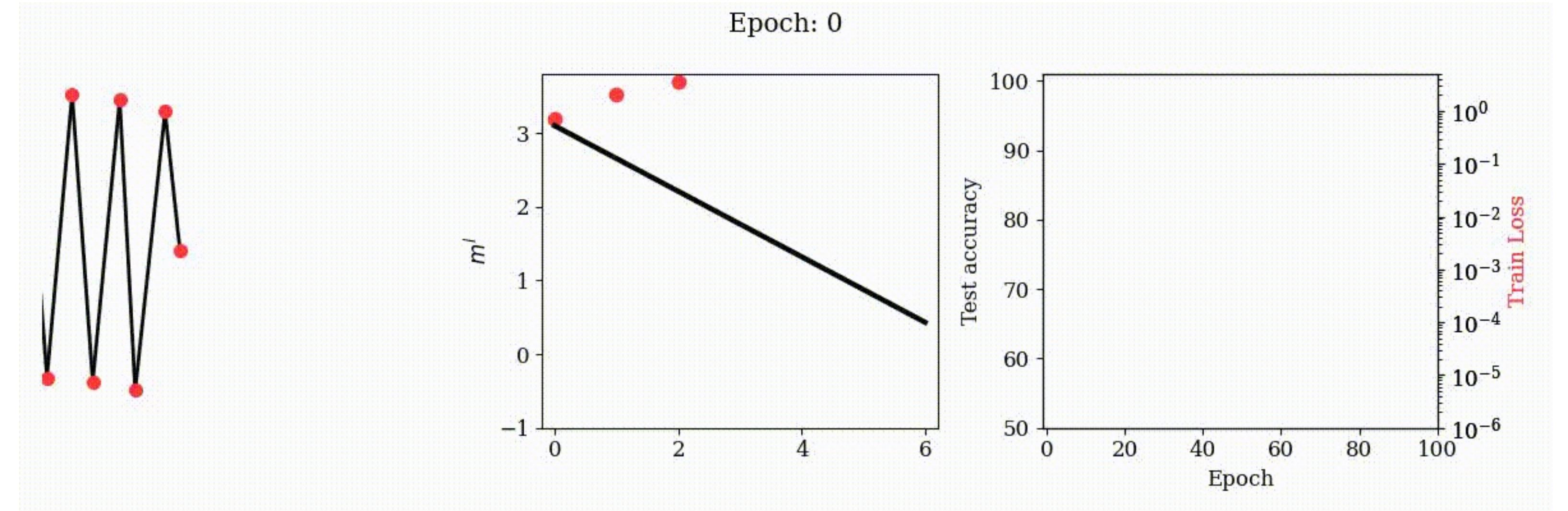
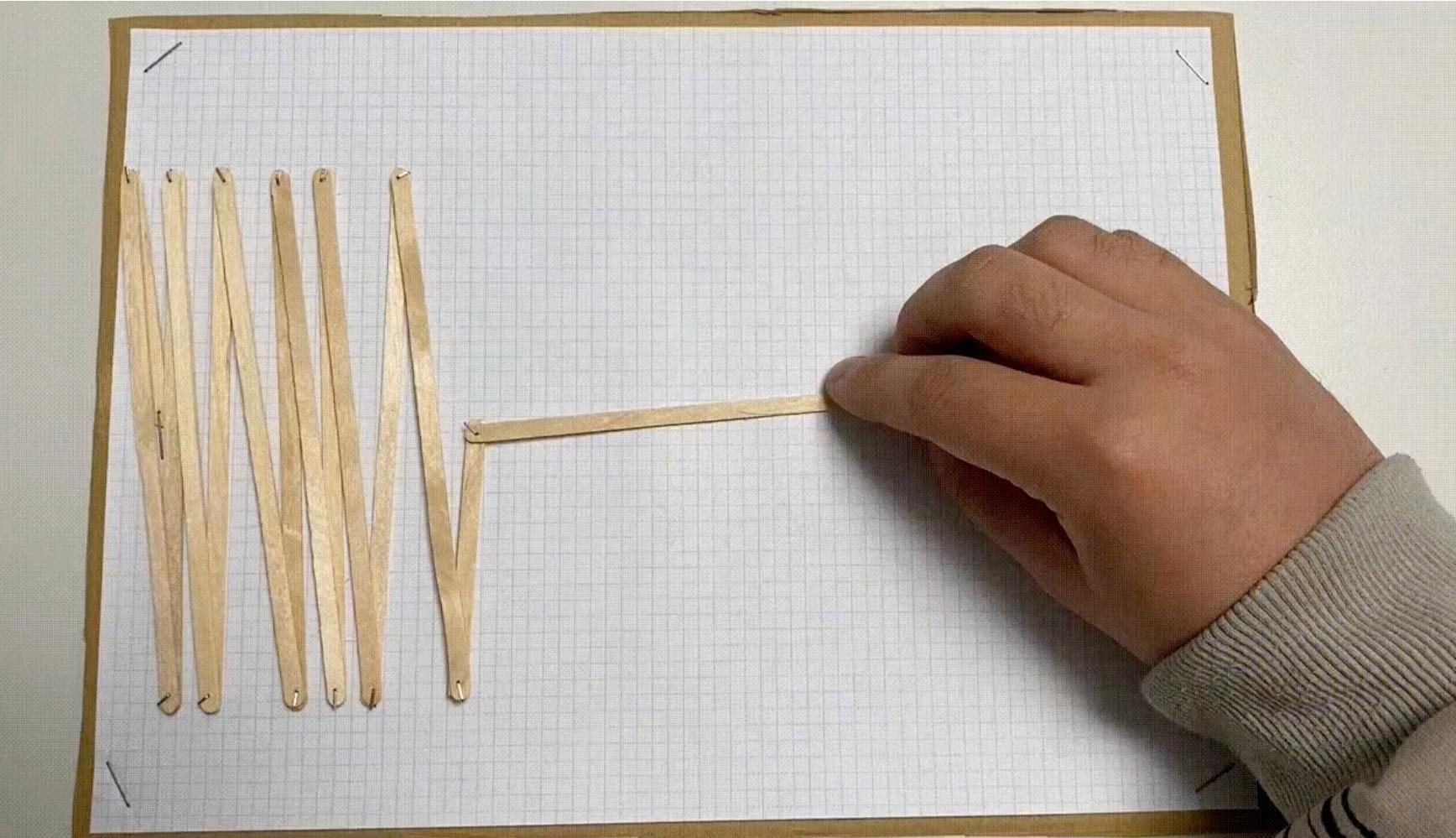


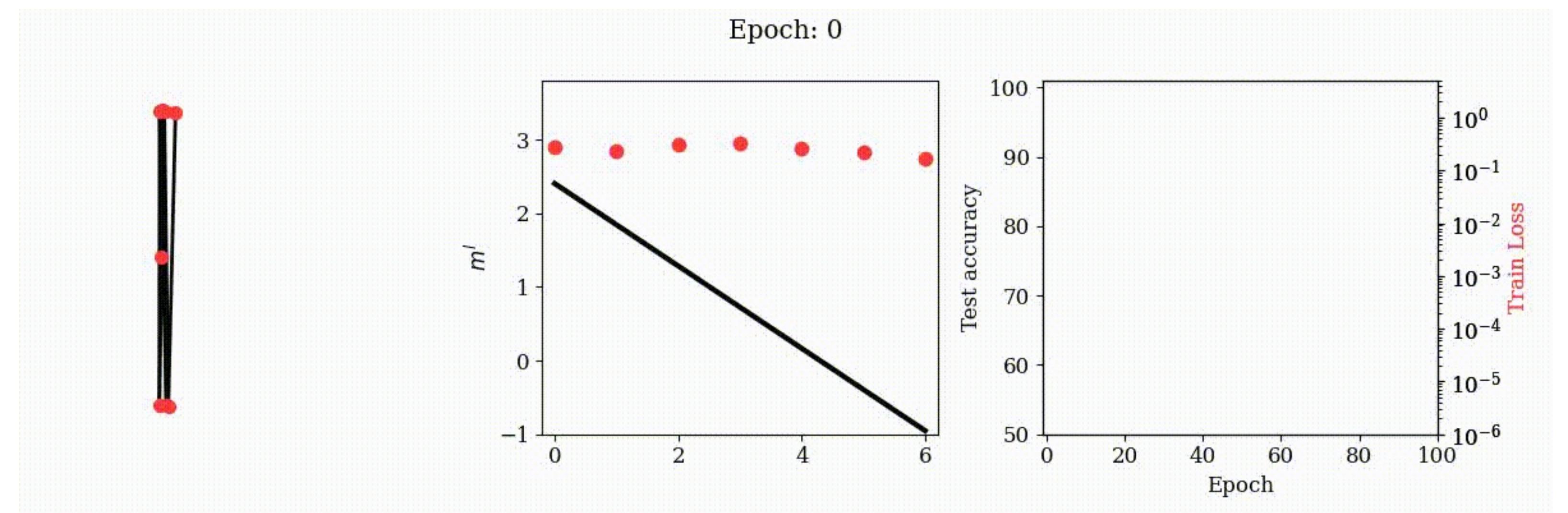
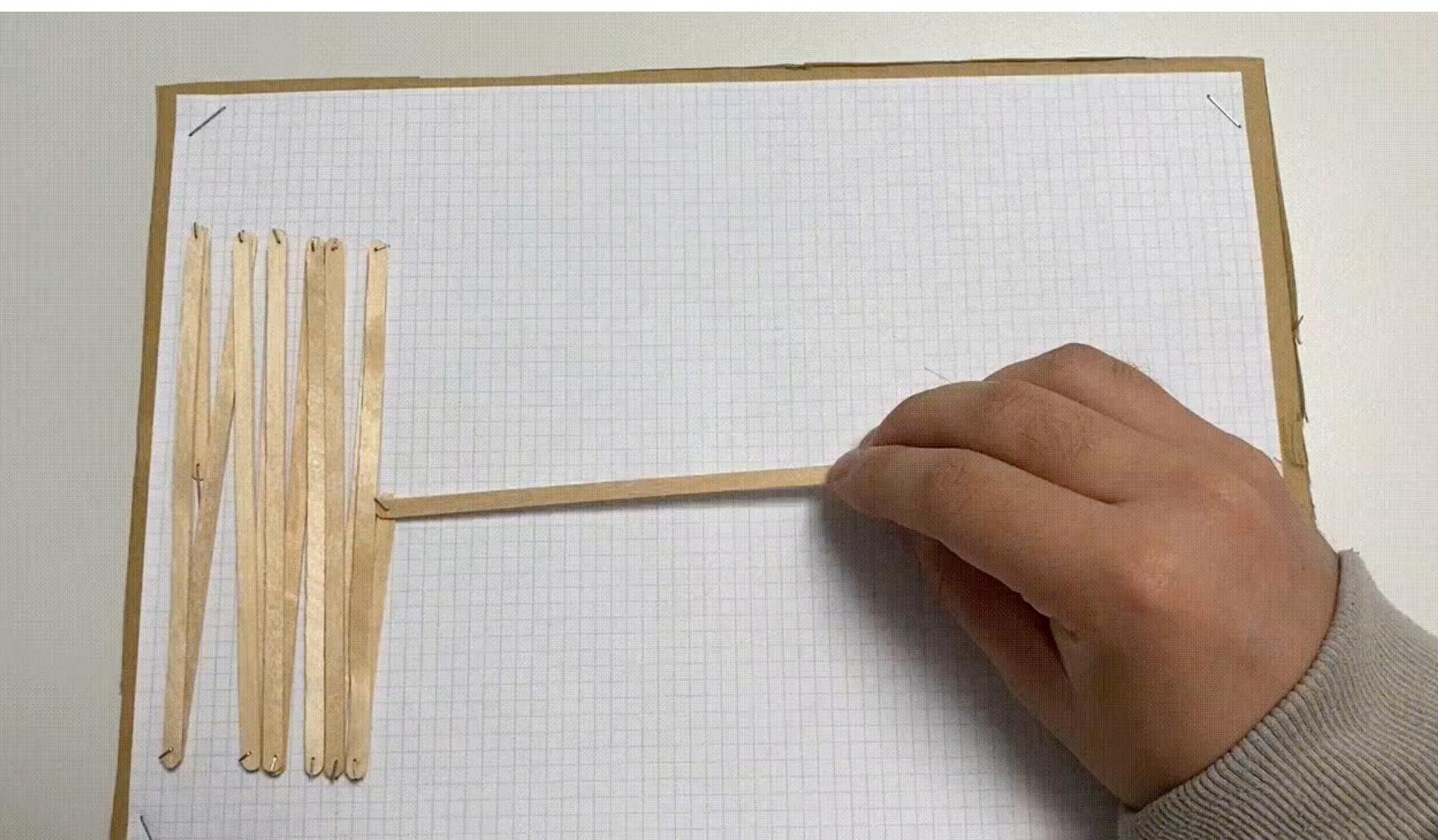
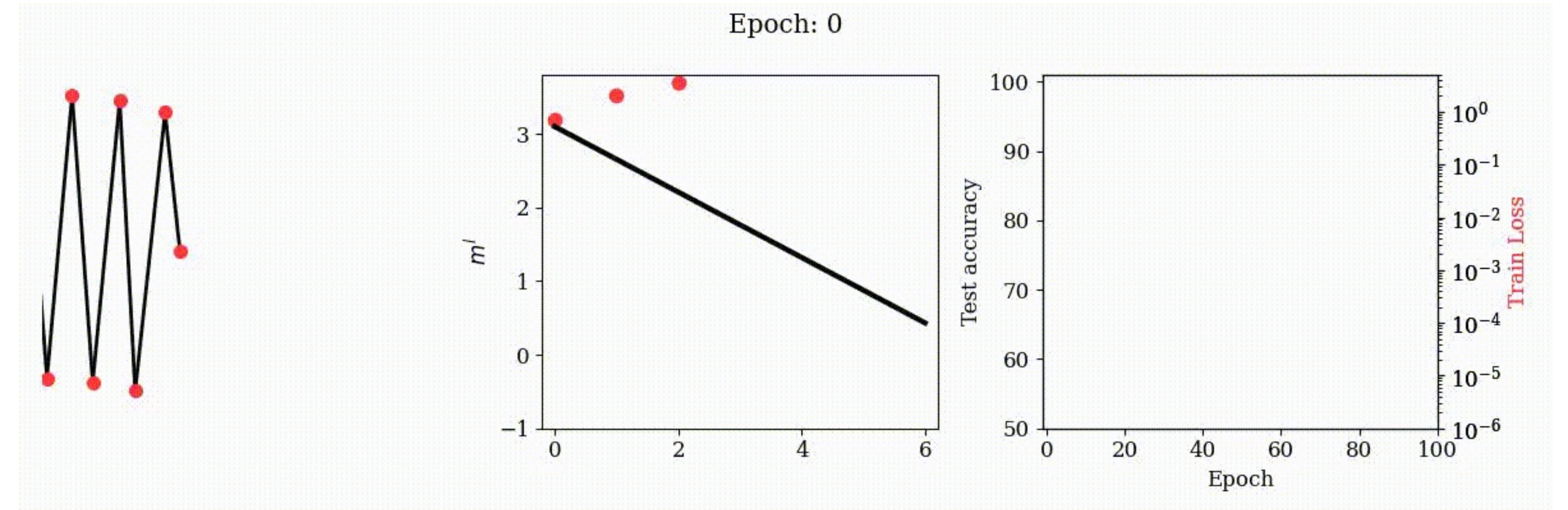
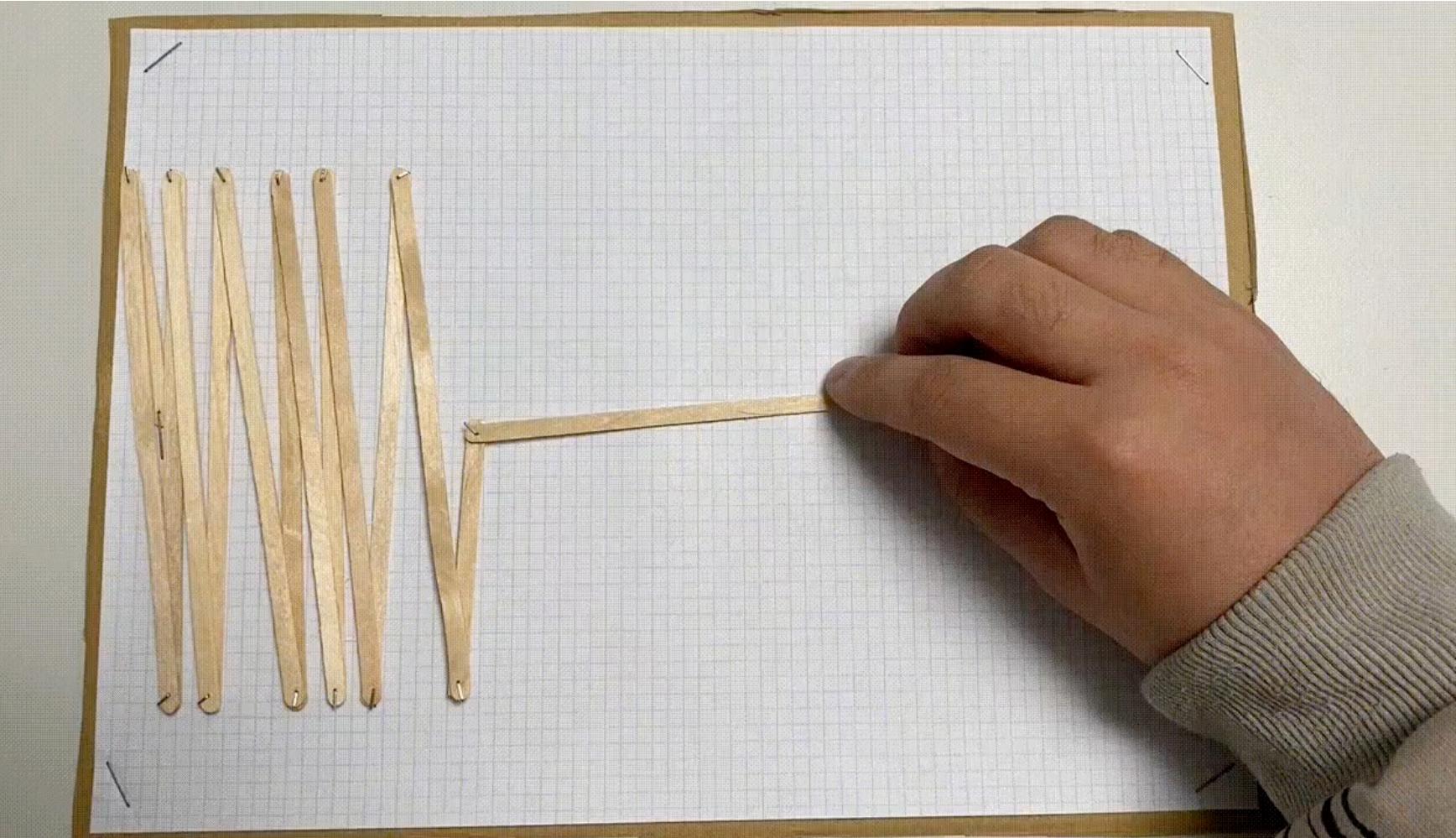
I AM NOT A BLACK BOX.

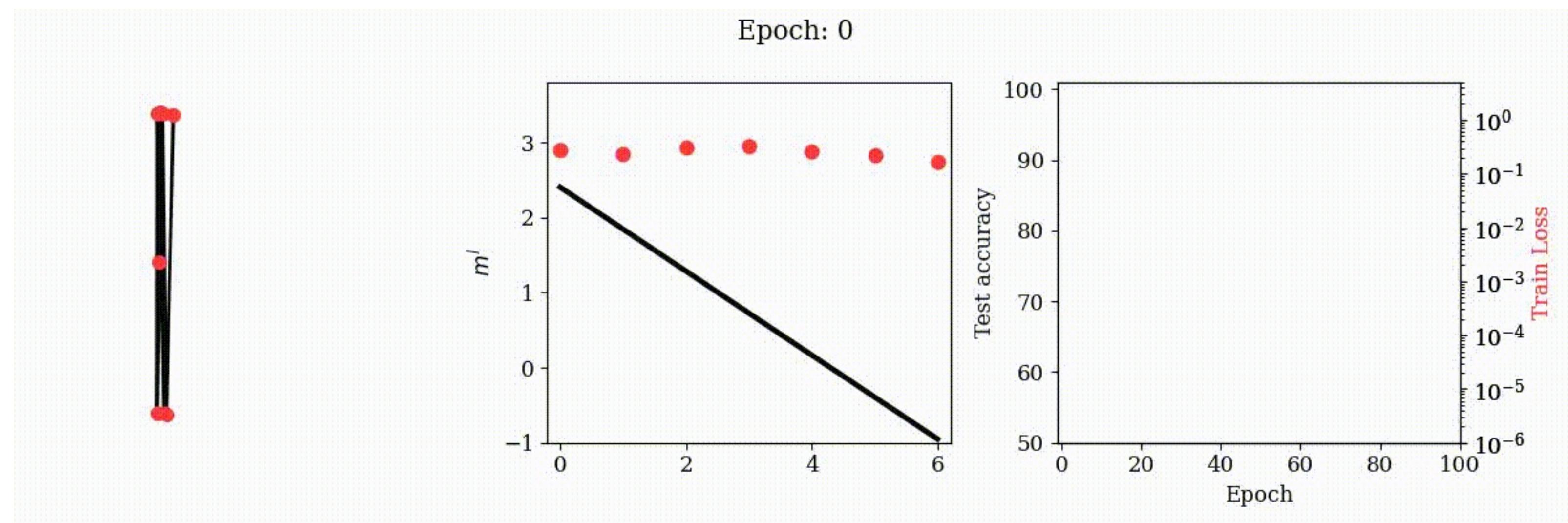
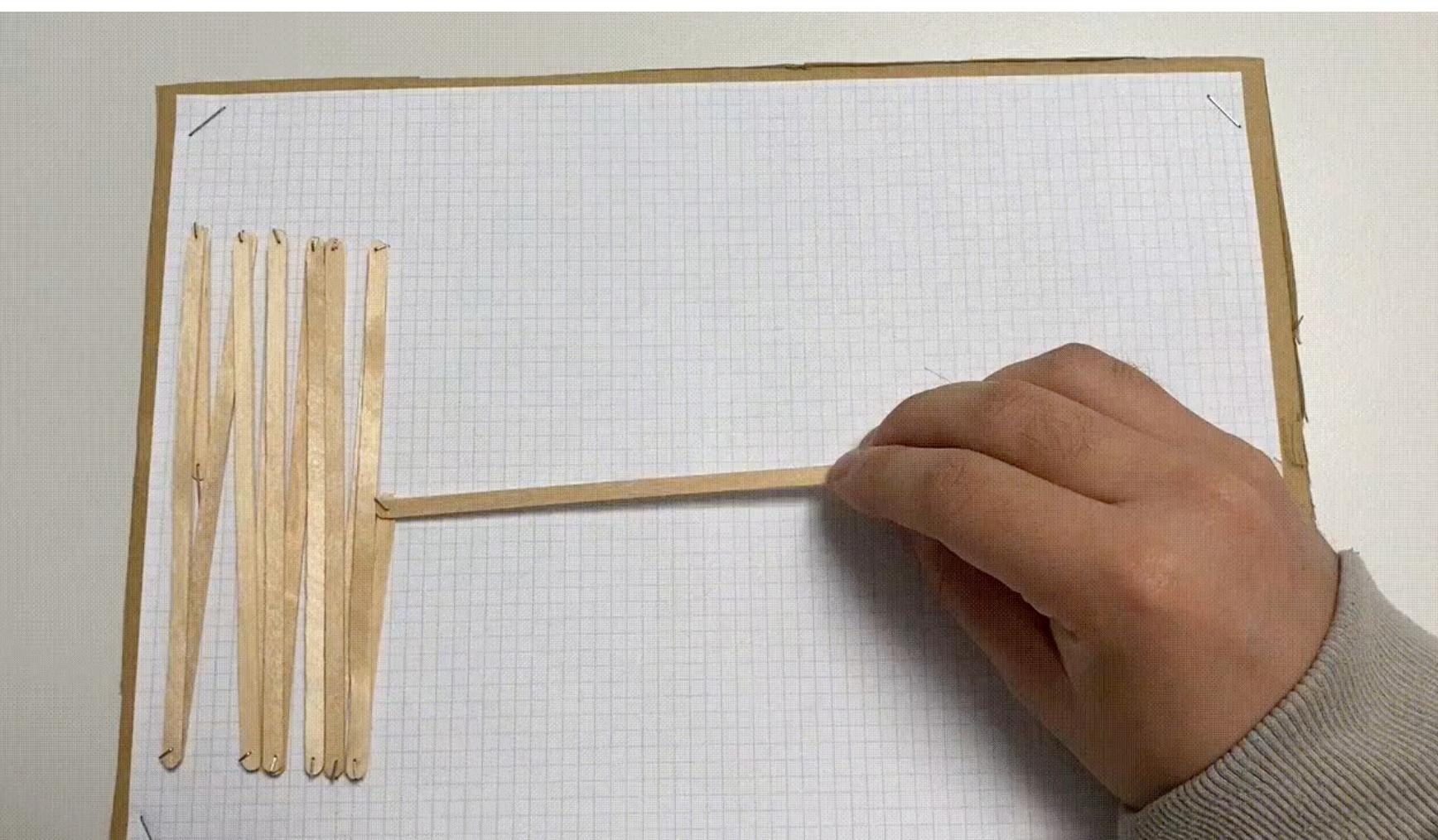
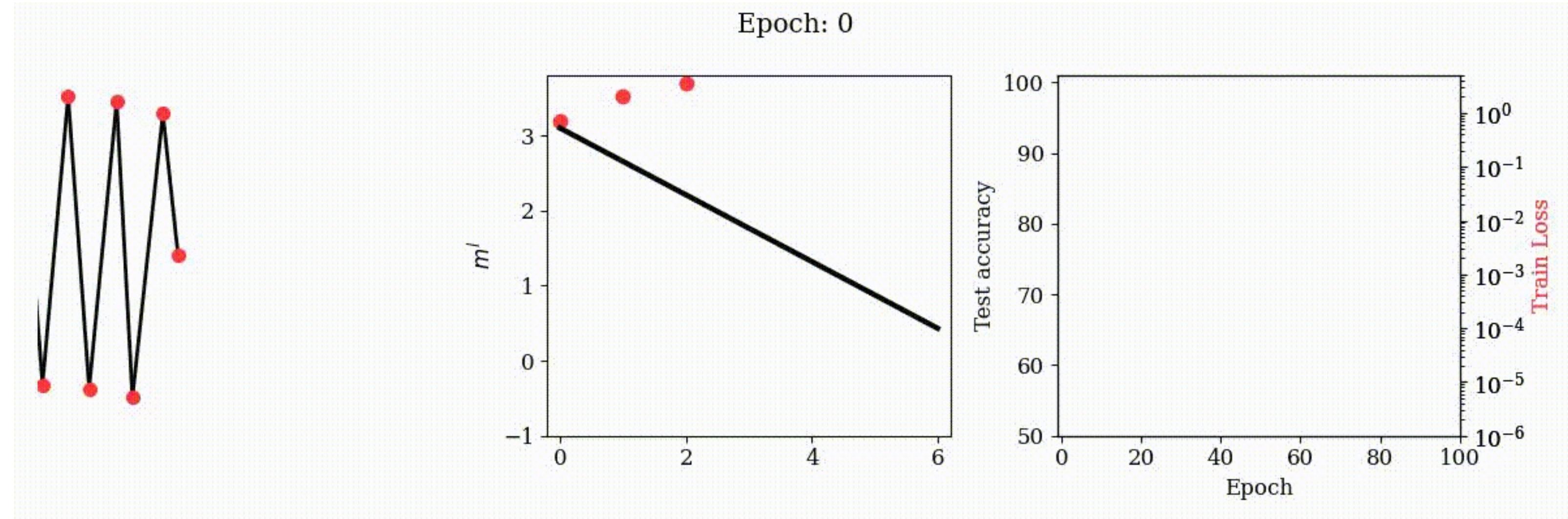
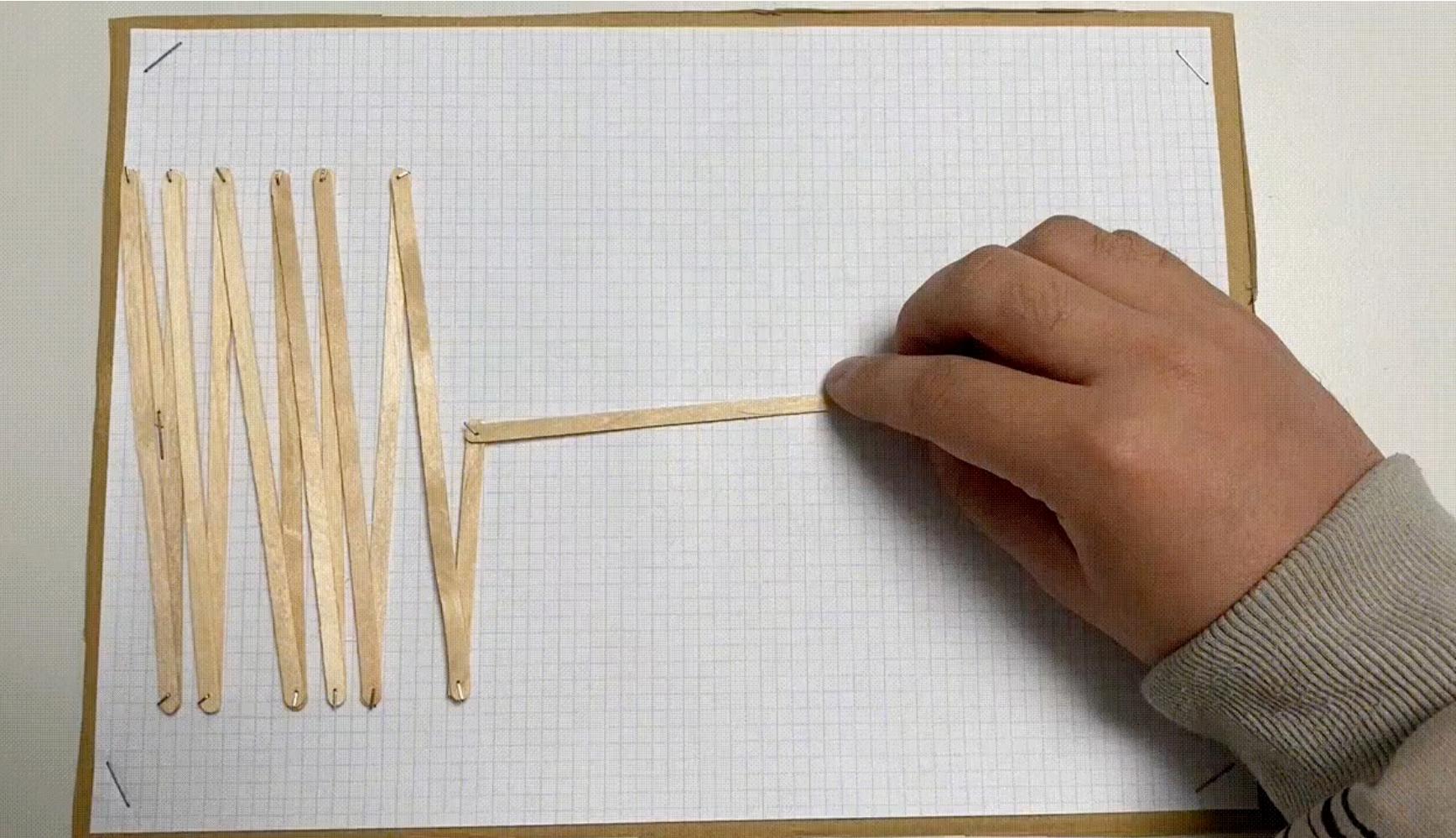
		<b>SCIENTIFIC THEORIES</b>
<b>Appetizer</b>	3	First principles of phenomenological?
		<b>PHENOMENOLOGICAL MODELS</b>
<b>First Course</b>	13	Springs, blocks, and folding rulers
		<b>FIRST-PRINCIPLE ATTEMPTS</b>
<b>Second Course</b>	26	Deep mean field theory and all that
		<b>ENDNOTES</b>
<b>Dessert</b>	34	Where I quote Boltzmann

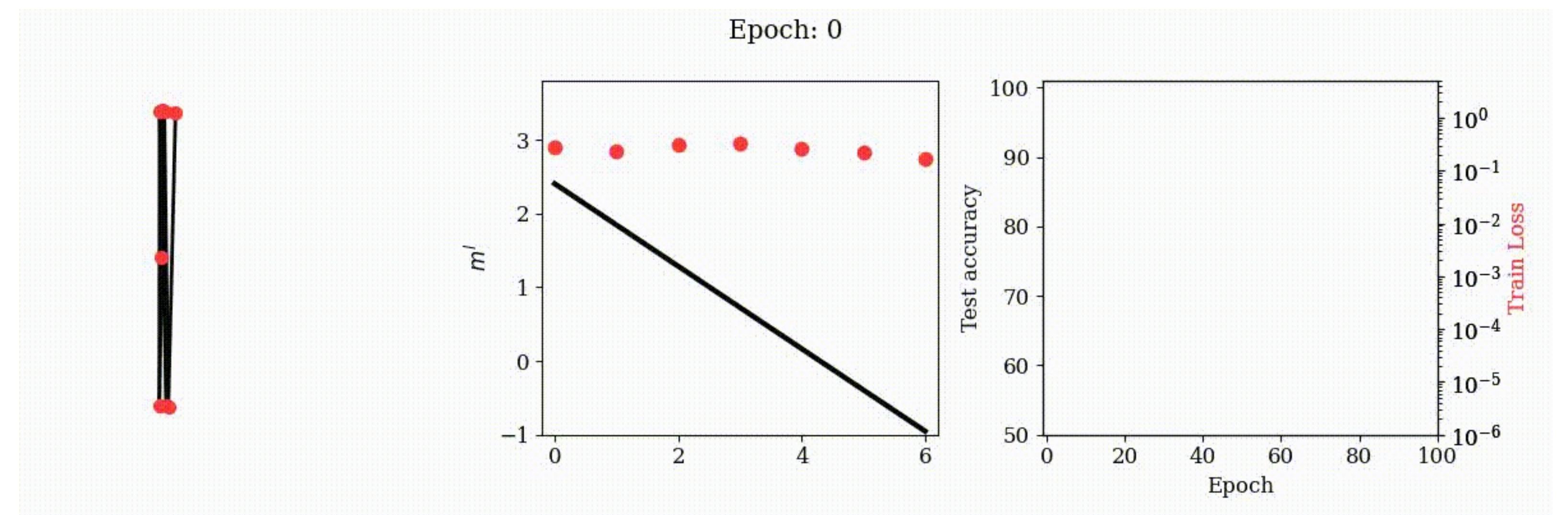
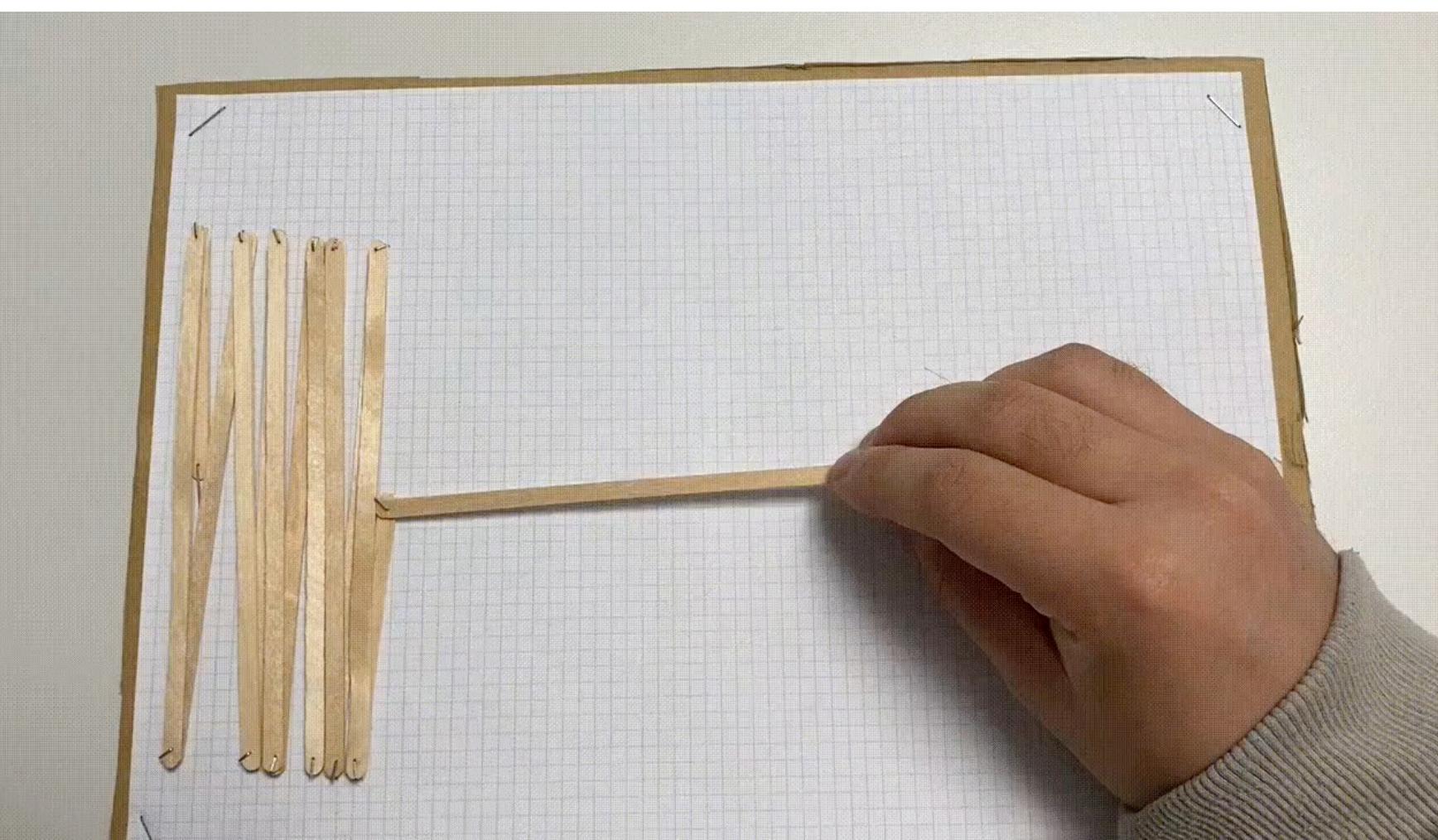
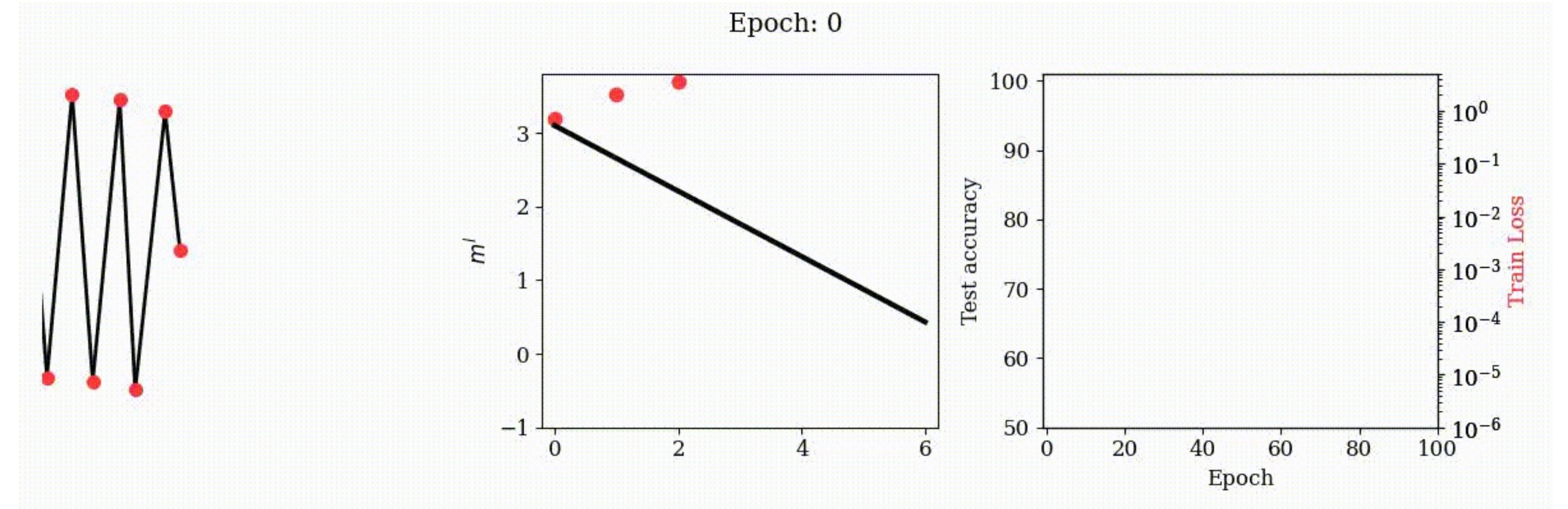
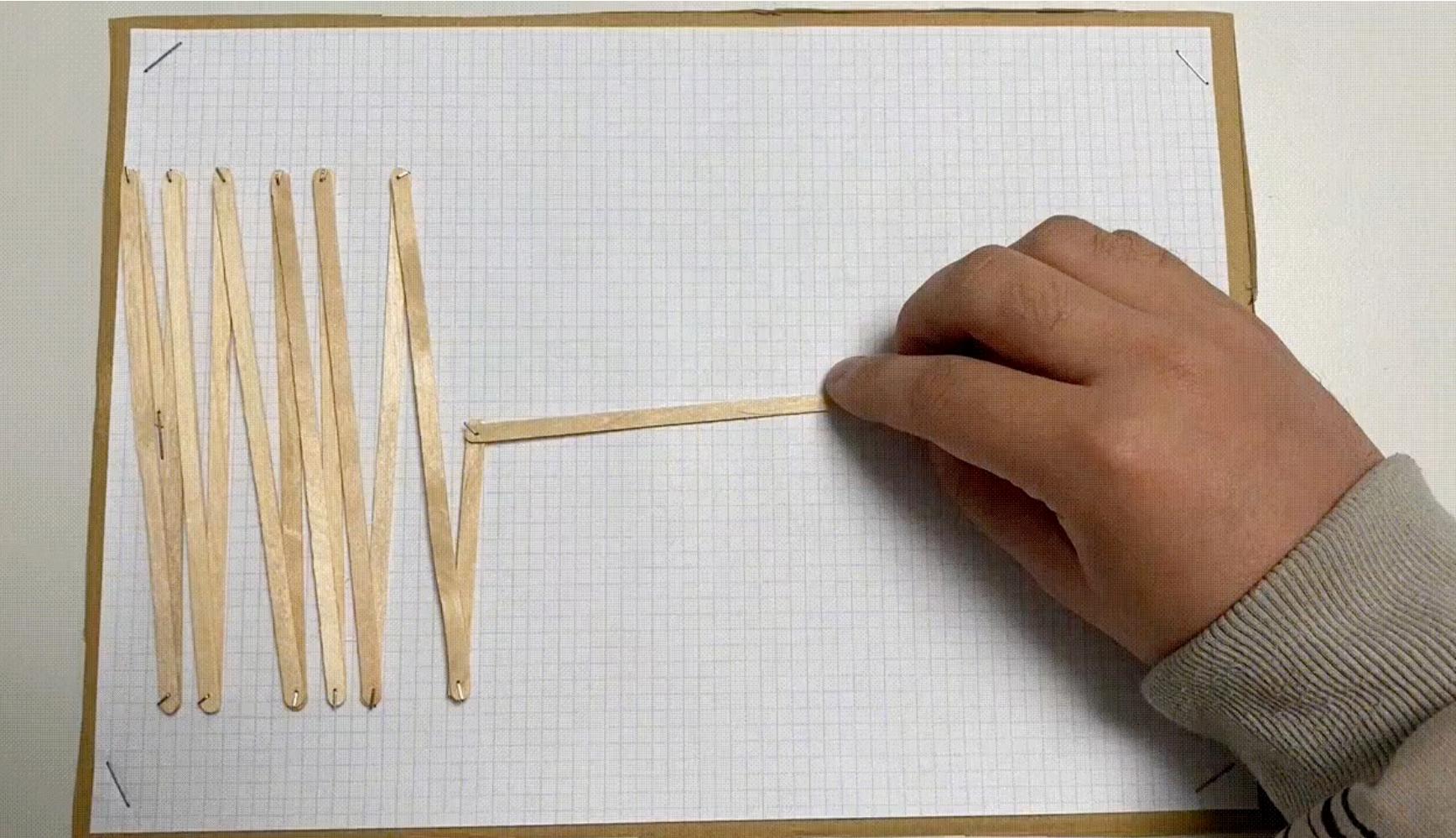
TRUE & USEFUL MACROSCOPIC  
STATEMENTS ABOUT DEEP NETS?



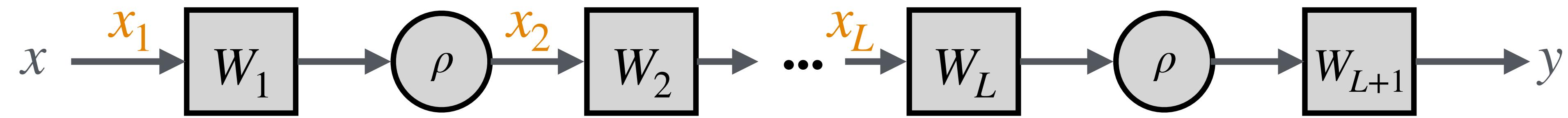




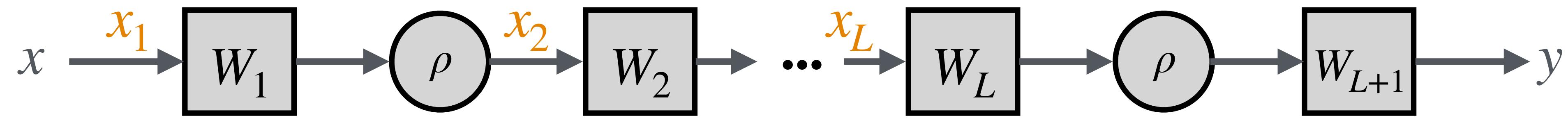




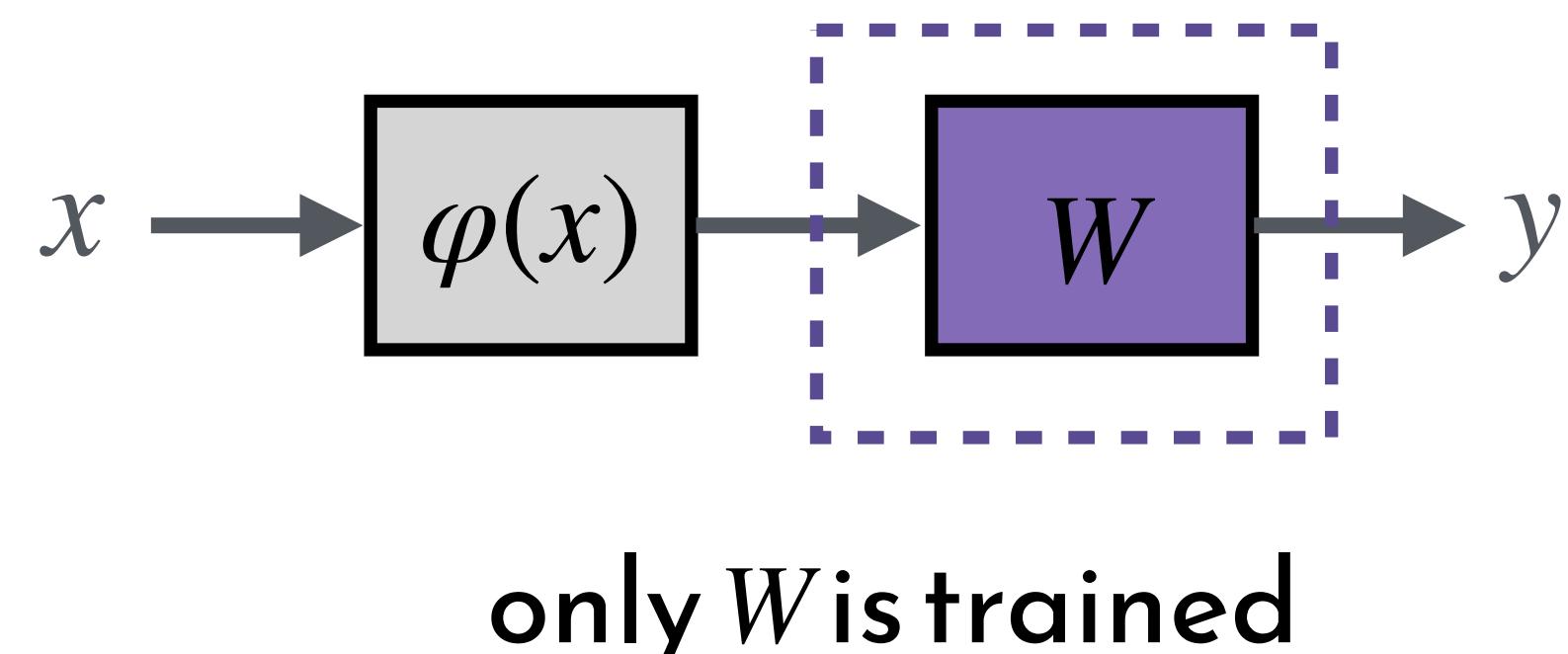
# FEATURE LEARNING



# FEATURE LEARNING

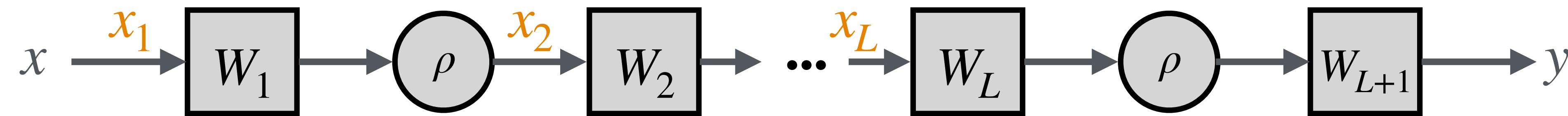


Non-feature learning (old)

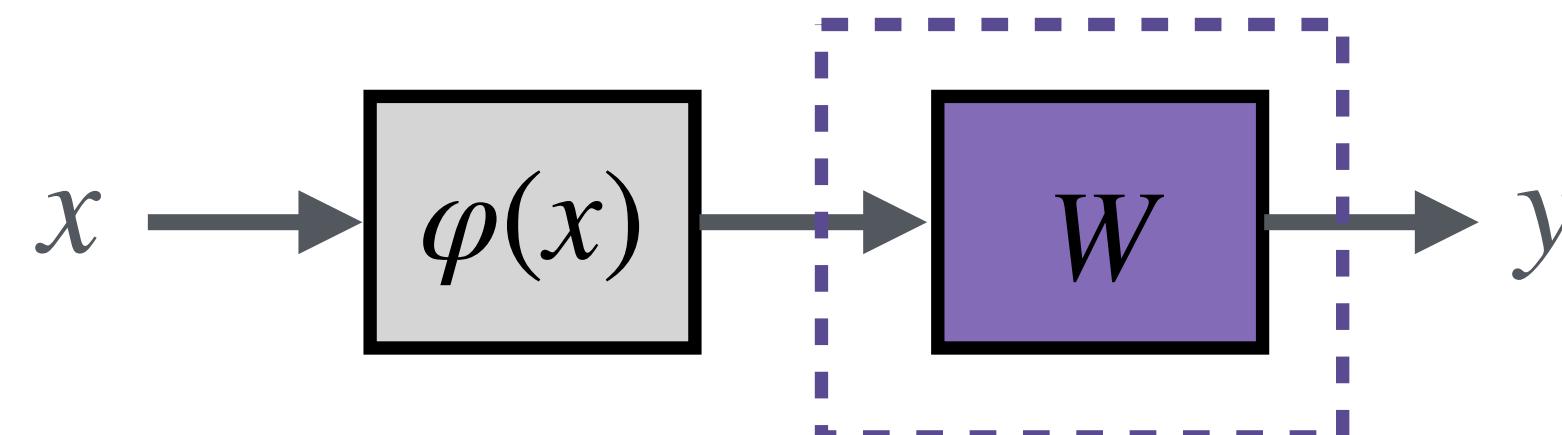


$$f_W(x) = W \circ \varphi(x)$$

# FEATURE LEARNING



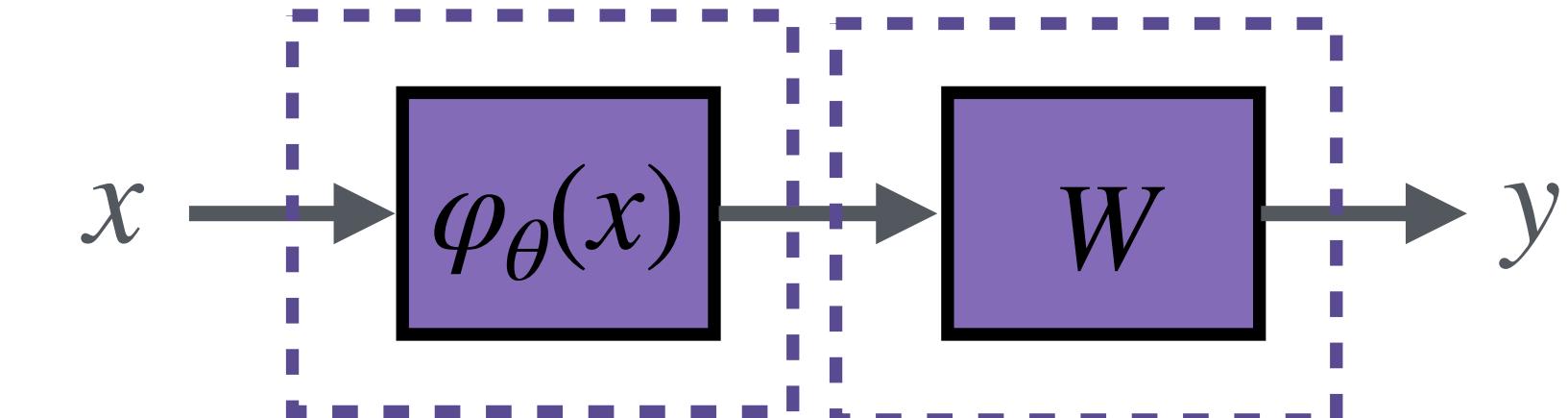
Non-feature learning (old)



only  $W$  is trained

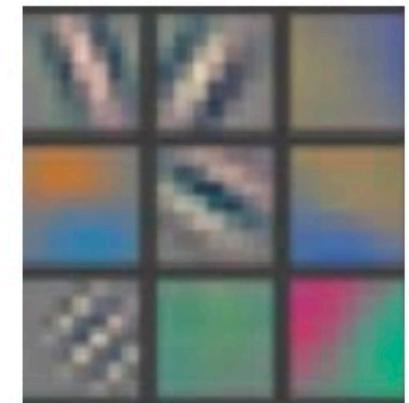
$$f_W(x) = W \circ \varphi(x)$$

Feature learning (new)

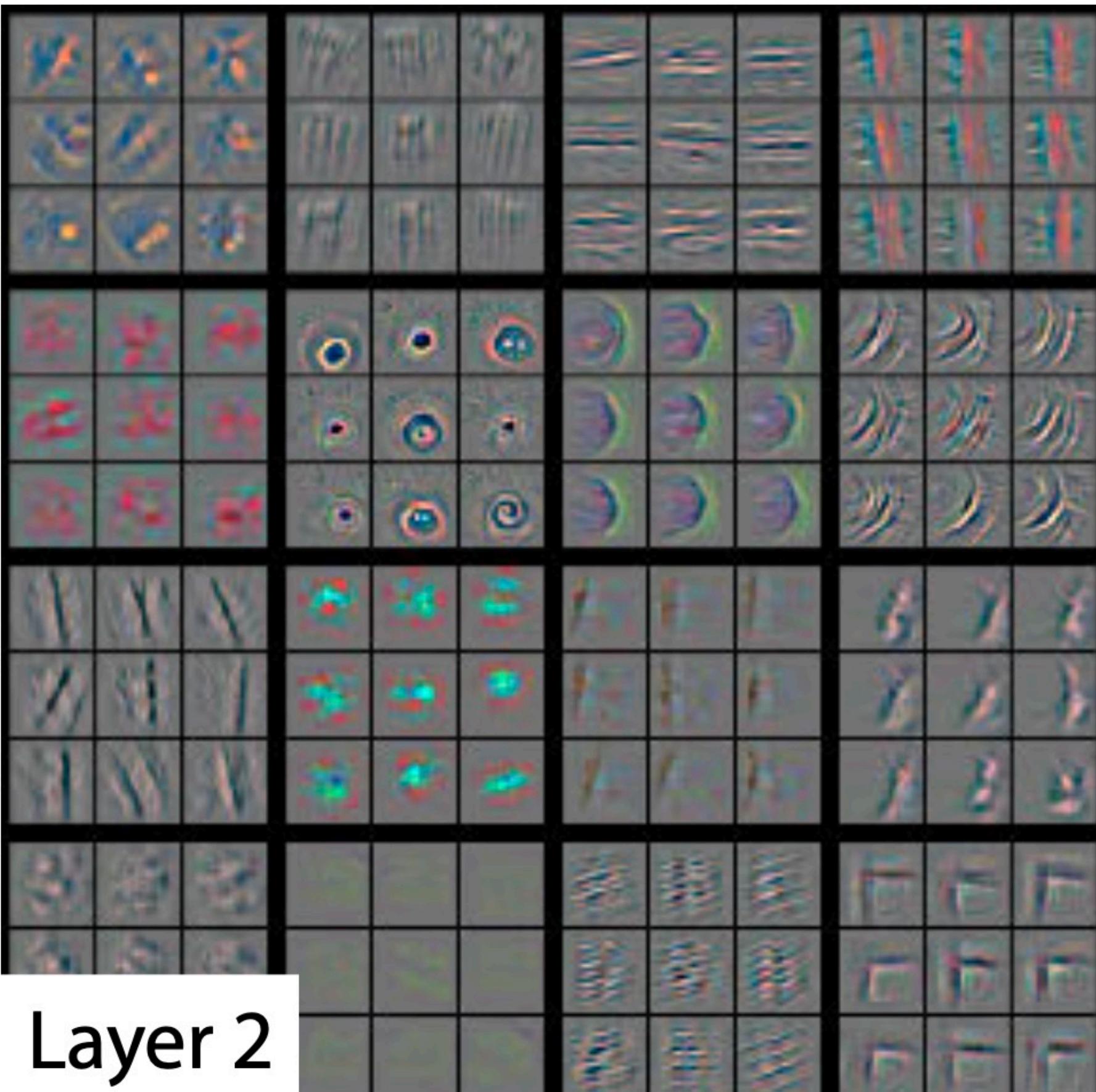


everything learns

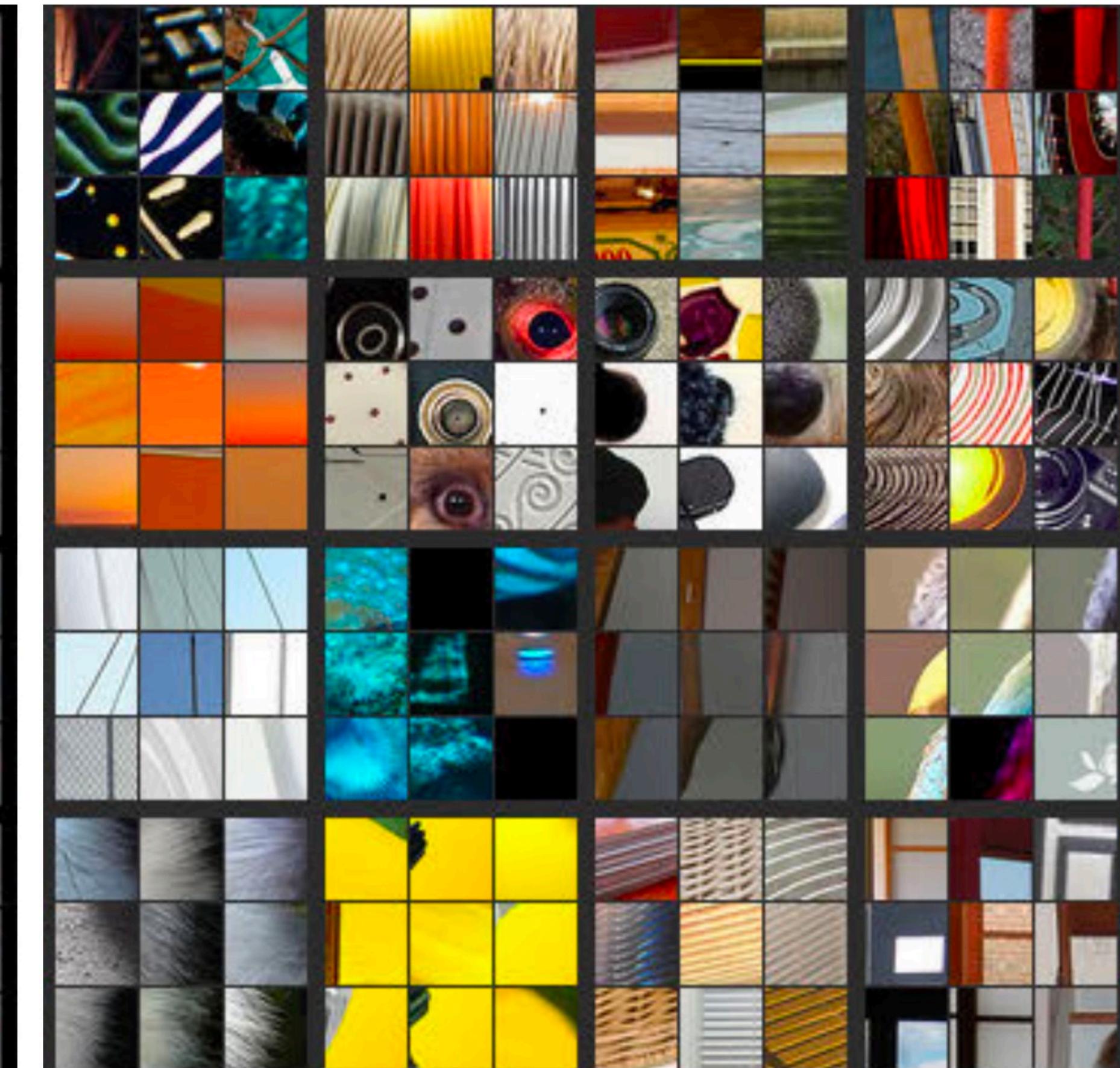
$$f_{W,\theta}(x) = W \circ \varphi_\theta(x)$$

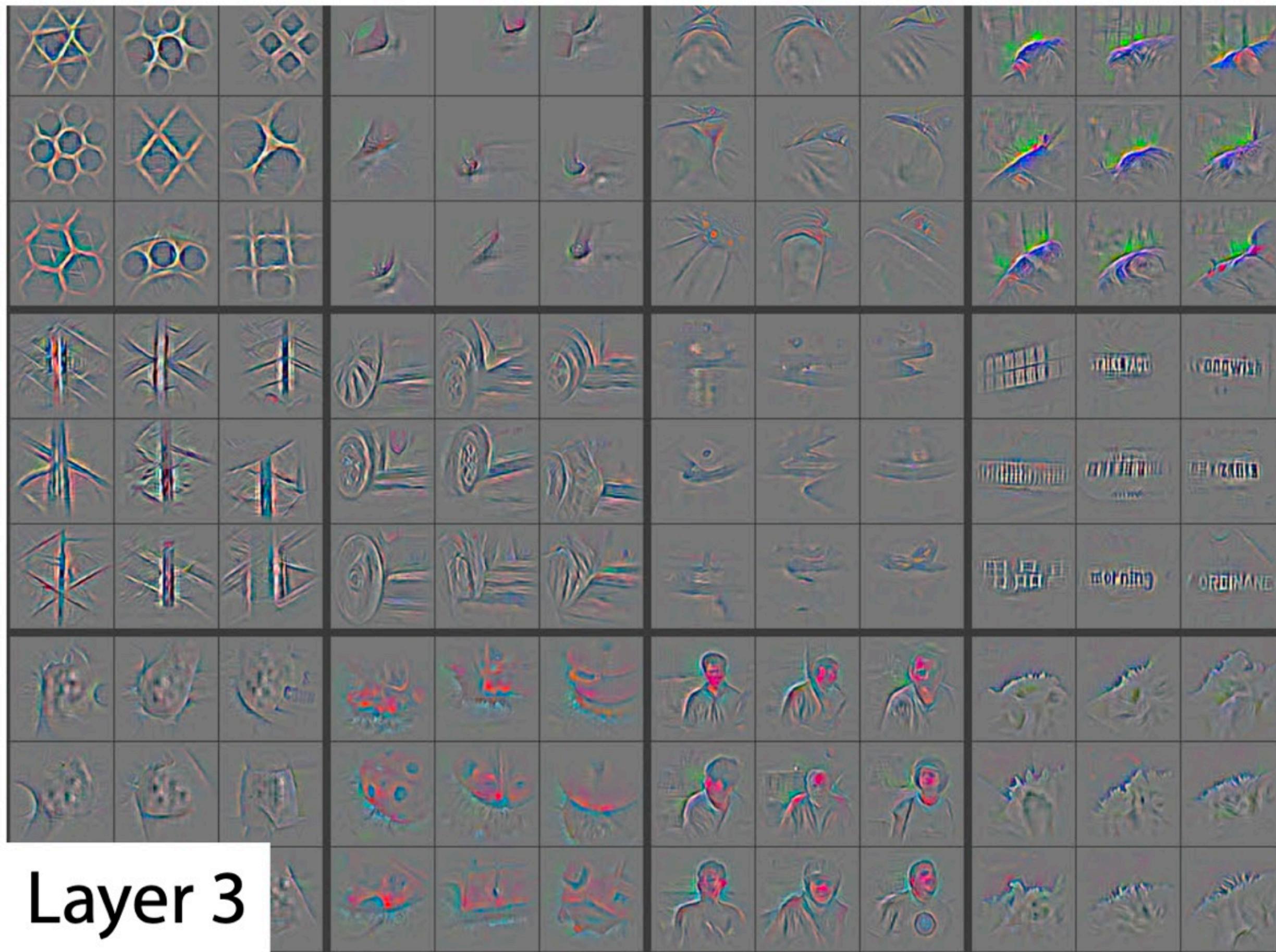


Layer 1

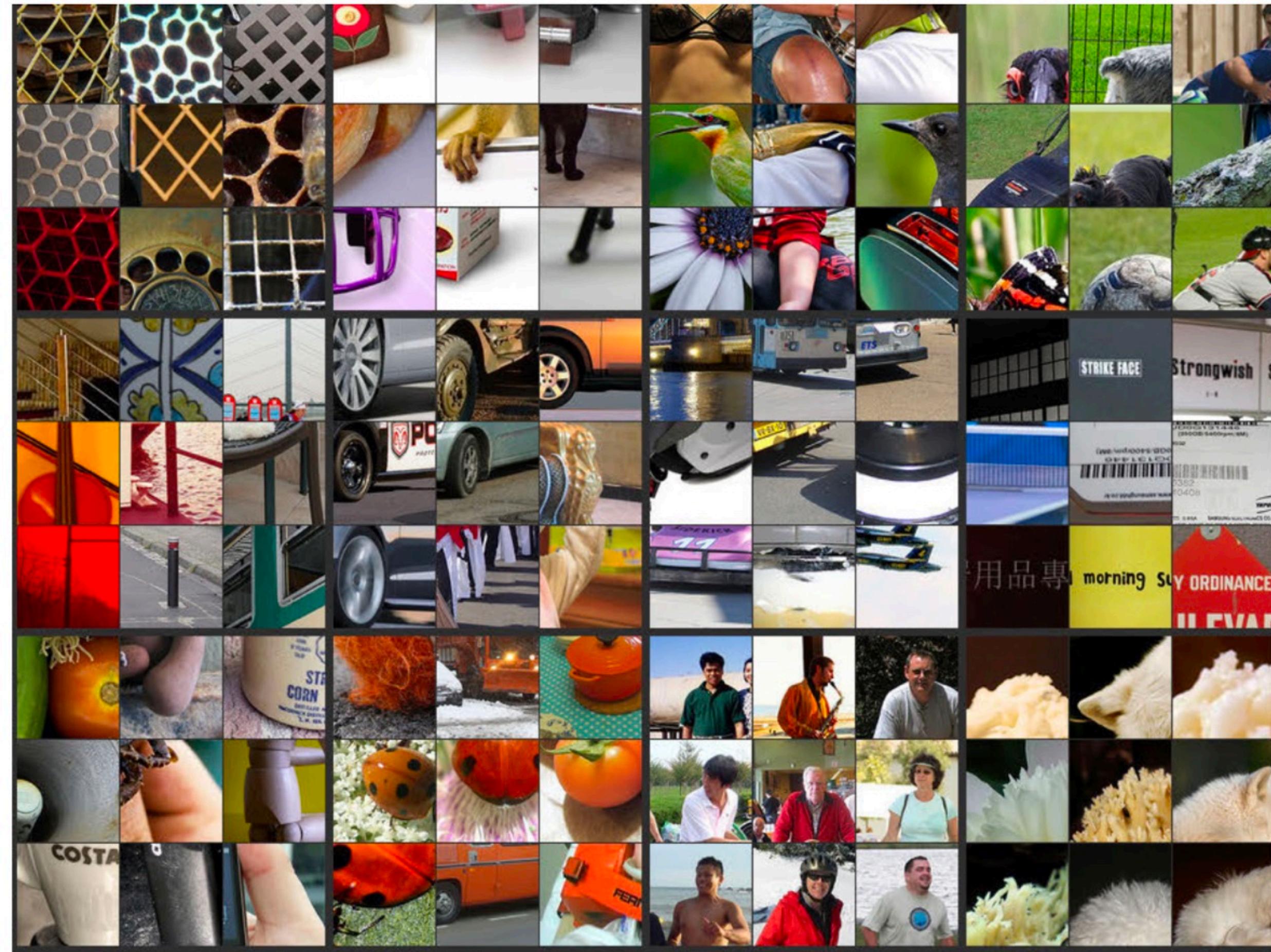


Layer 2





# Layer 3



---

# Feature Learning in Deep Neural Networks – Studies on Speech Recognition Tasks

---

**Dong Yu, Michael L. Seltzer, Jinyu Li<sup>1</sup>, Jui-Ting Huang<sup>1</sup>, Frank Seide<sup>2</sup>**

Microsoft Research, Redmond, WA 98052

<sup>1</sup>Microsoft Corporation, Redmond, WA 98052

<sup>2</sup>Microsoft Research Asia, Beijing, P.R.C.

{dongyu, mseltzer, jinyu, jthuang, fseide}@microsoft.com

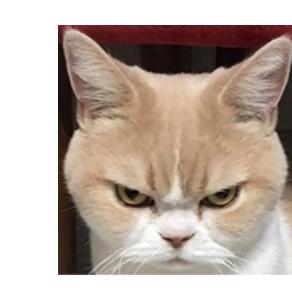
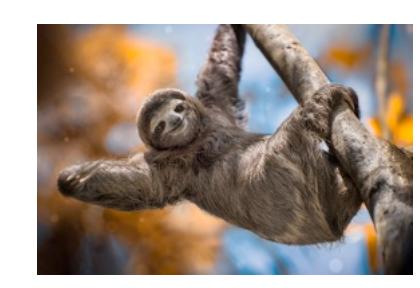
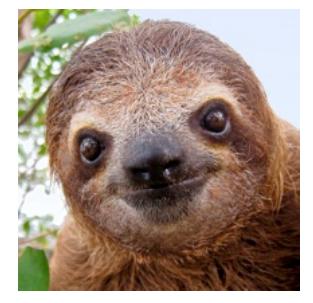
Chen et al. 2022

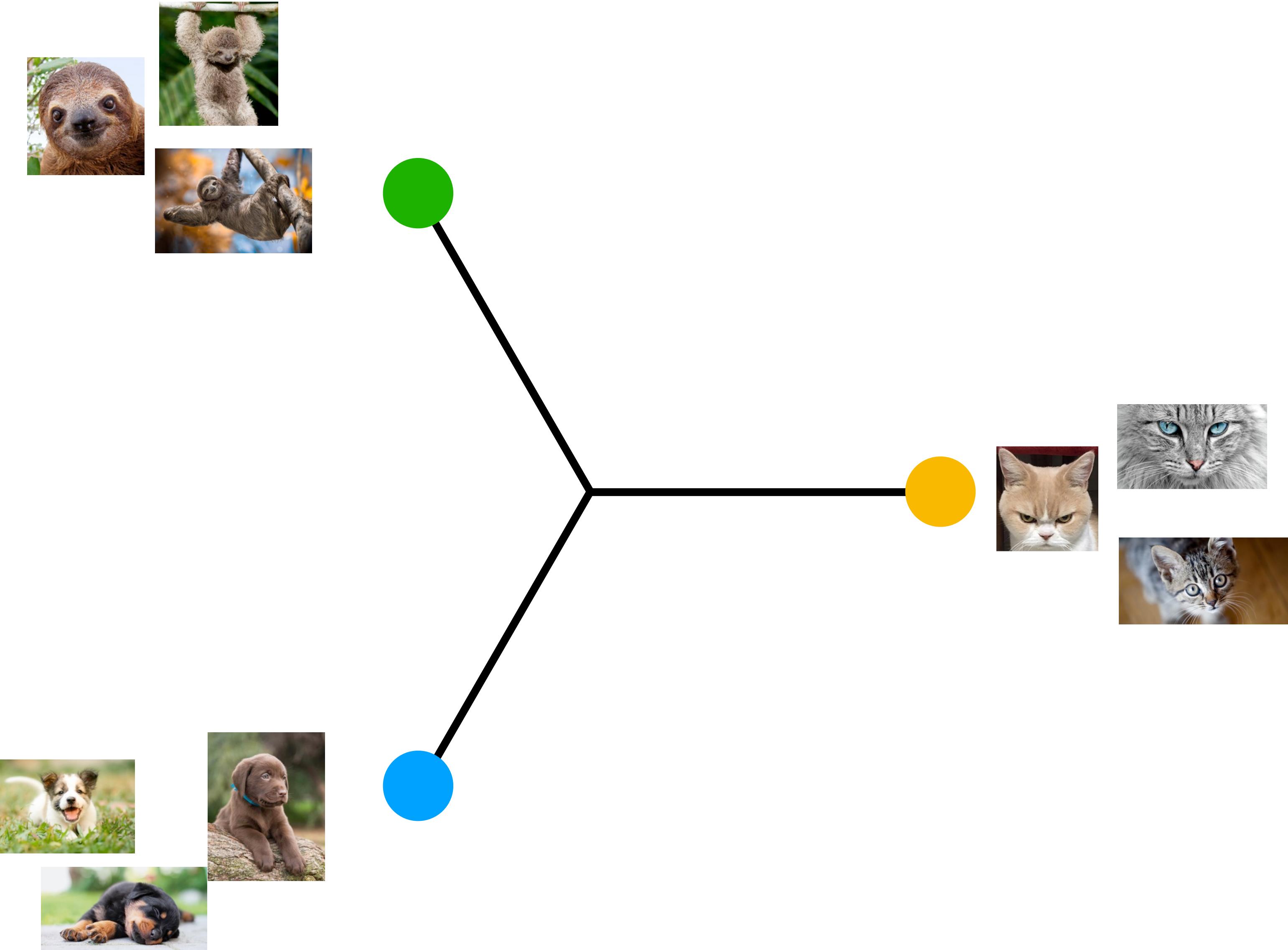
Nichani et al. 2023

Cui et al. 2024

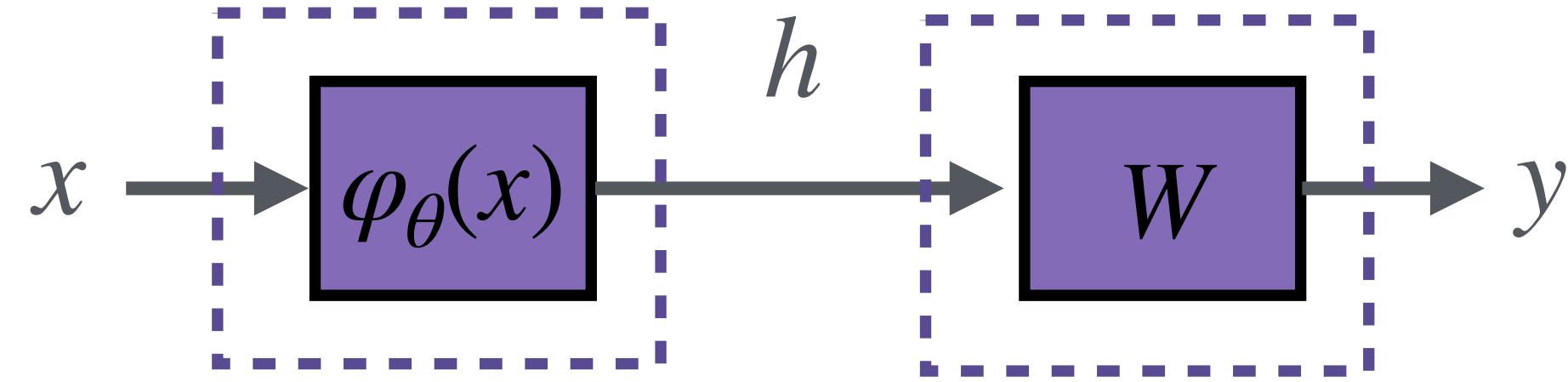
A series of works from Misha Belkin, Cengiz, and many others...

...





# NEURAL COLLAPSE

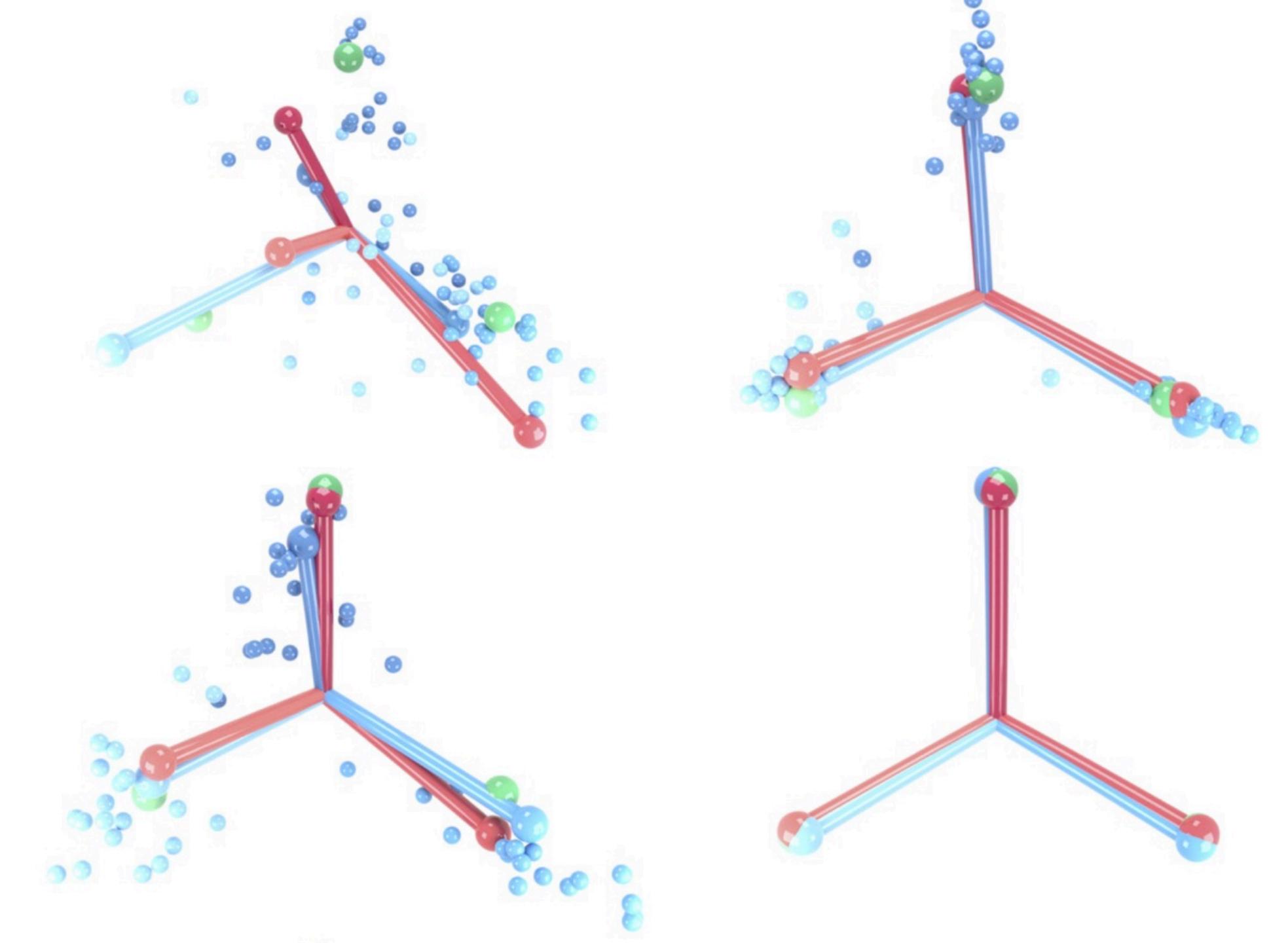


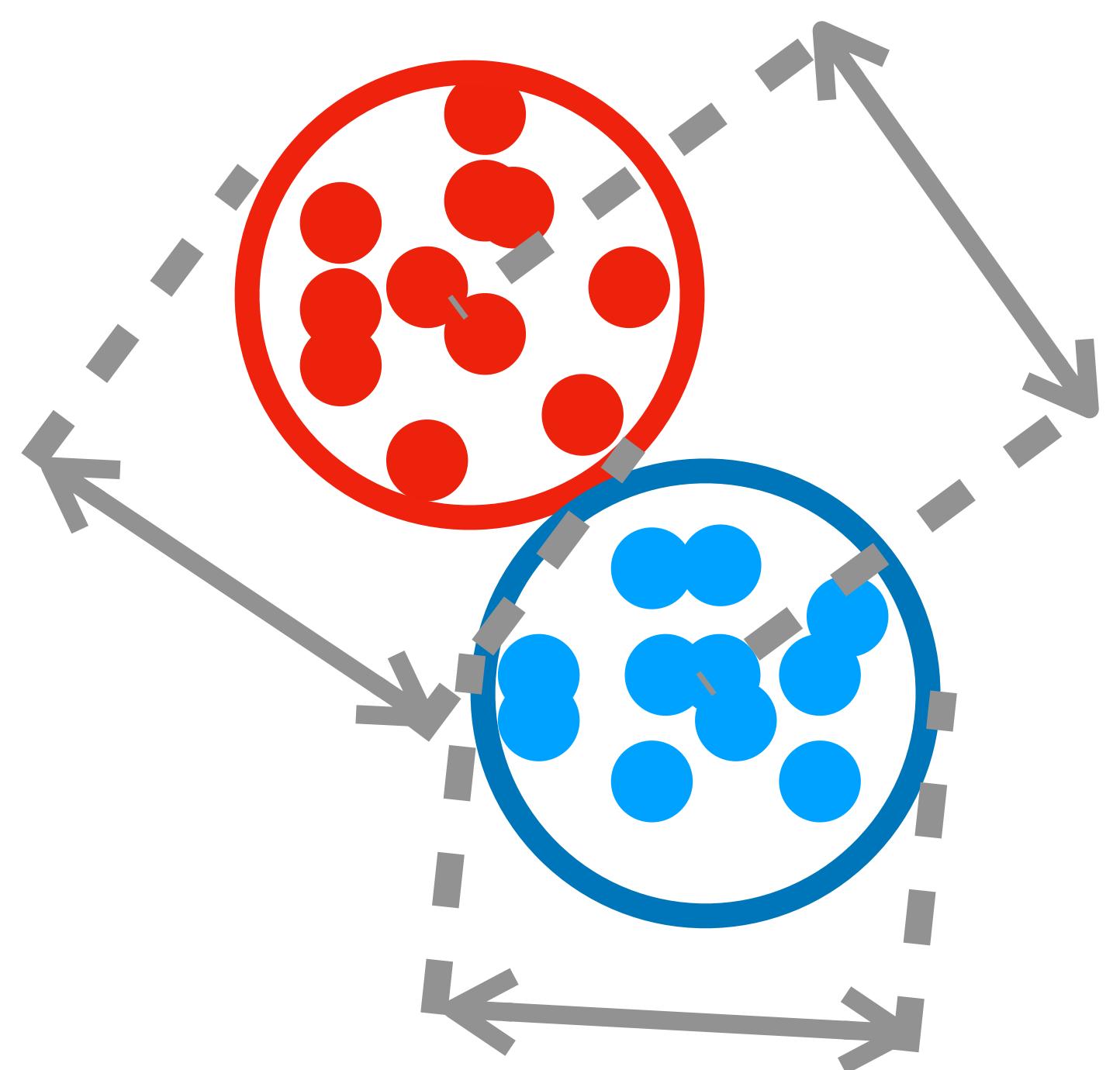
$K$  classes, each with  $N_k$  examples

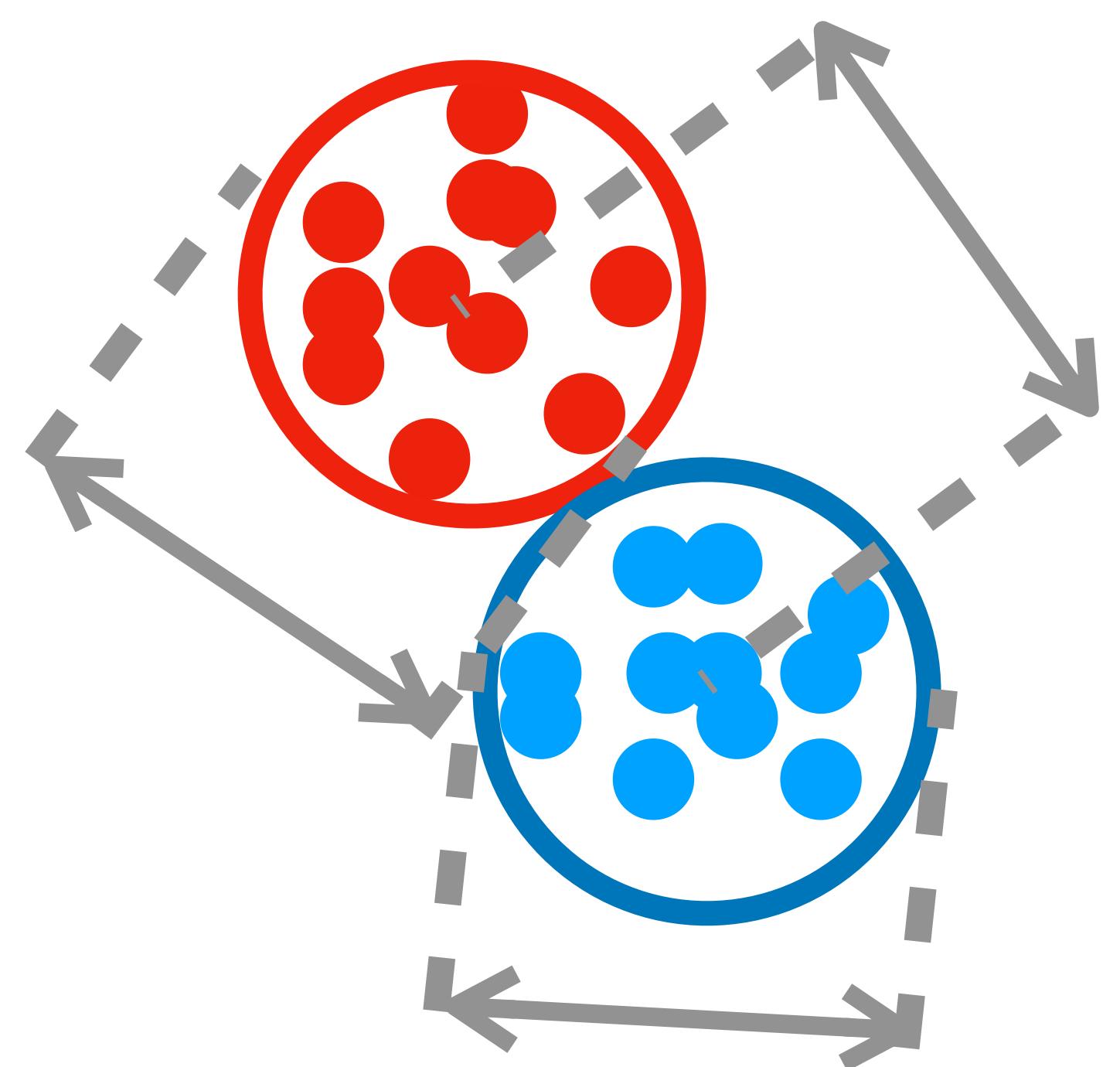
$$\Sigma_k = \text{Cov}(h_{kn})_{n=1}^{N_k} \quad \mu_k = \text{Ave}_n(h_{kn})$$

$$\Sigma_w = \text{Ave}_k(\Sigma_k) \quad \Sigma_b = \text{Cov}(\mu_k)$$

$$\text{Tr}(\Sigma_w \Sigma_b^+) \rightarrow 0$$

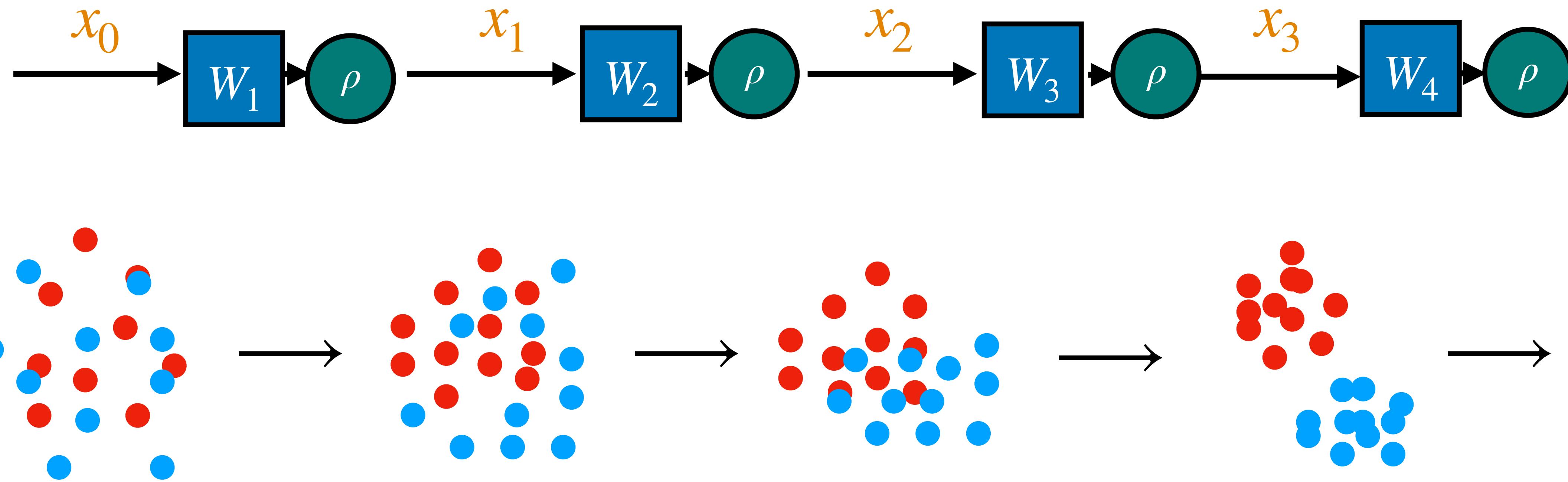






$$D = \log \frac{\text{within-class variance}}{\text{between-class variance}}$$

# DEEP NEURAL COLLAPSE



---

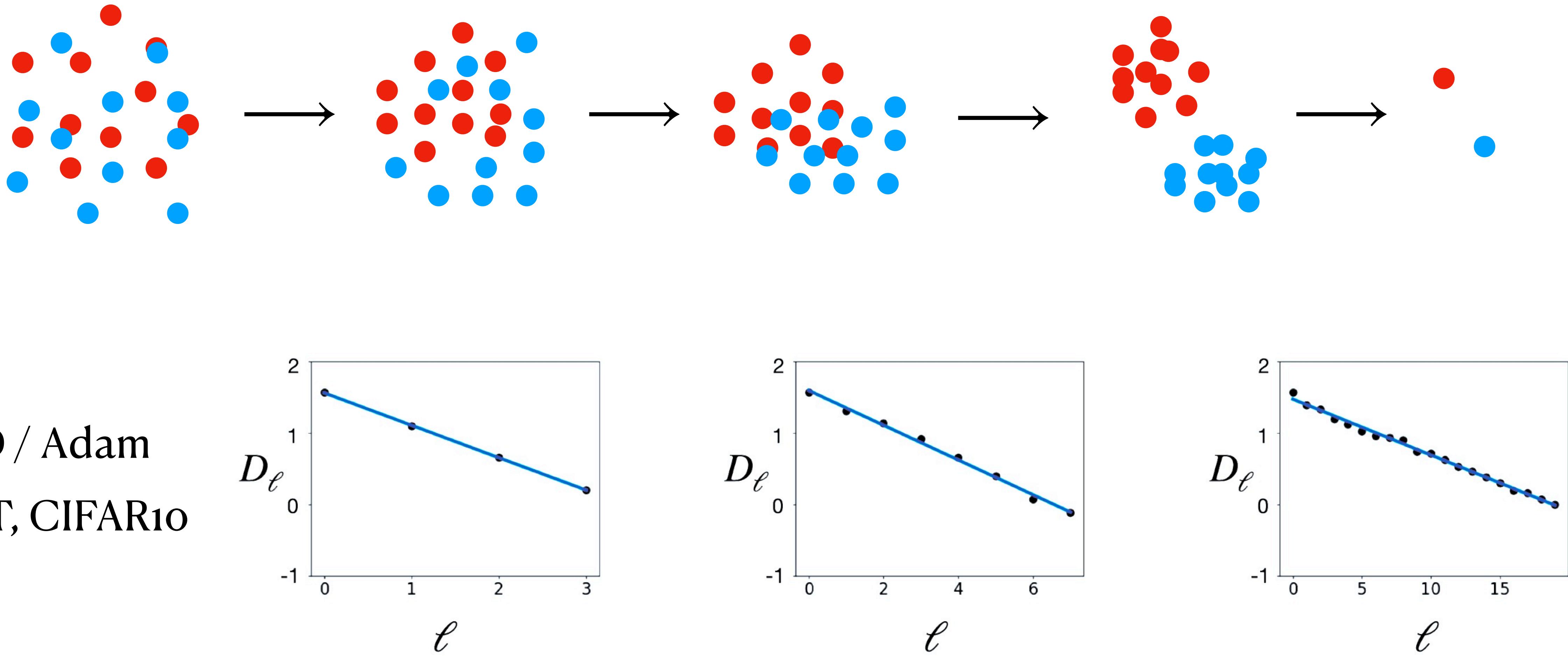
[Papyan, Han, Donoho, (PNAS 2022)]

[Súkeník, Mondelli, Lampert, Nuemr (NeurIPS 2023)]

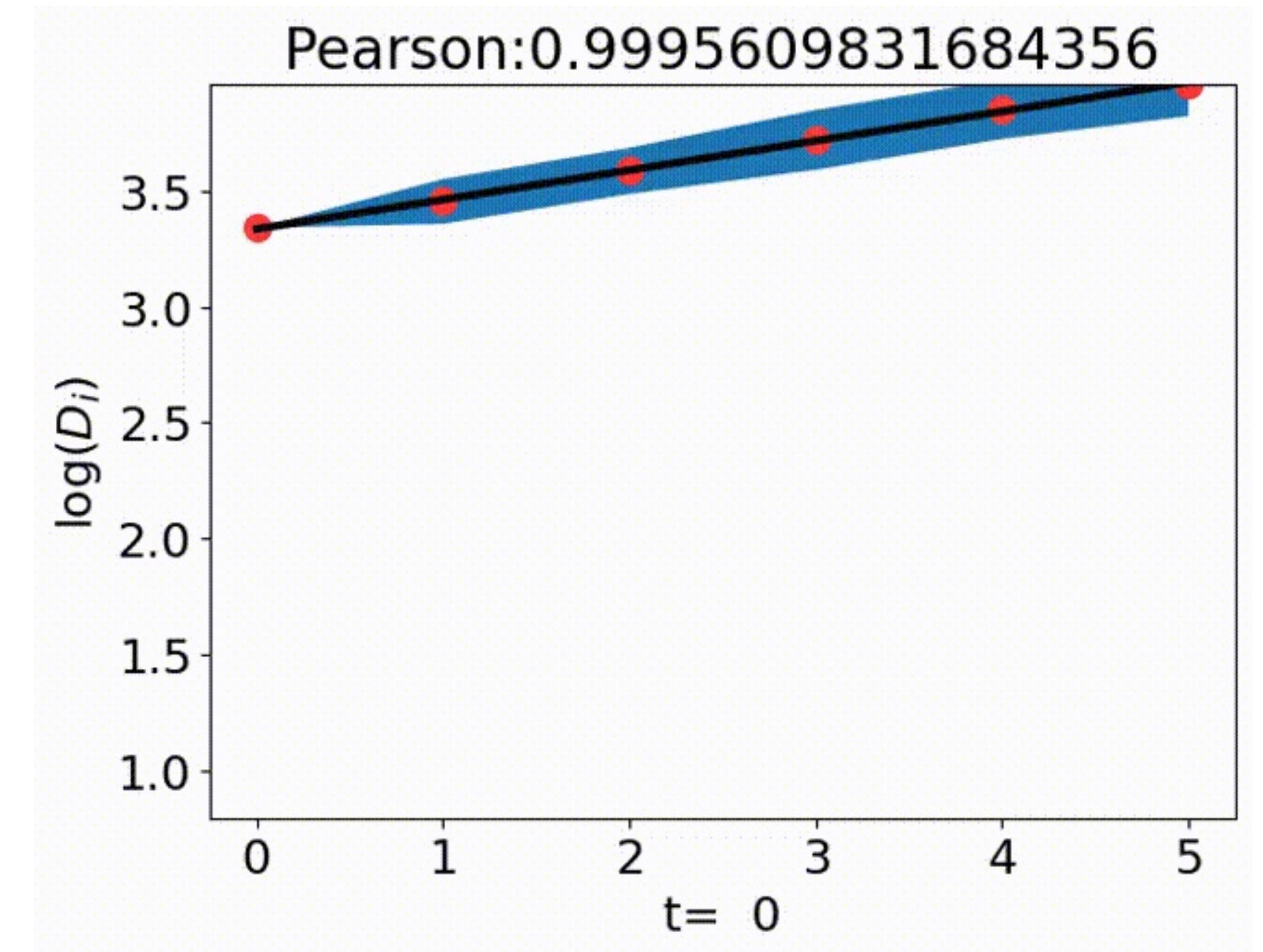
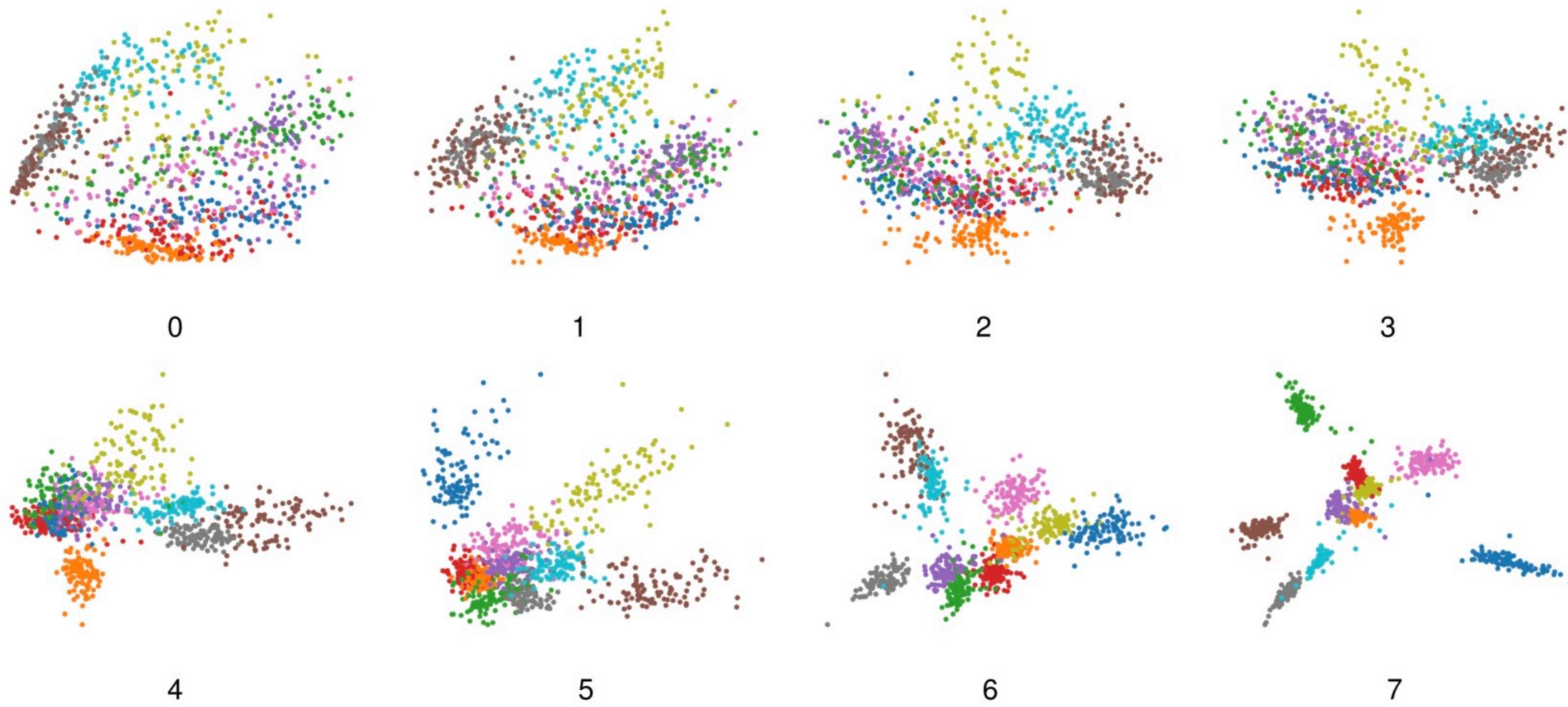
[Rangamani et al. (ICML 2023)]

# A LAW OF DATA SEPARATION

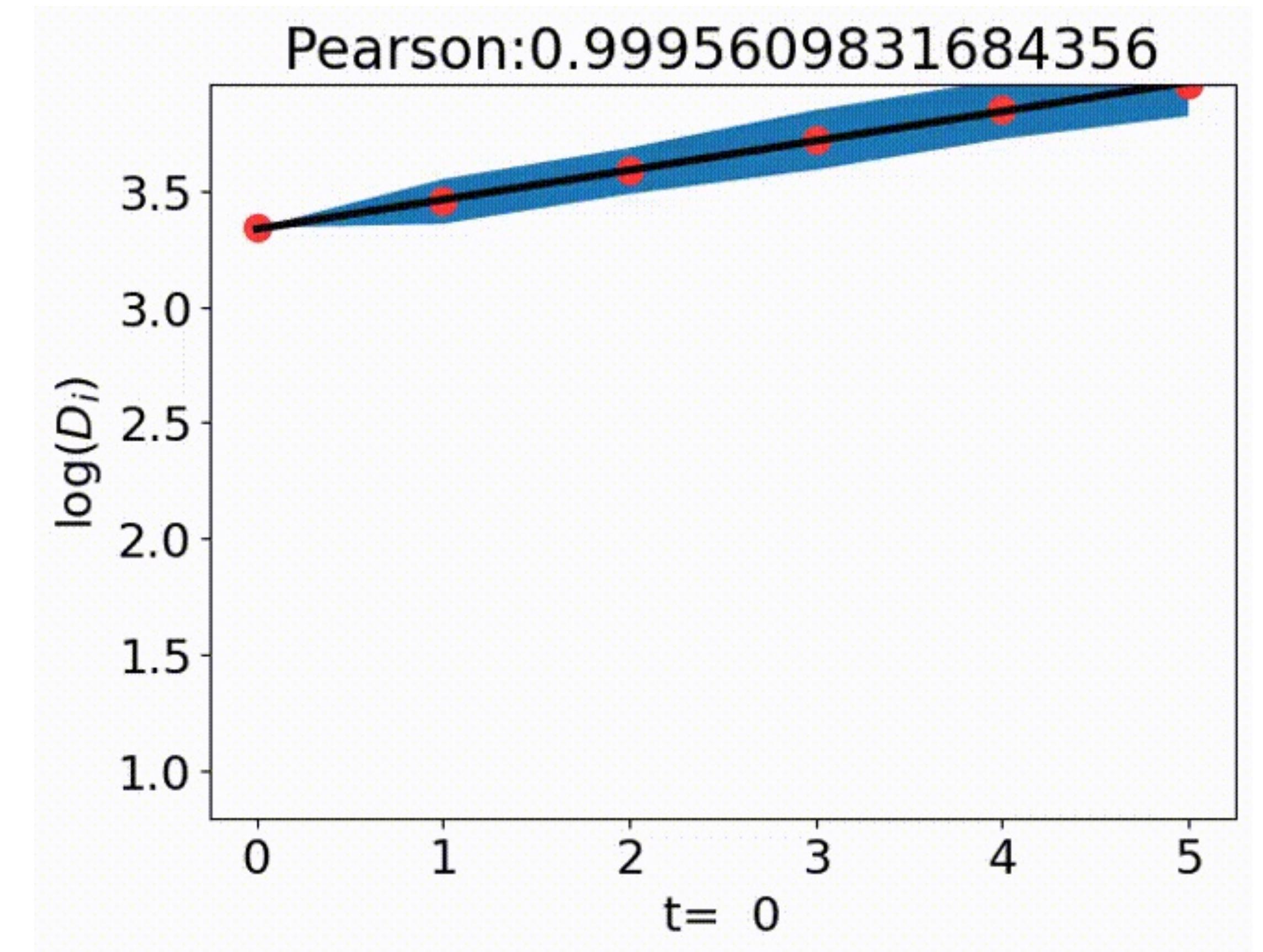
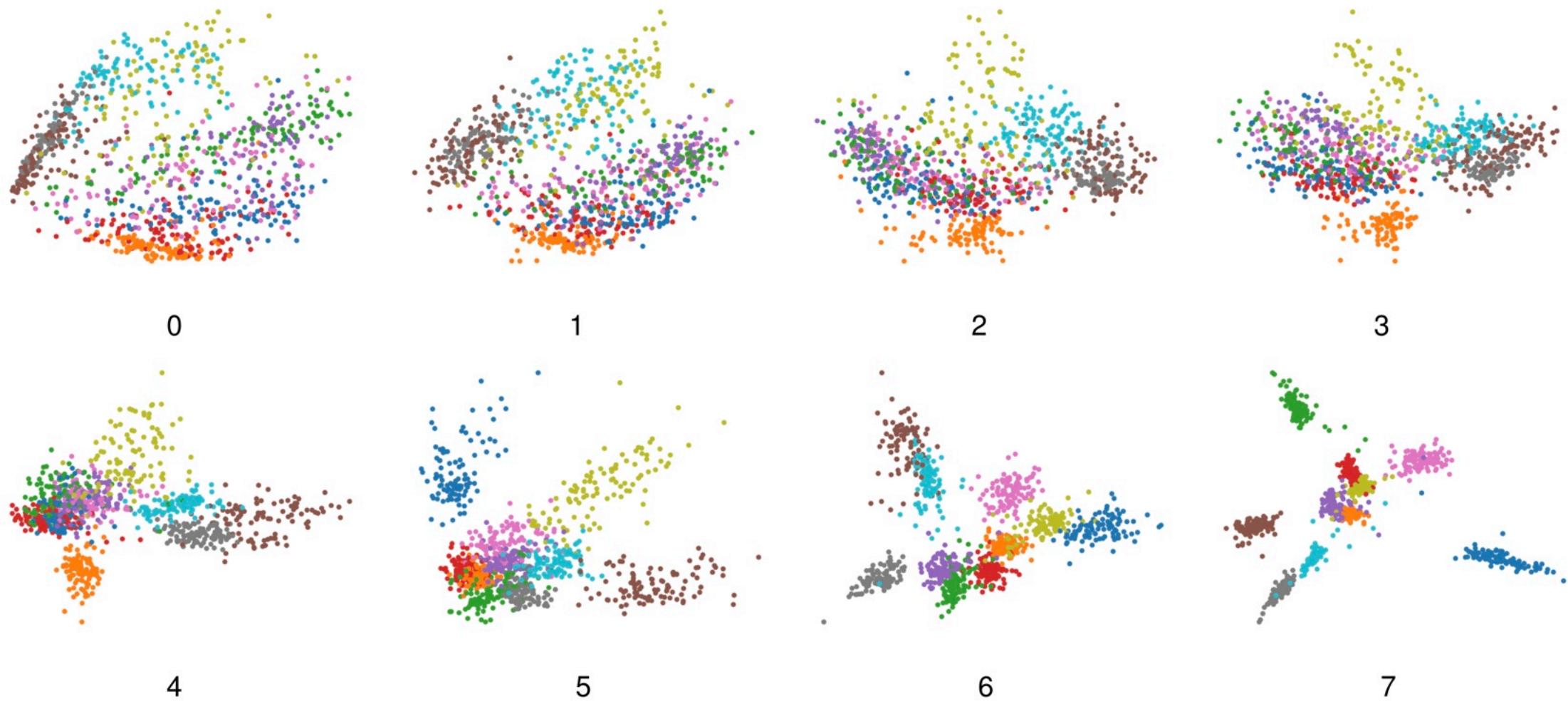
HE, SU; PNAS 2023



$$D_\ell = \text{trace } \Sigma_w \Sigma_b^+ \quad \text{or} \quad D_\ell = \text{trace } \Sigma_w / \text{trace} \Sigma_b$$

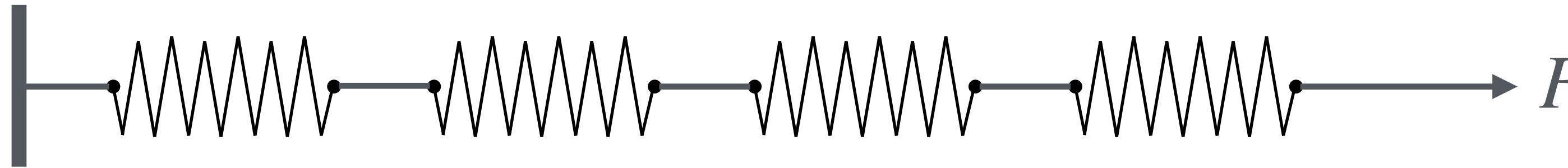


$$D_\ell = \text{trace } \Sigma_w \Sigma_b^+ \quad \text{or} \quad D_\ell = \text{trace } \Sigma_w / \text{trace} \Sigma_b$$



# A 1<sup>ST</sup> MECHANICAL ANALOGY

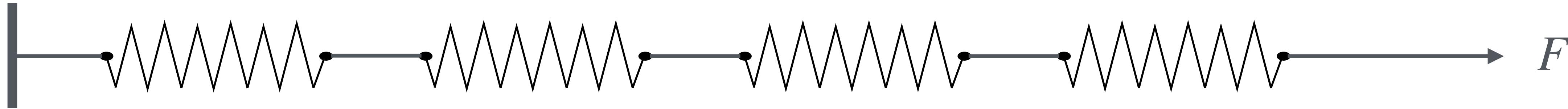
EQUISEPARATION IS SOMEHOW A LINEAR PHENOMENON



$$h_L = \frac{LF}{k} \qquad w_\ell = \frac{F}{k}$$

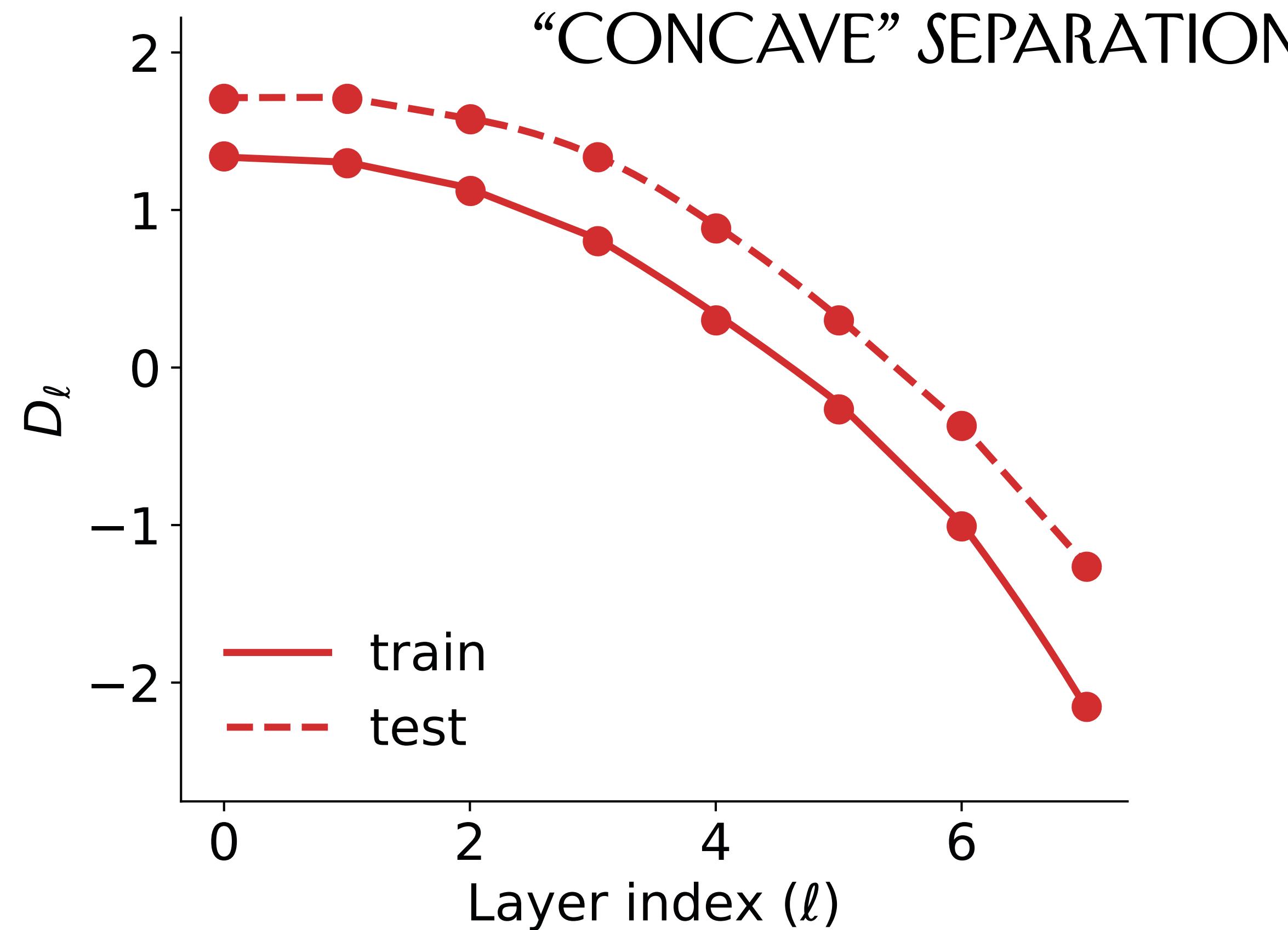
# A MECHANICAL ANALOGY

EQUISEPARATION IS SOMEHOW A LINEAR PHENOMENON



$$h_L = \frac{LF}{k} \qquad w_\ell = \frac{F}{k}$$

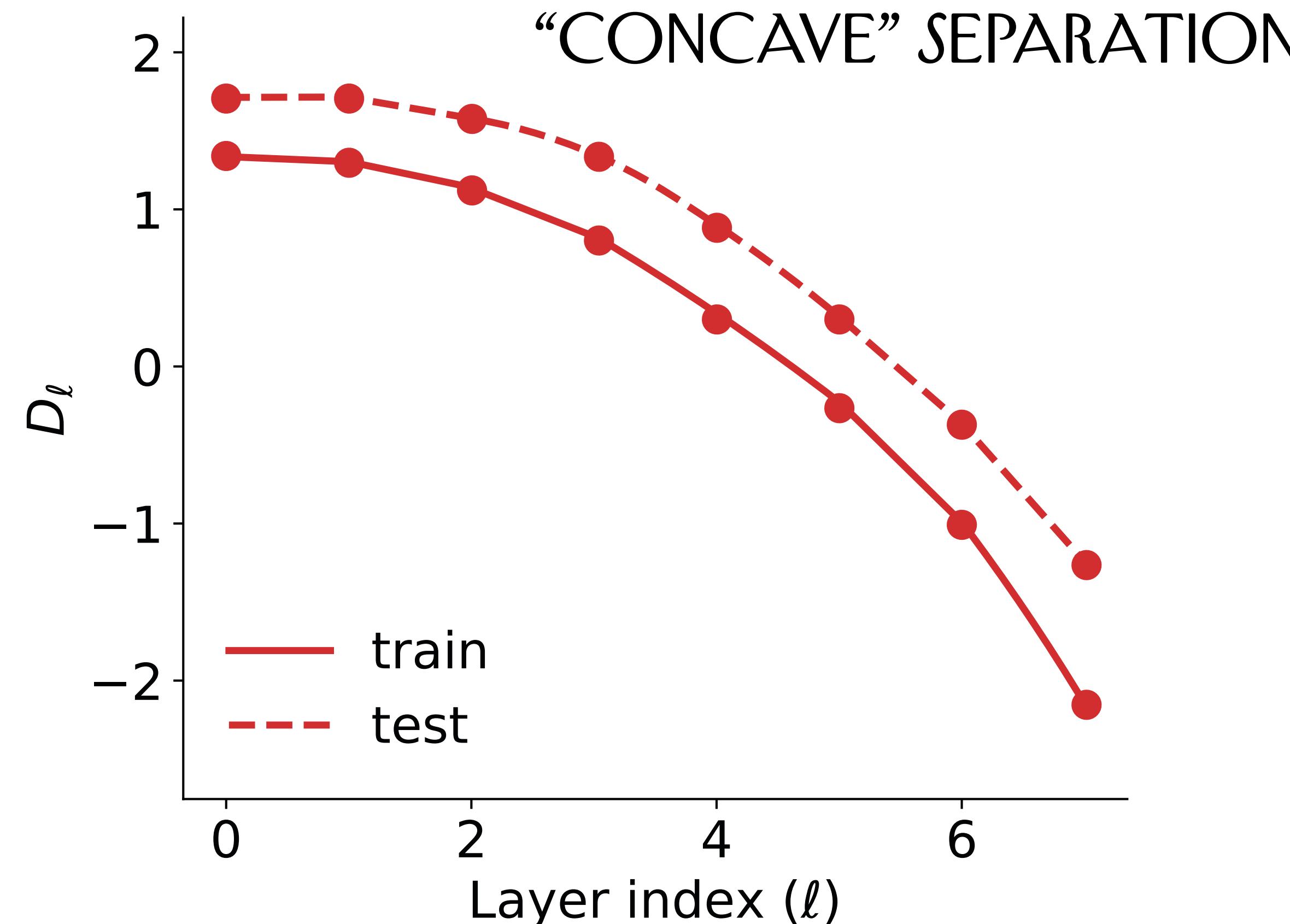
# A LAW OF DATA SEPARATION?



8-layer ReLU ResNet | MNIST | 10000  
examples | full batch | SGD | lr = 0.0001,

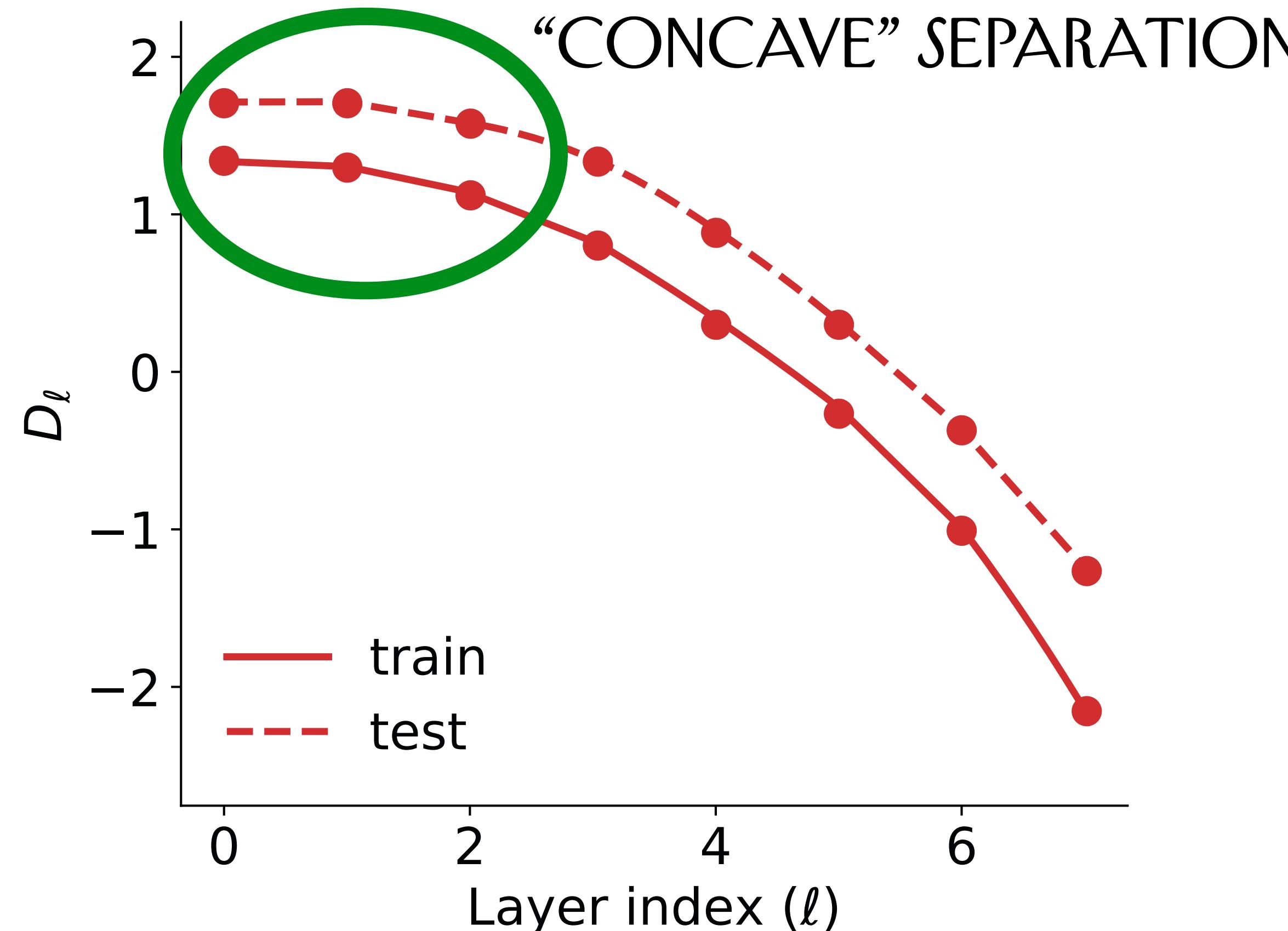
acc = 82.49%

# A LAW OF DATA SEPARATION?



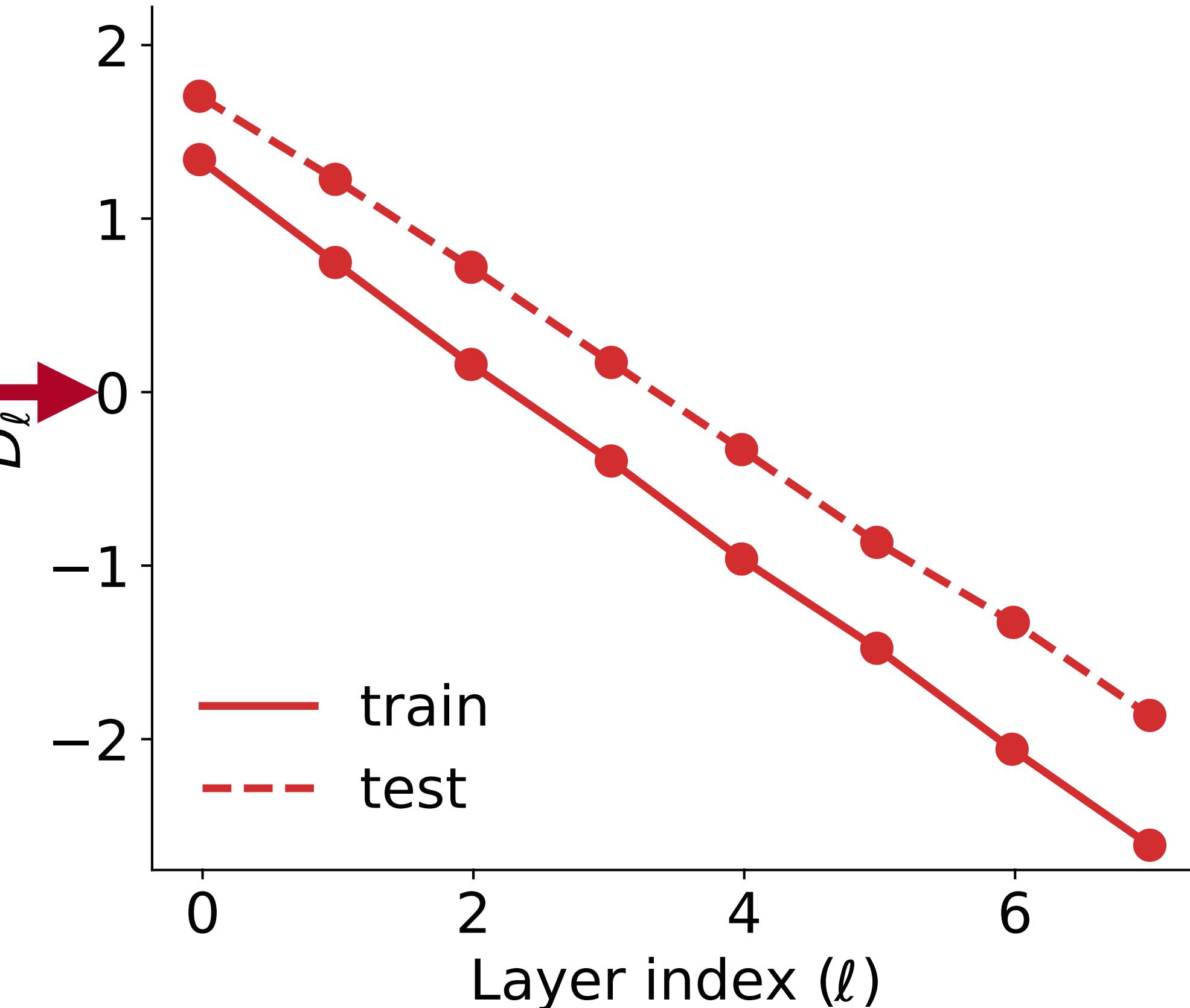
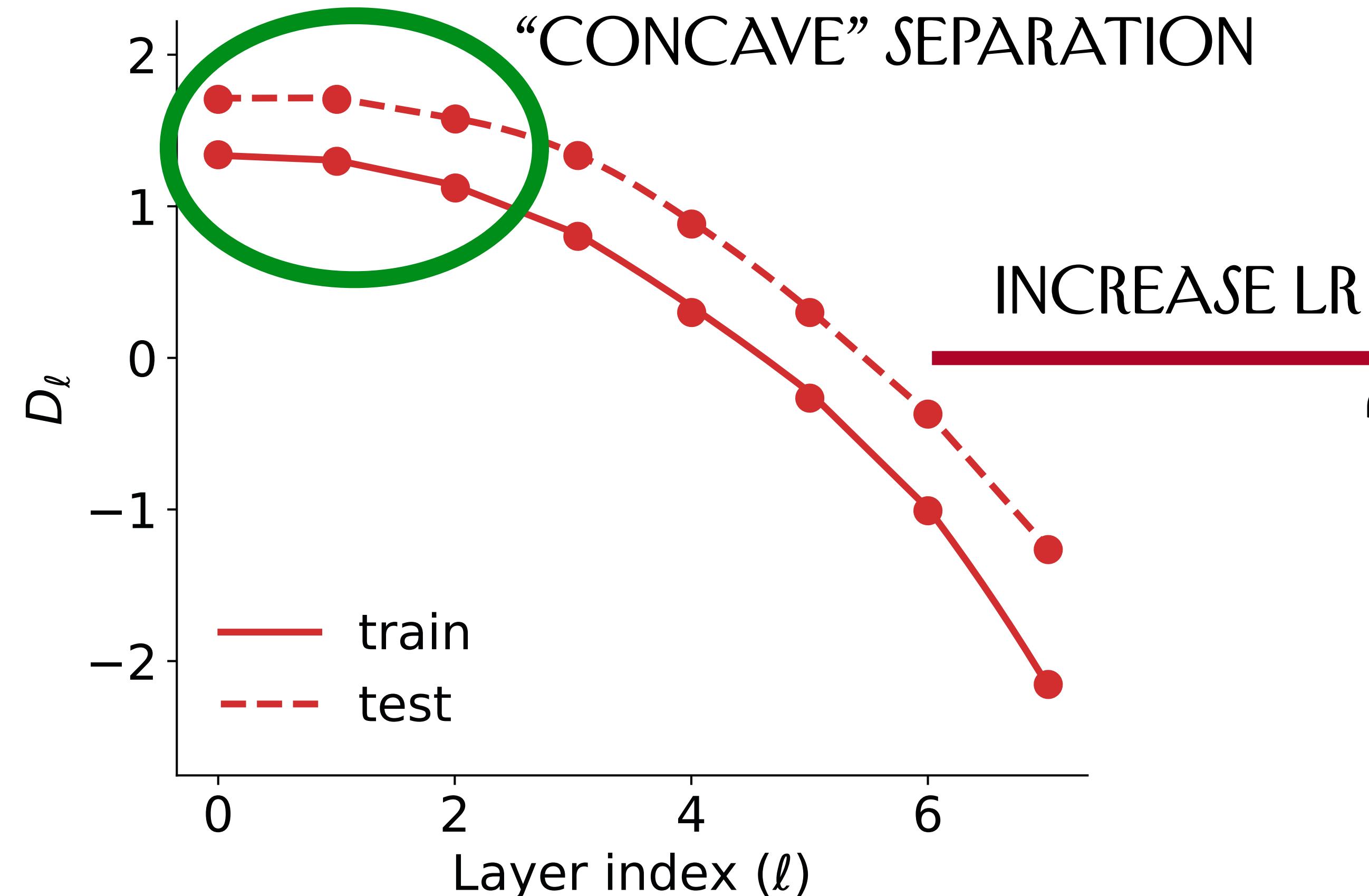
8-layer ReLU ResNet | MNIST | 10000  
examples | full batch | SGD | lr = 0.0001,  
acc = 82.49%

# A LAW OF DATA SEPARATION?



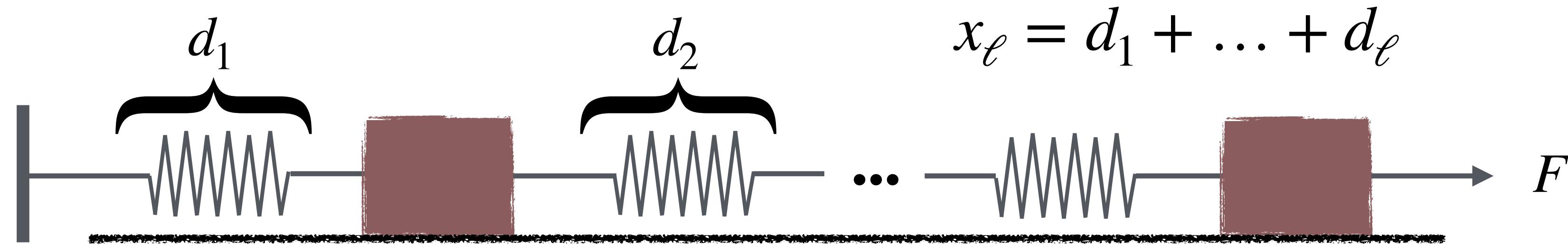
8-layer ReLU ResNet | MNIST | 10000  
examples | full batch | SGD | lr = 0.0001,  
acc = 82.49%

# A LAW OF DATA SEPARATION?

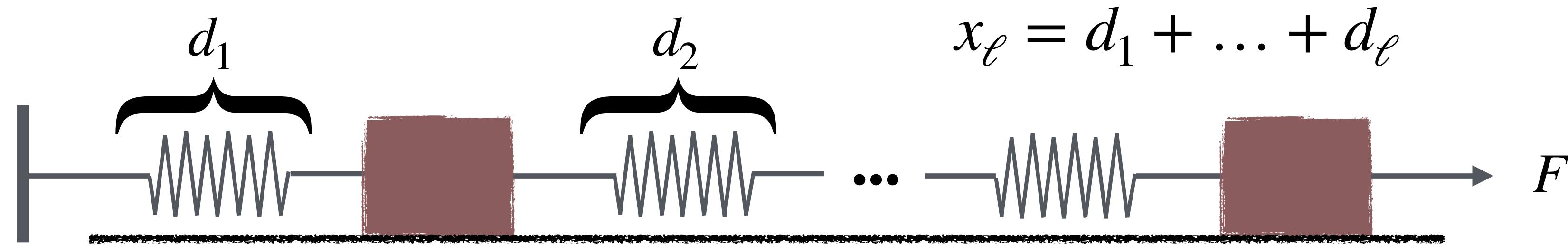


8-layer ReLU ResNet | MNIST | 10000  
examples | full batch | SGD | lr = 0.0001,  
acc = 82.49%

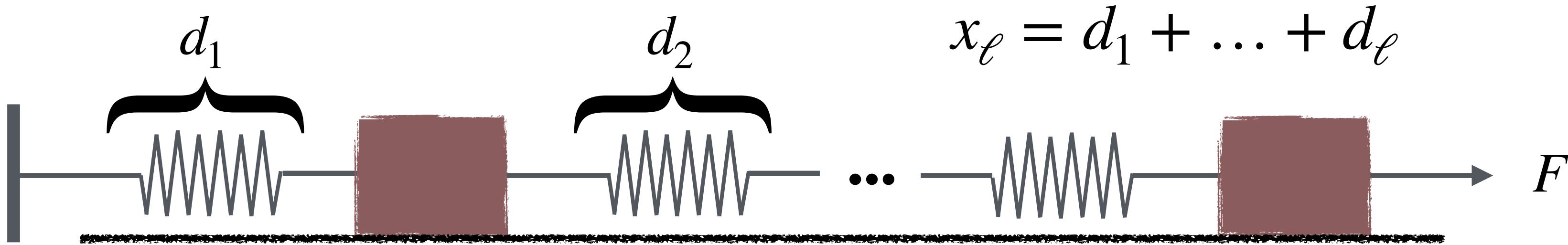
$$\ddot{x}_\ell = k(d_{\ell+1} - d_\ell) - \gamma \dot{x}_\ell - f_\ell$$



$$\ddot{x}_\ell = k(d_{\ell+1} - d_\ell) - \gamma \dot{x}_\ell - f_\ell$$

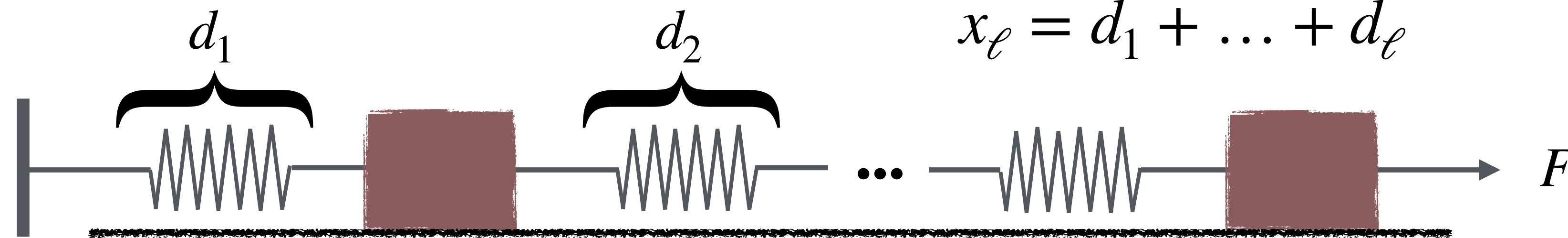


$$\ddot{x}_\ell = k(d_{\ell+1} - d_\ell) - \gamma \dot{x}_\ell - f_\ell$$

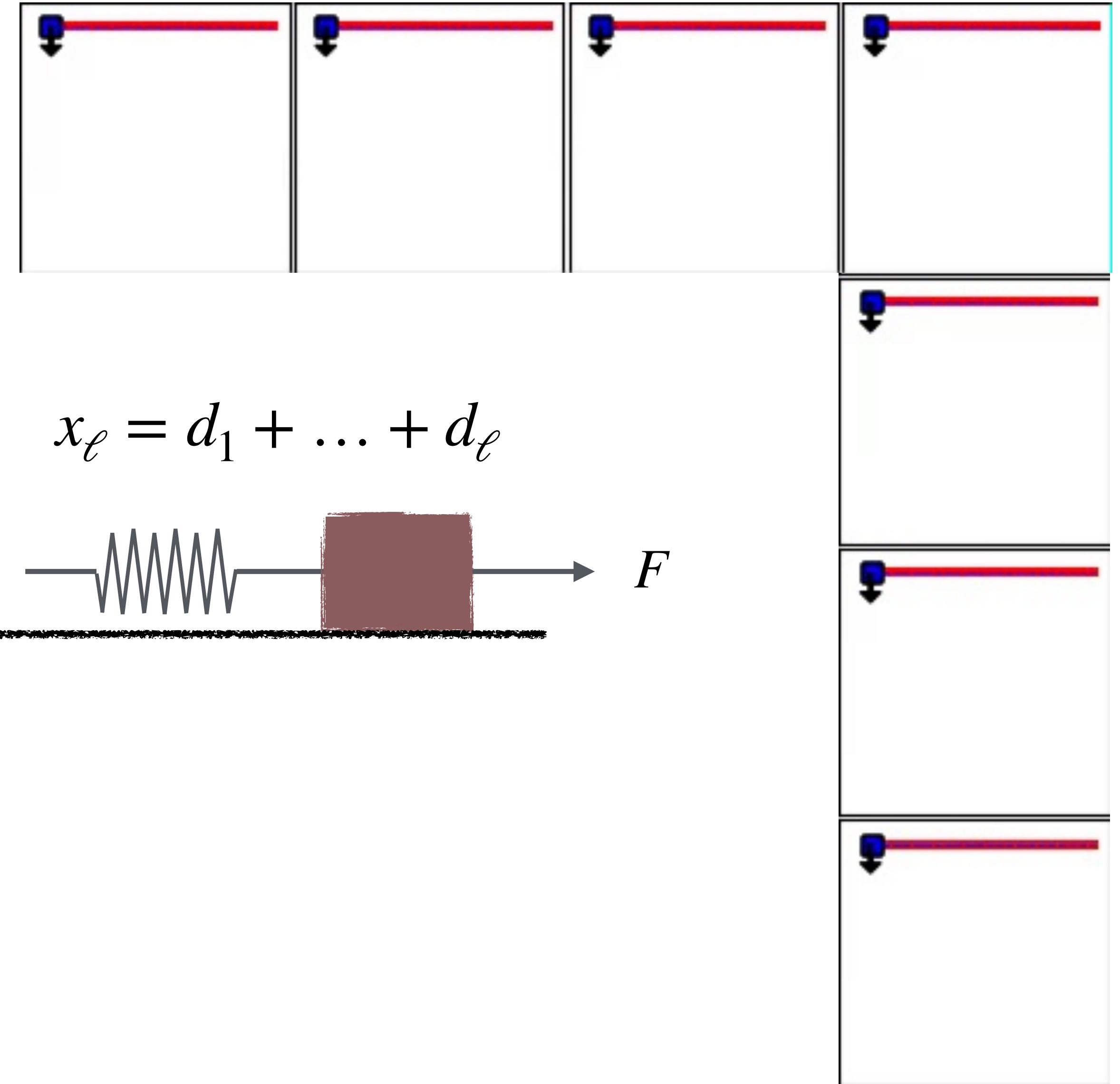


$$\ddot{x}_\ell = k(d_{\ell+1} - d_\ell) - \gamma \dot{x}_\ell - f_\ell + \sigma \xi_\ell$$

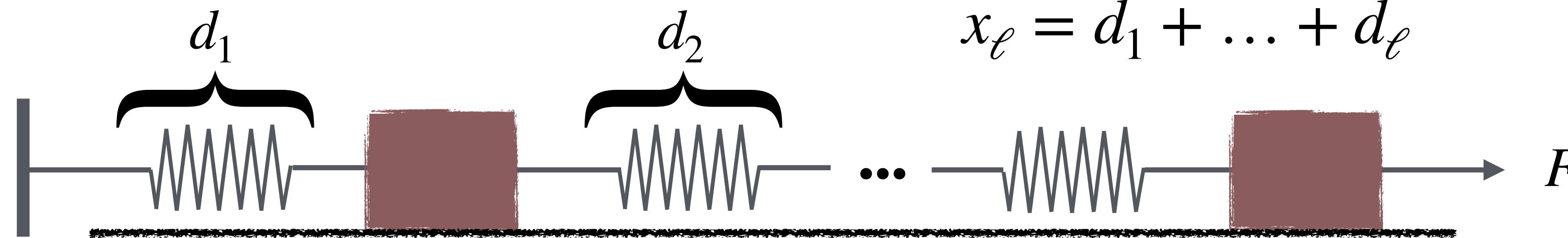
$$\ddot{x}_\ell = k(d_{\ell+1} - d_\ell) - \gamma \dot{x}_\ell - f_\ell$$



$$\ddot{x}_\ell = k(d_{\ell+1} - d_\ell) - \gamma \dot{x}_\ell - f_\ell + \sigma \xi_\ell$$



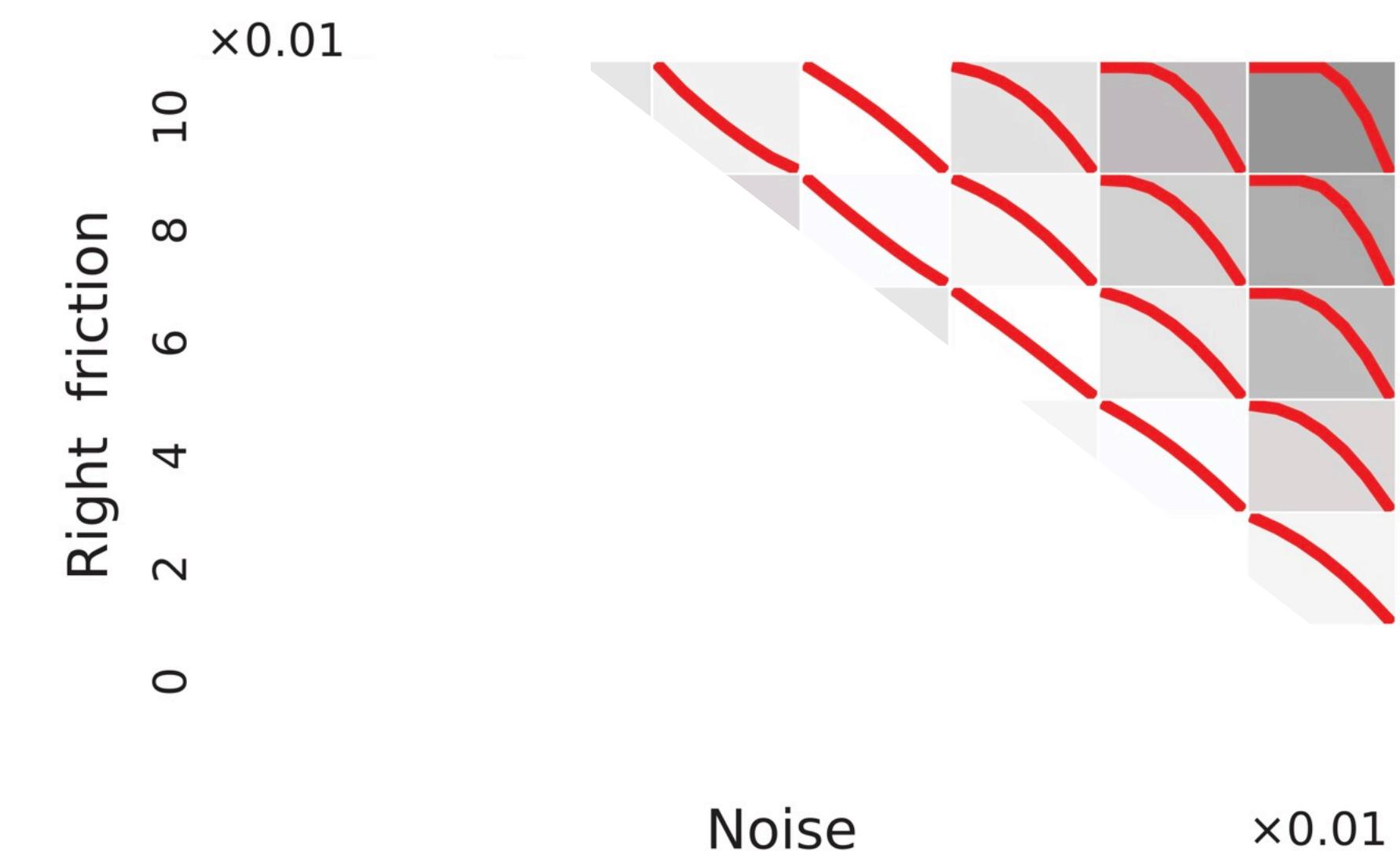
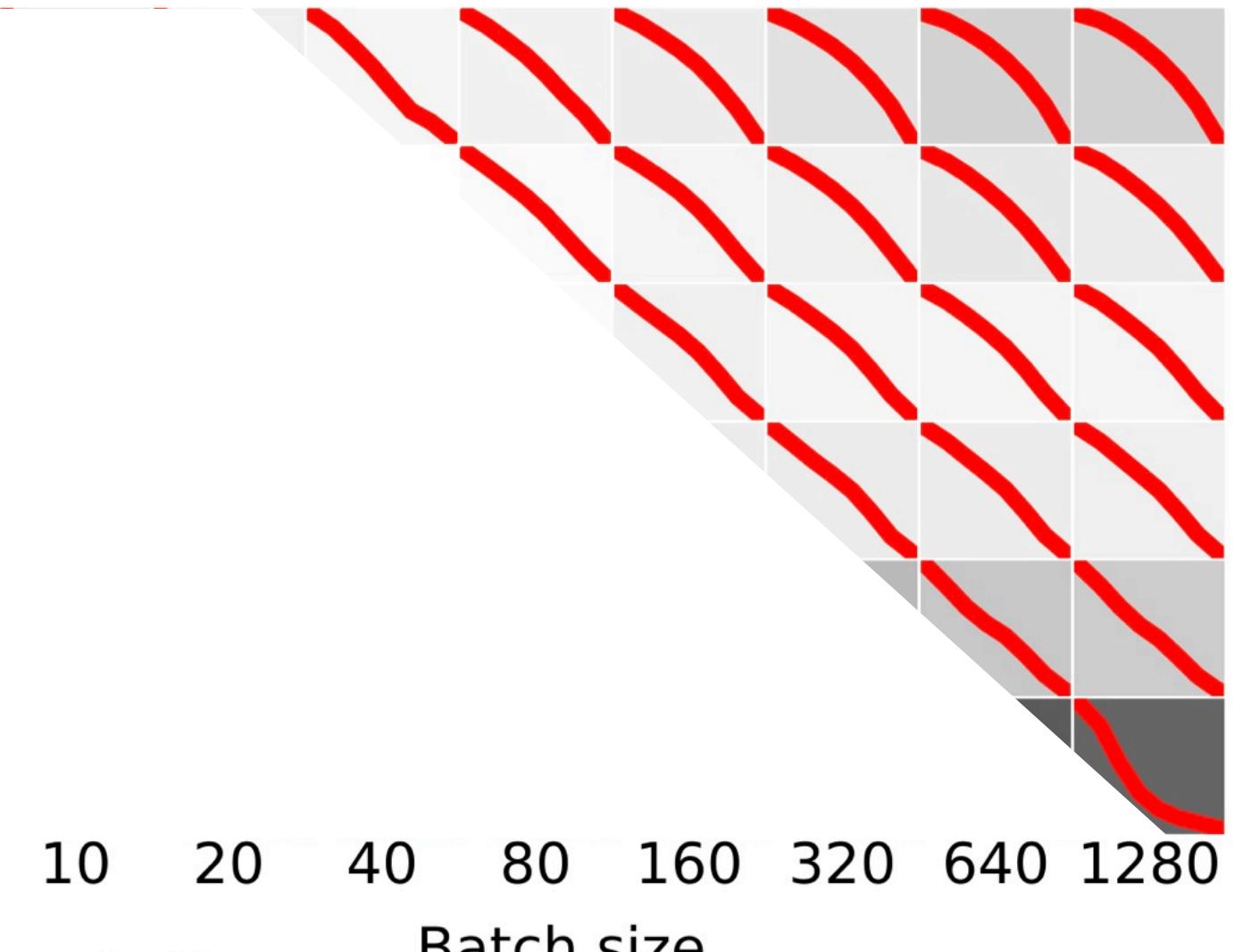
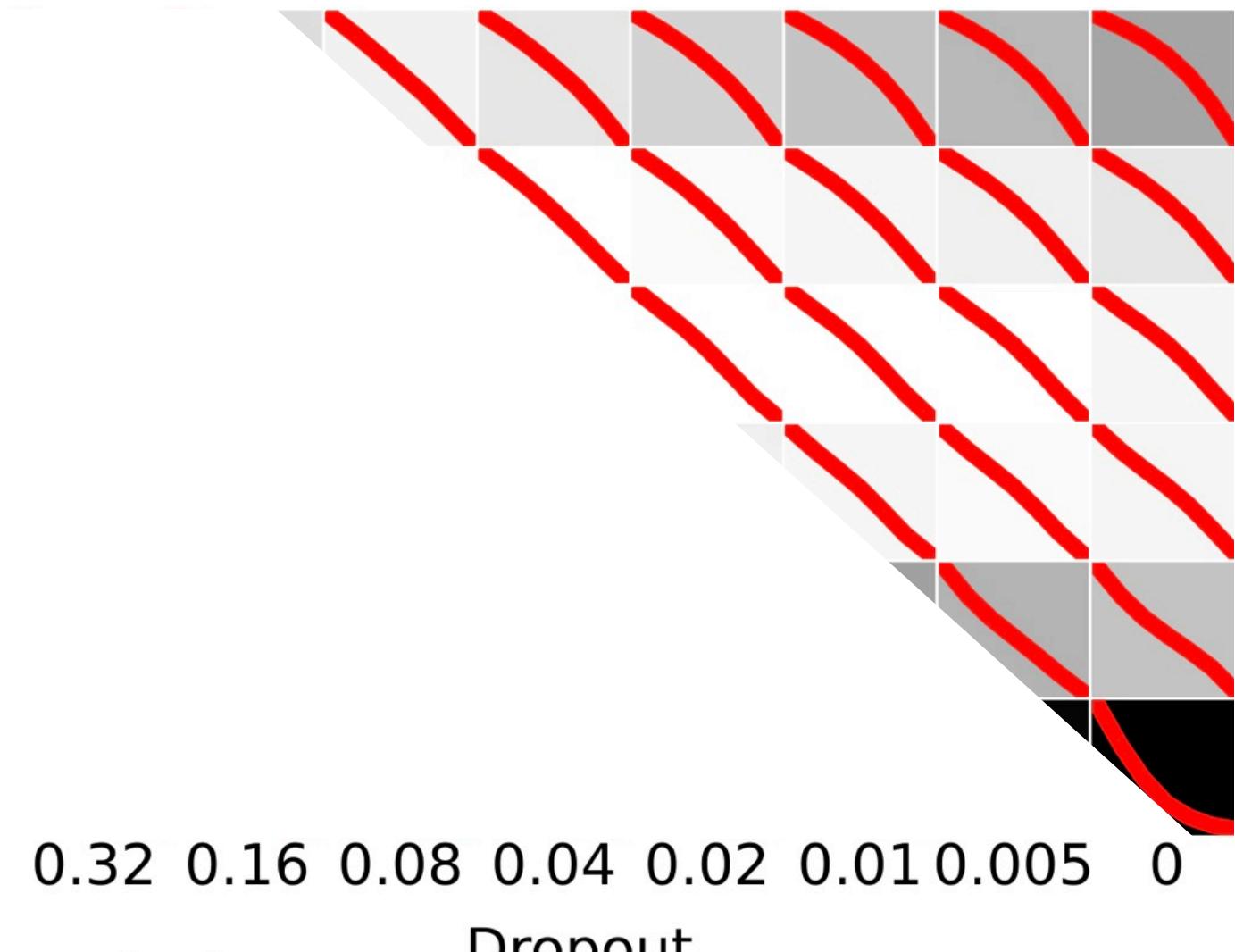
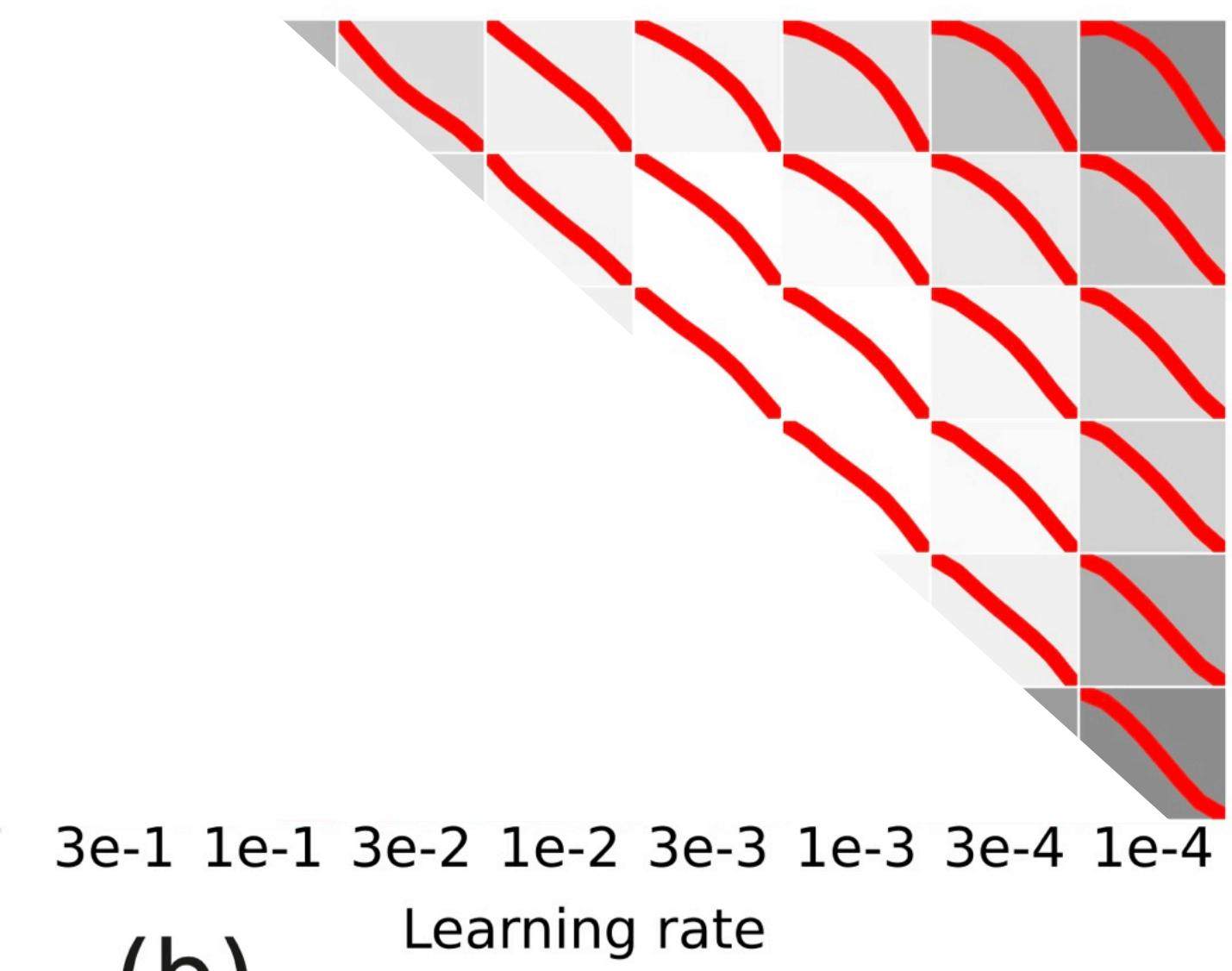
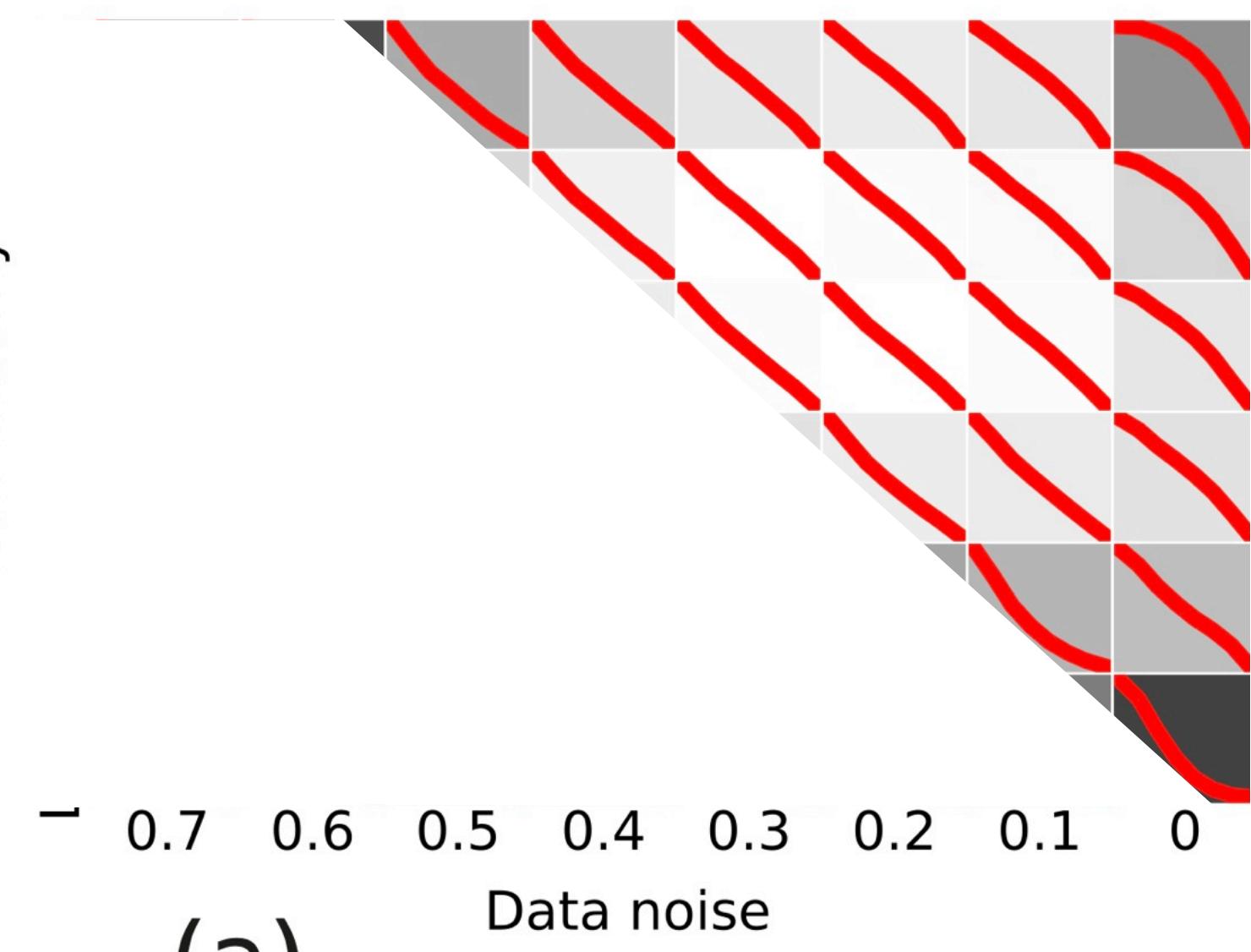
$$\ddot{x}_\ell = k(d_{\ell+1} - d_\ell) - \gamma \dot{x}_\ell - f_\ell$$

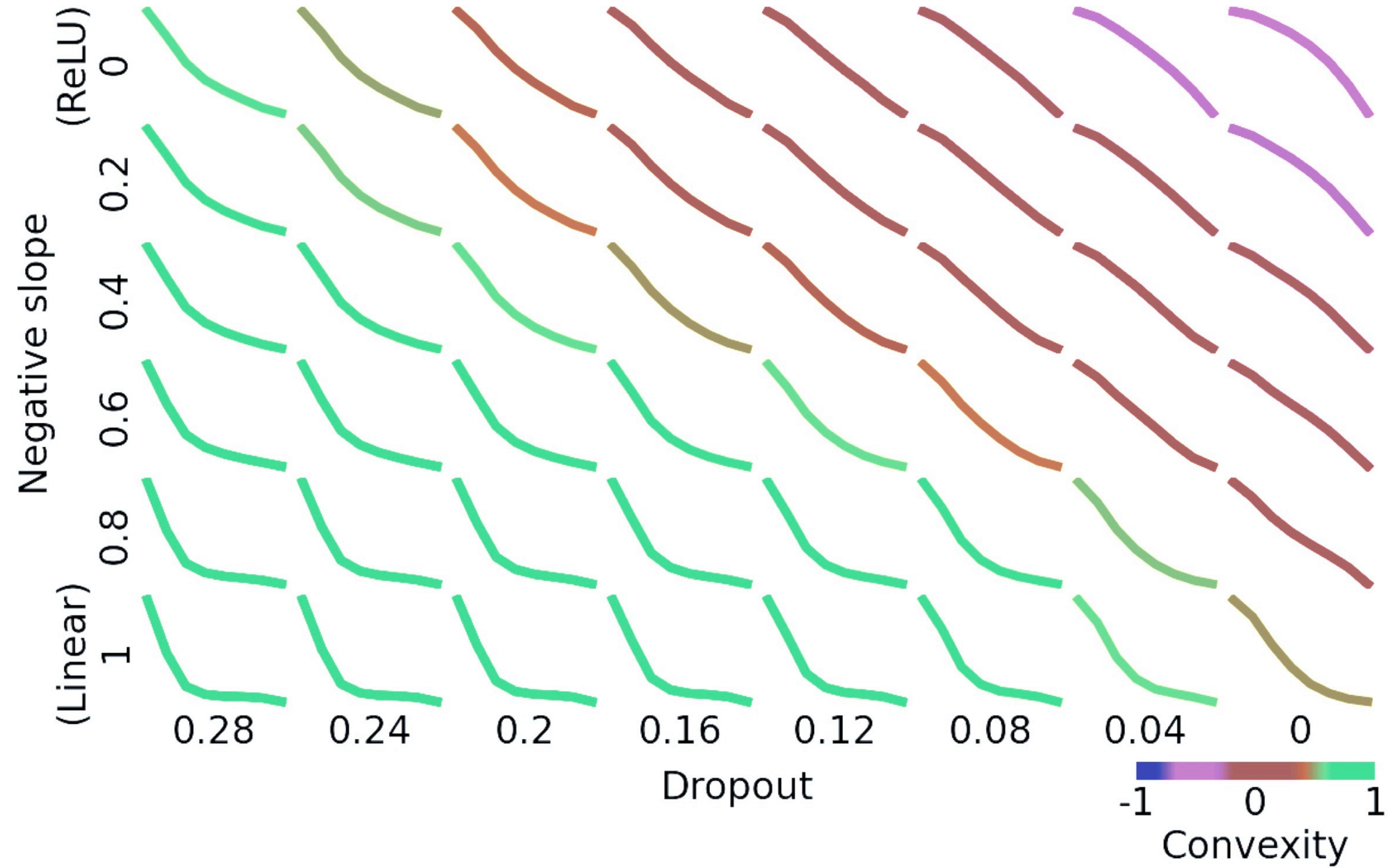


$$\ddot{x}_\ell = k(d_{\ell+1} - d_\ell) - \gamma \dot{x}_\ell - f_\ell + \sigma \xi_\ell$$

A LOW DIMENSIONAL MODEL THAT WE CAN SOLVE!

Non-linearity

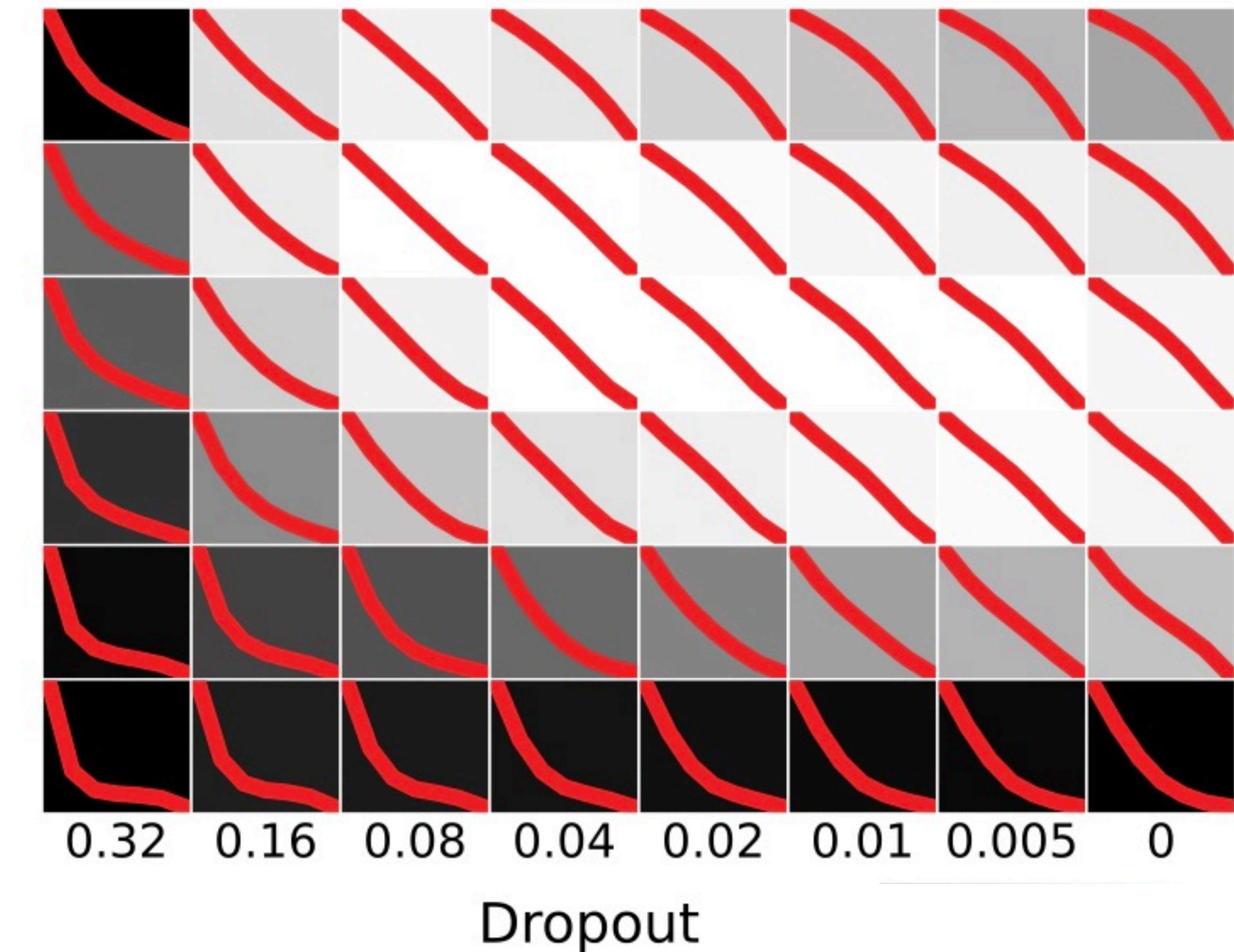
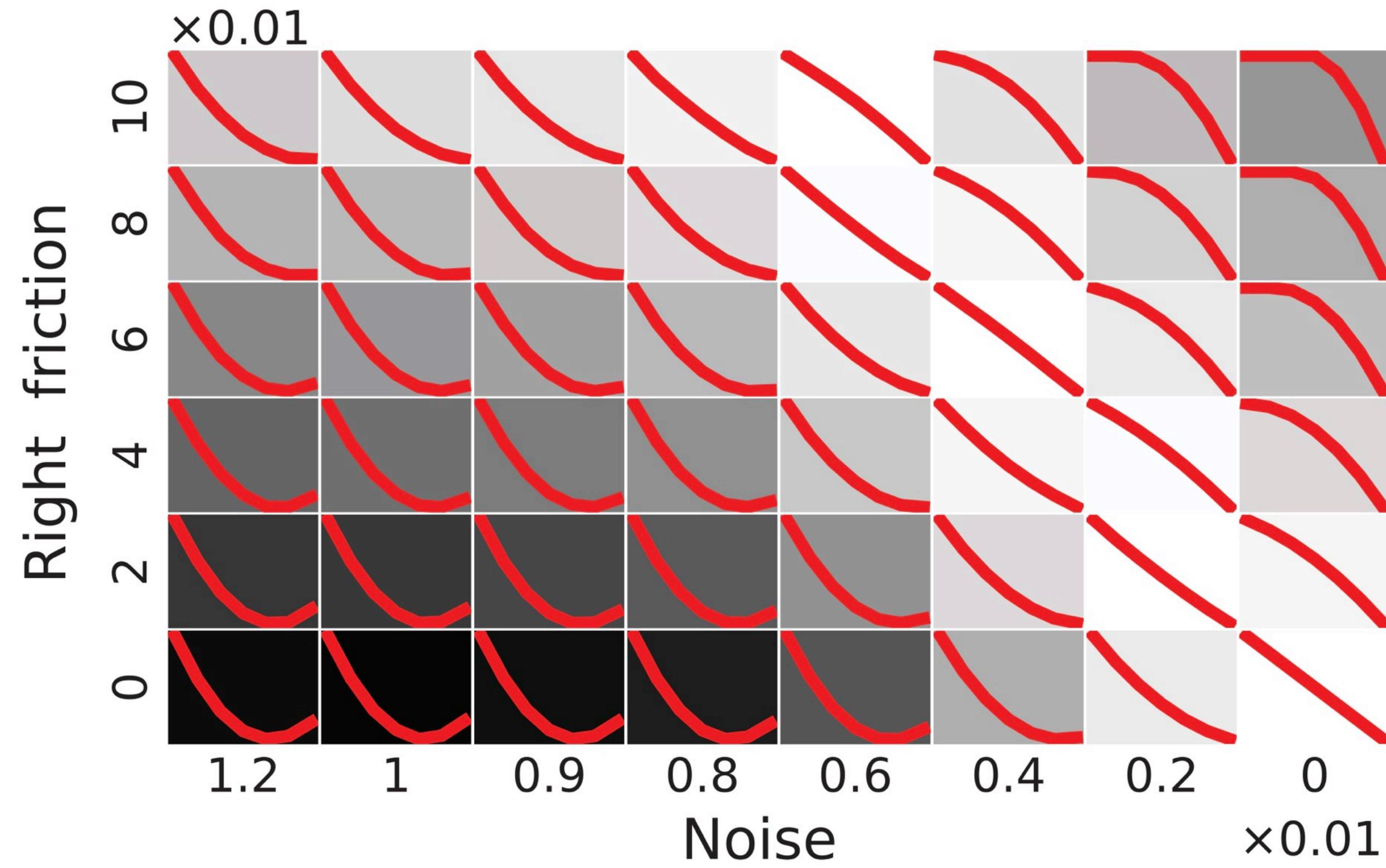




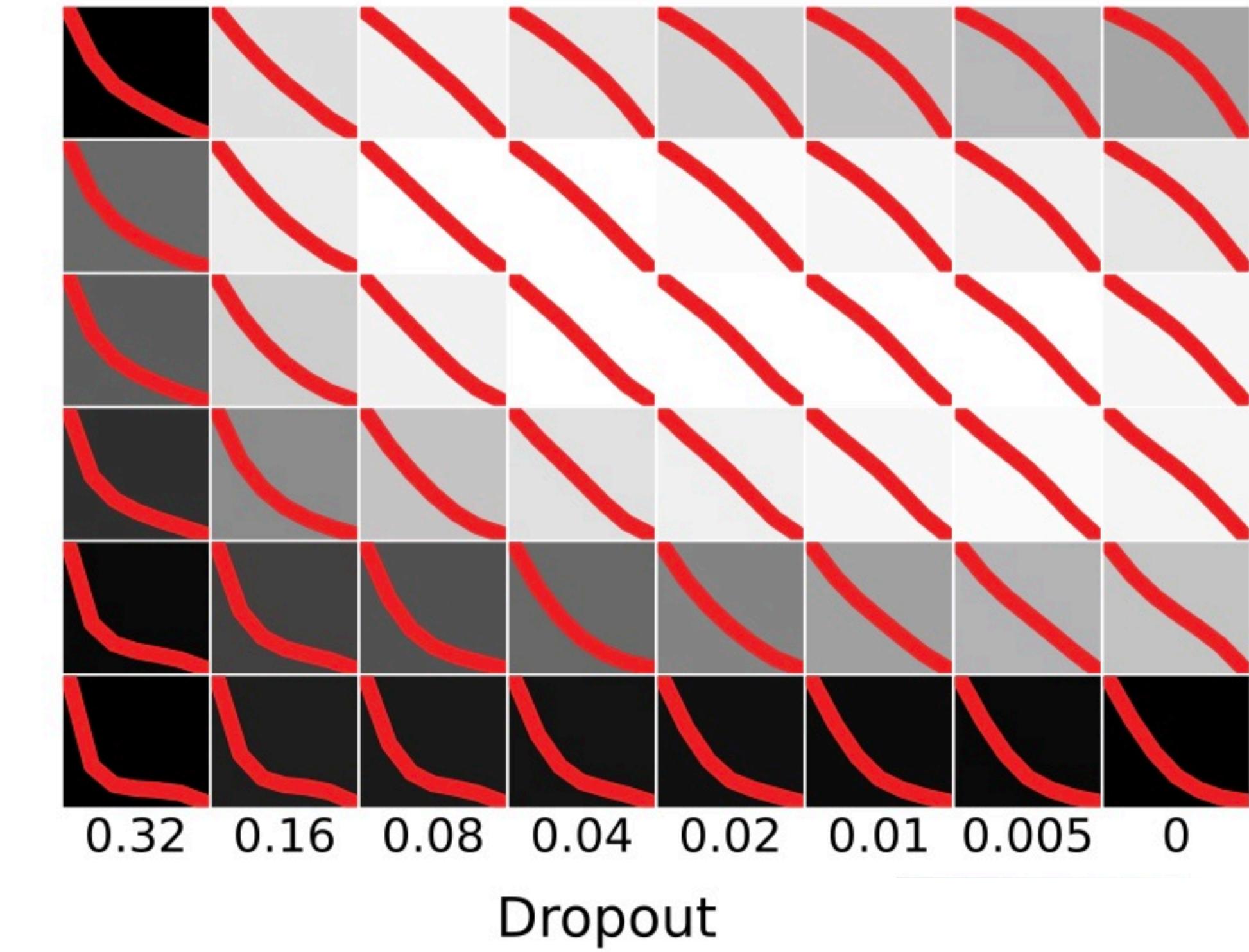
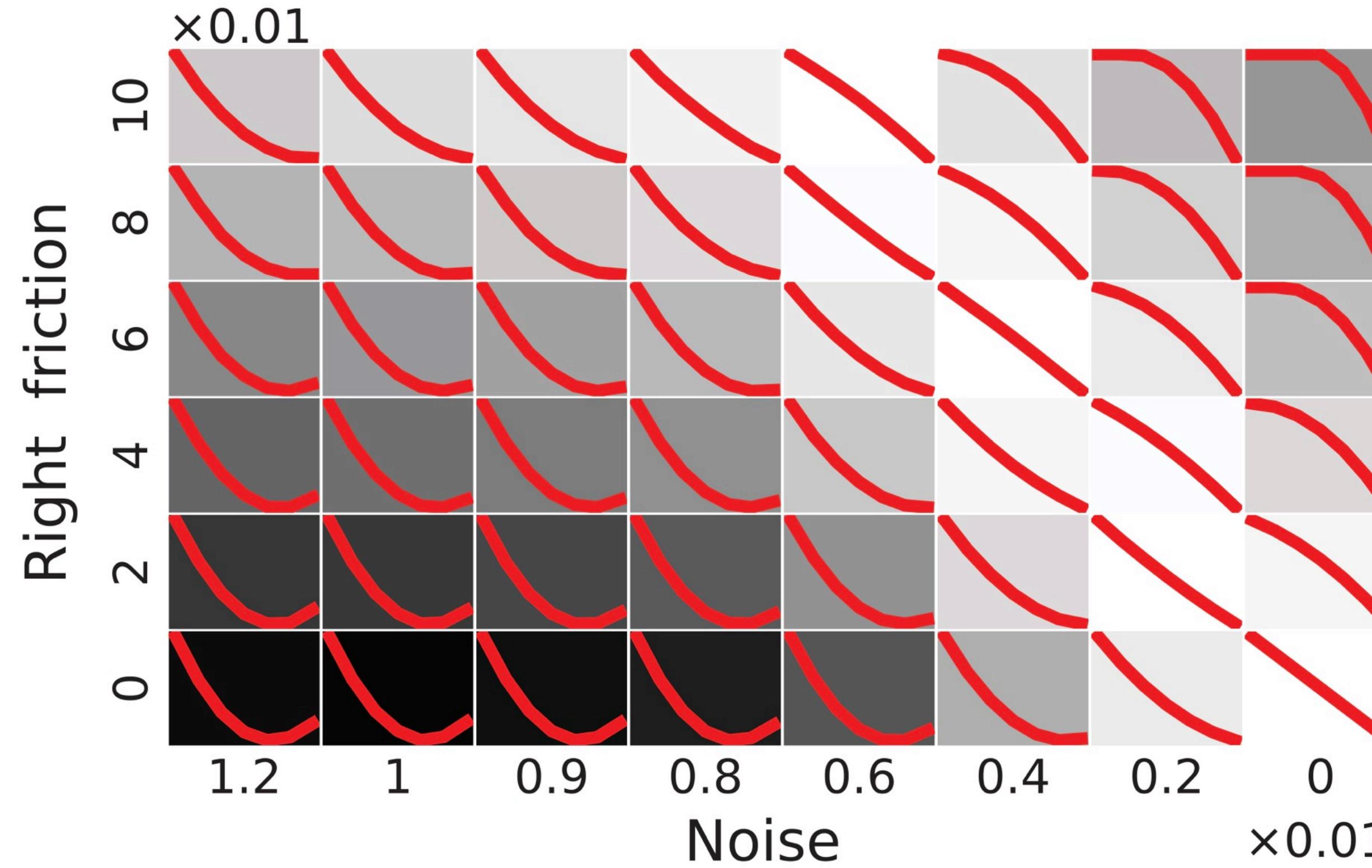
# FORWARD VS BACKWARD PASS



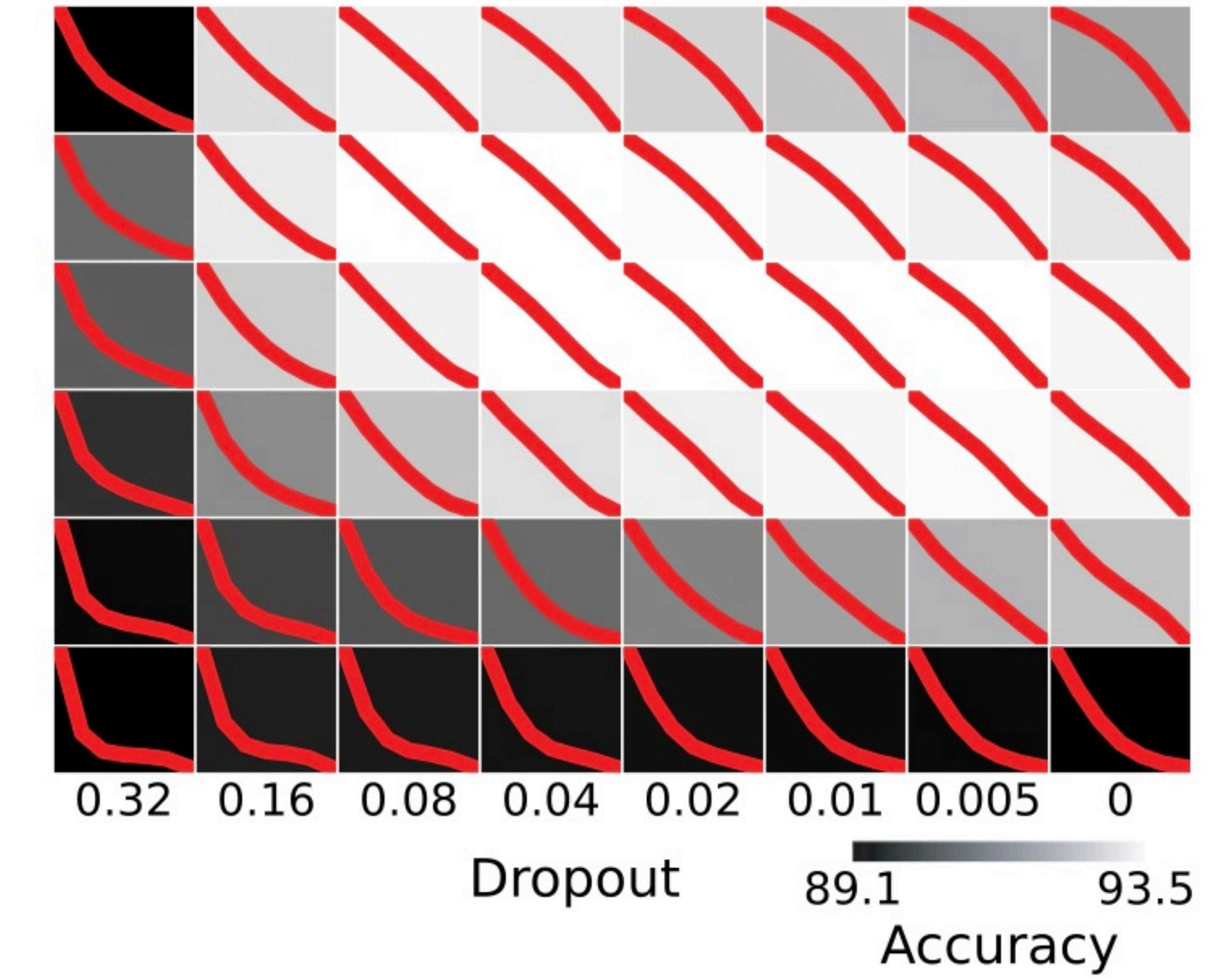
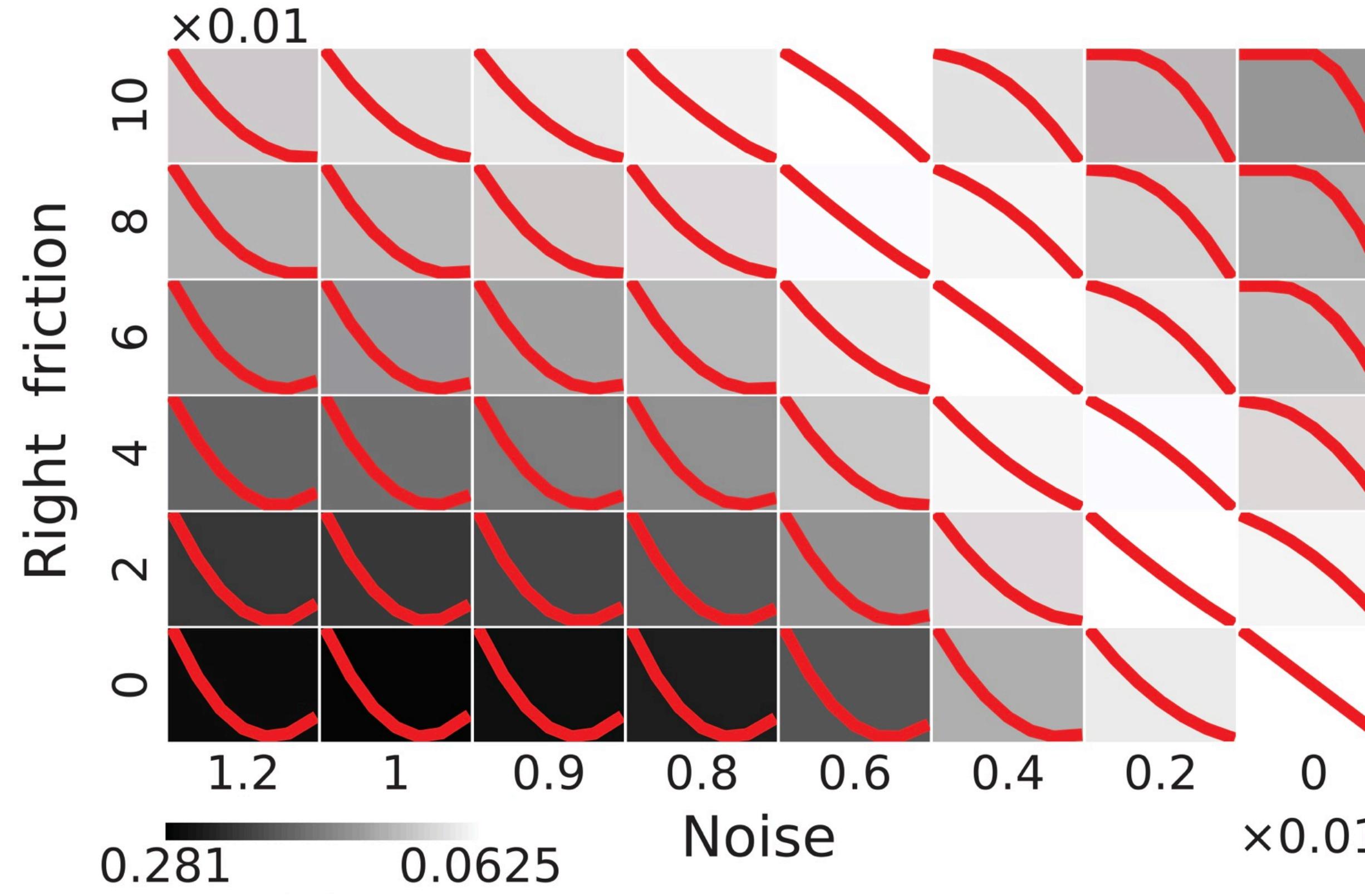
$$F_\ell = \begin{cases} k(w_{\ell+1} - w_\ell) & \text{if } \dot{h}_\ell = 0 \text{ and } -\mu_r \leq k(w_{\ell+1} - w_\ell) \leq \mu_\ell \\ -\mu_\ell & \text{if } \dot{h}_\ell > 0 \text{ or } (\dot{h}_\ell = 0 \text{ and } k(w_{\ell+1} - w_\ell) > \mu_\ell) \\ \mu_r & \text{if } \dot{h}_\ell < 0 \text{ or } (\dot{h}_\ell = 0 \text{ and } k(w_{\ell+1} - w_\ell) < -\mu_r) \end{cases}$$

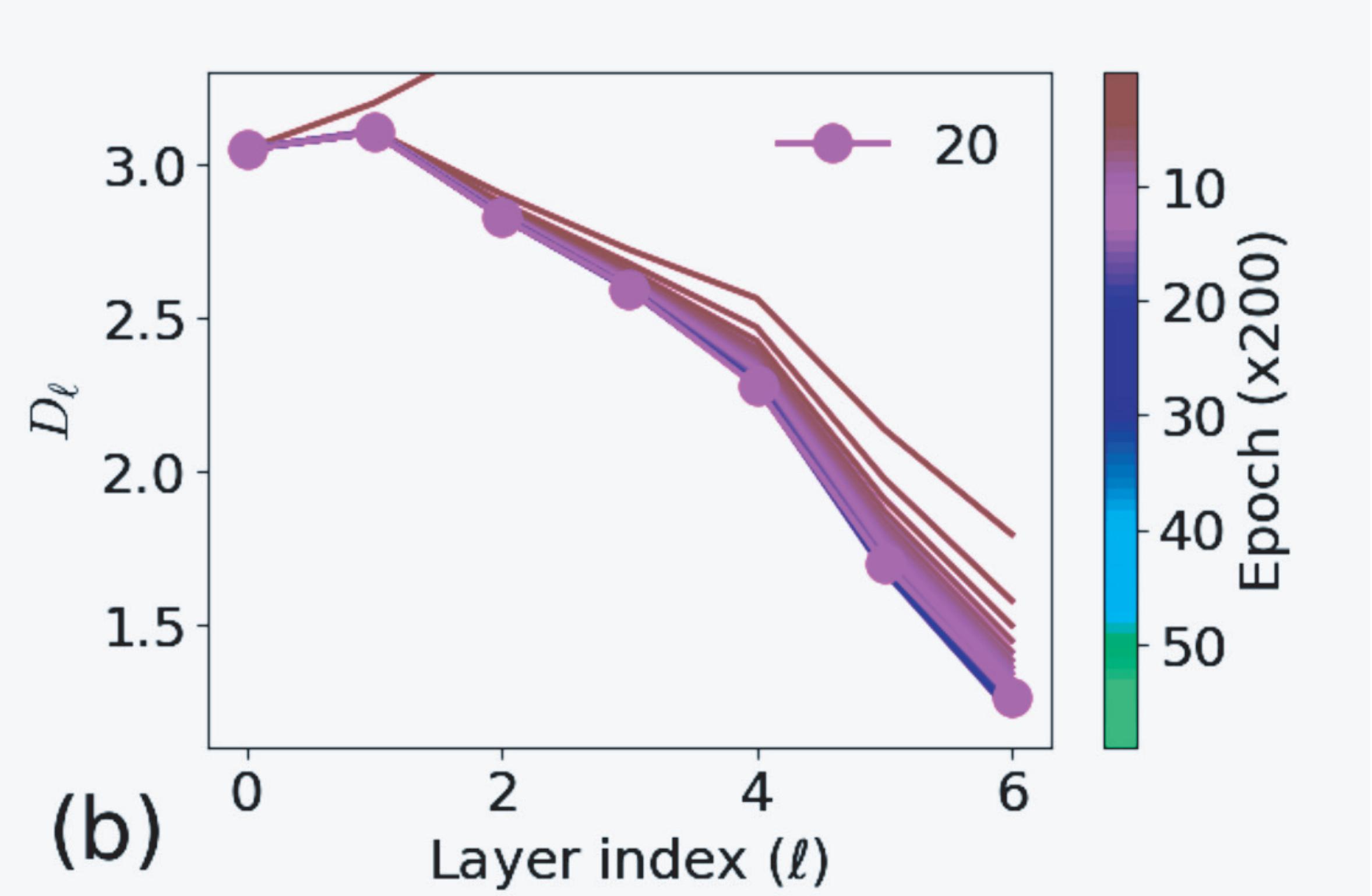
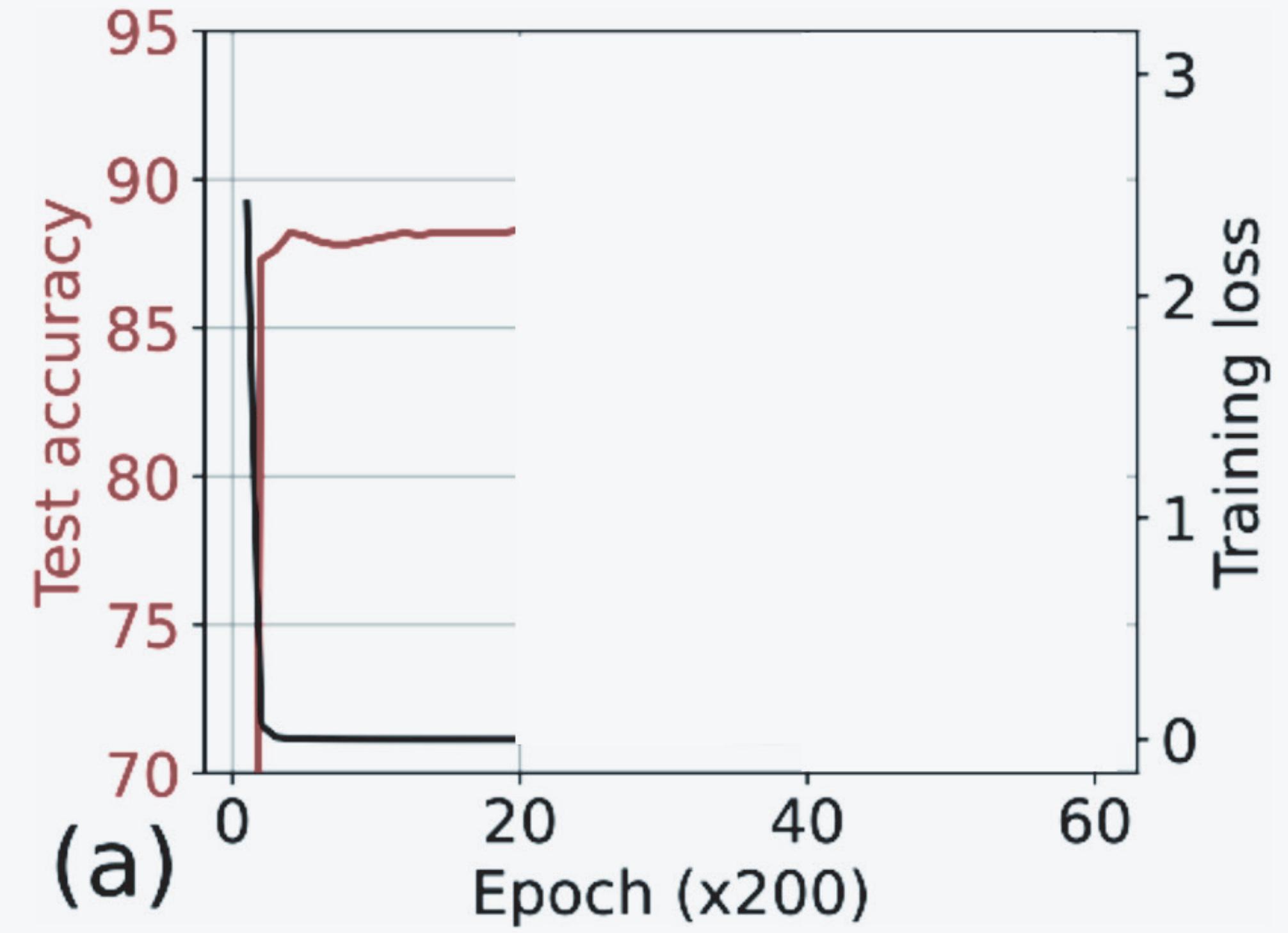


# OK SO WHAT?

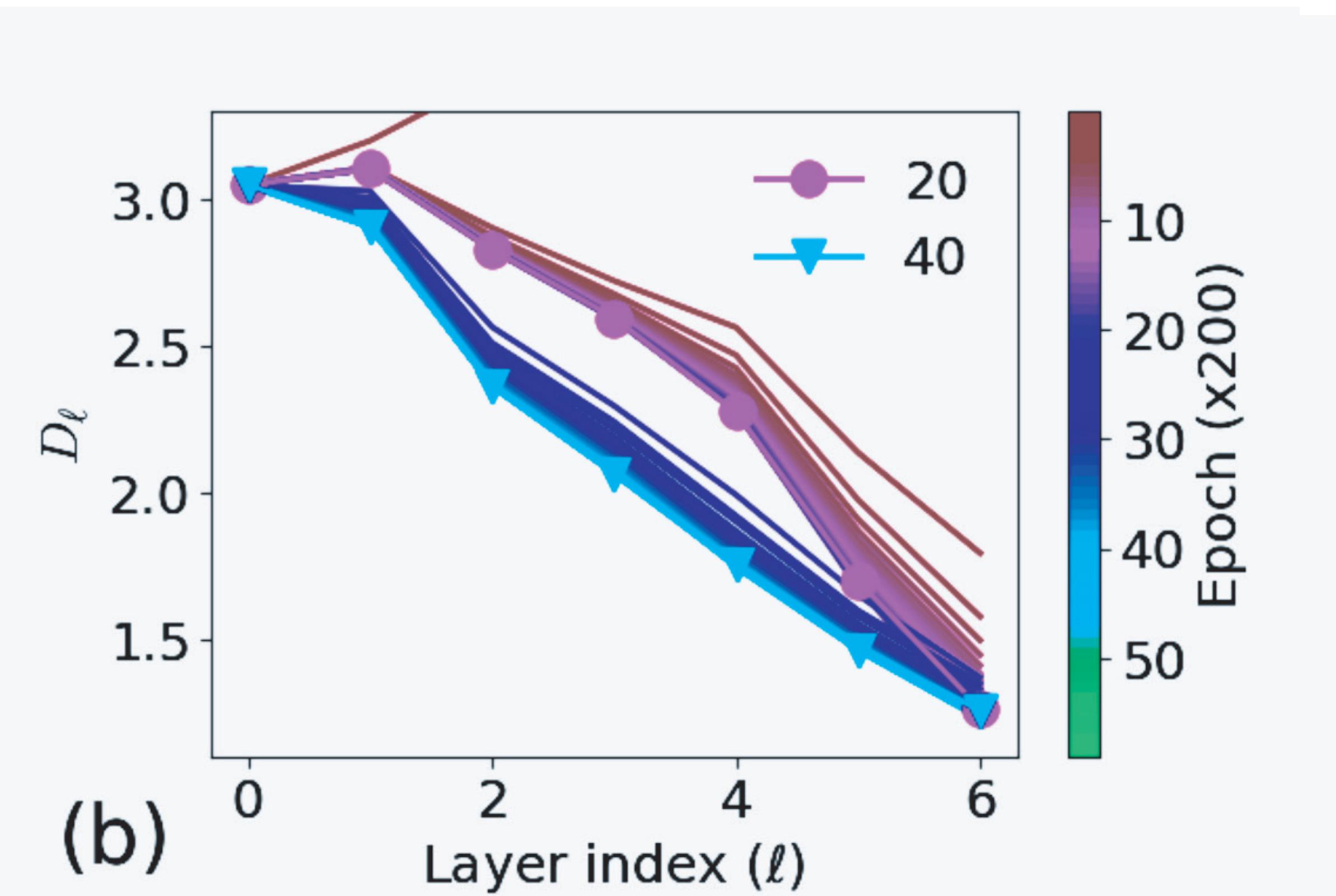
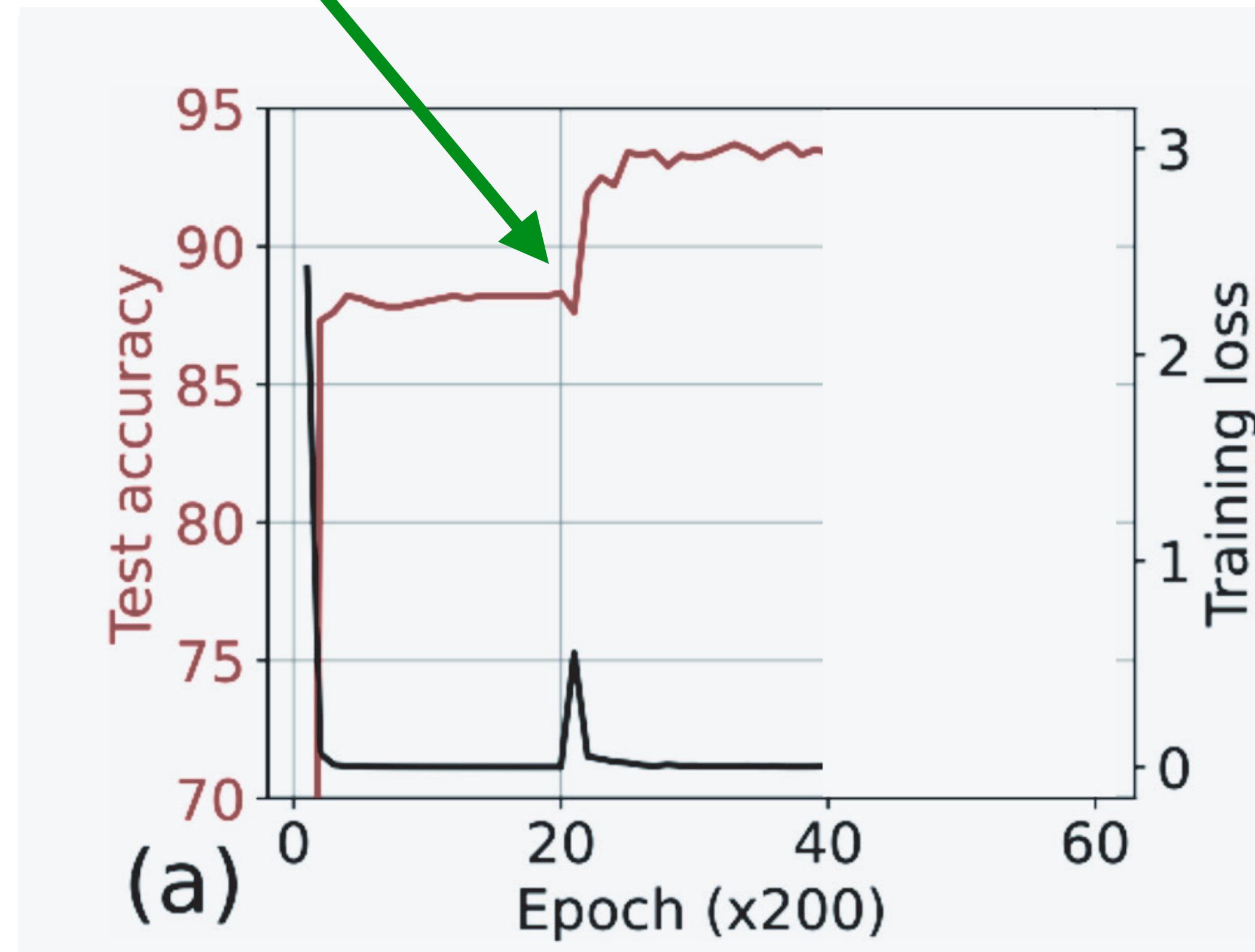


# OK SO WHAT?



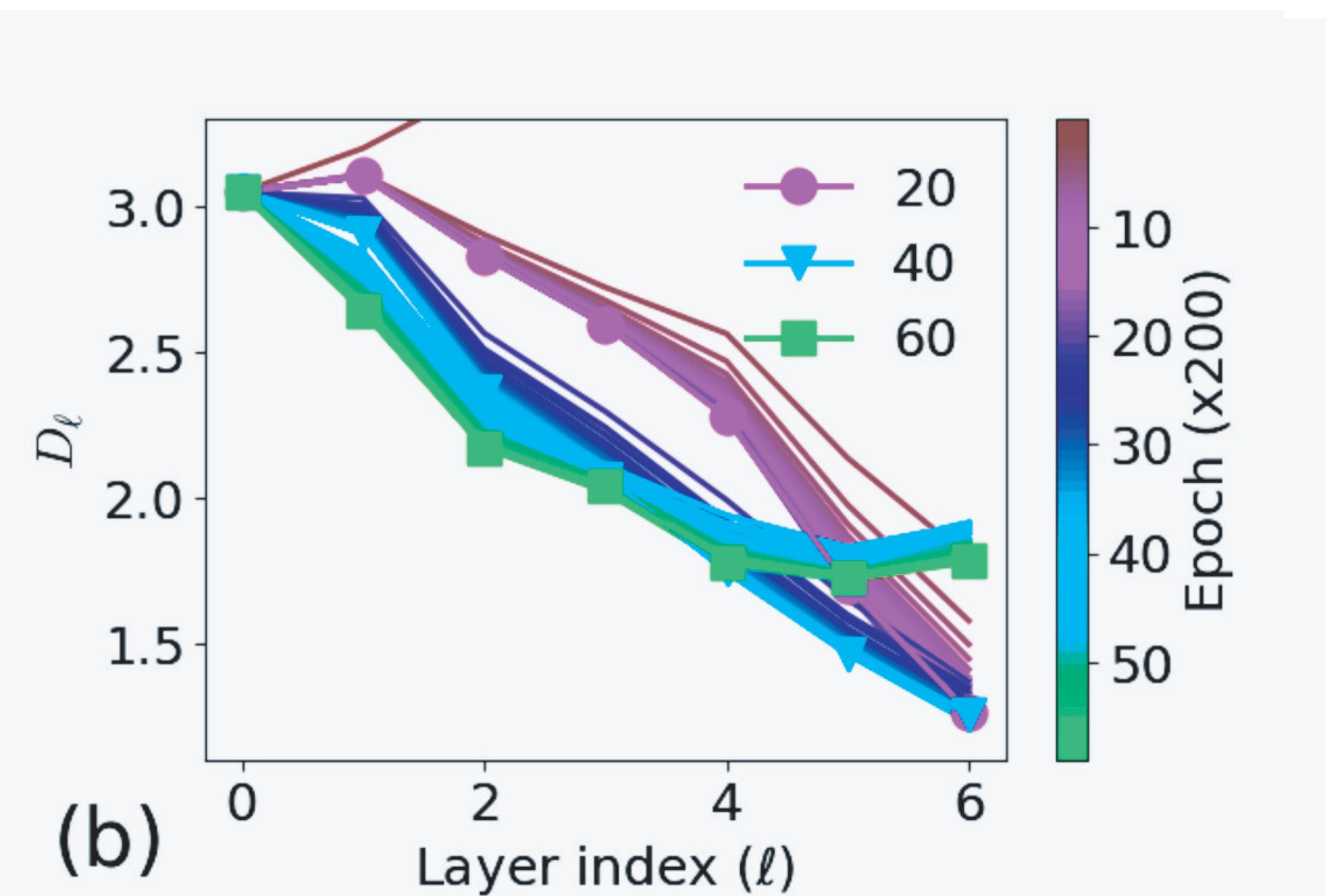
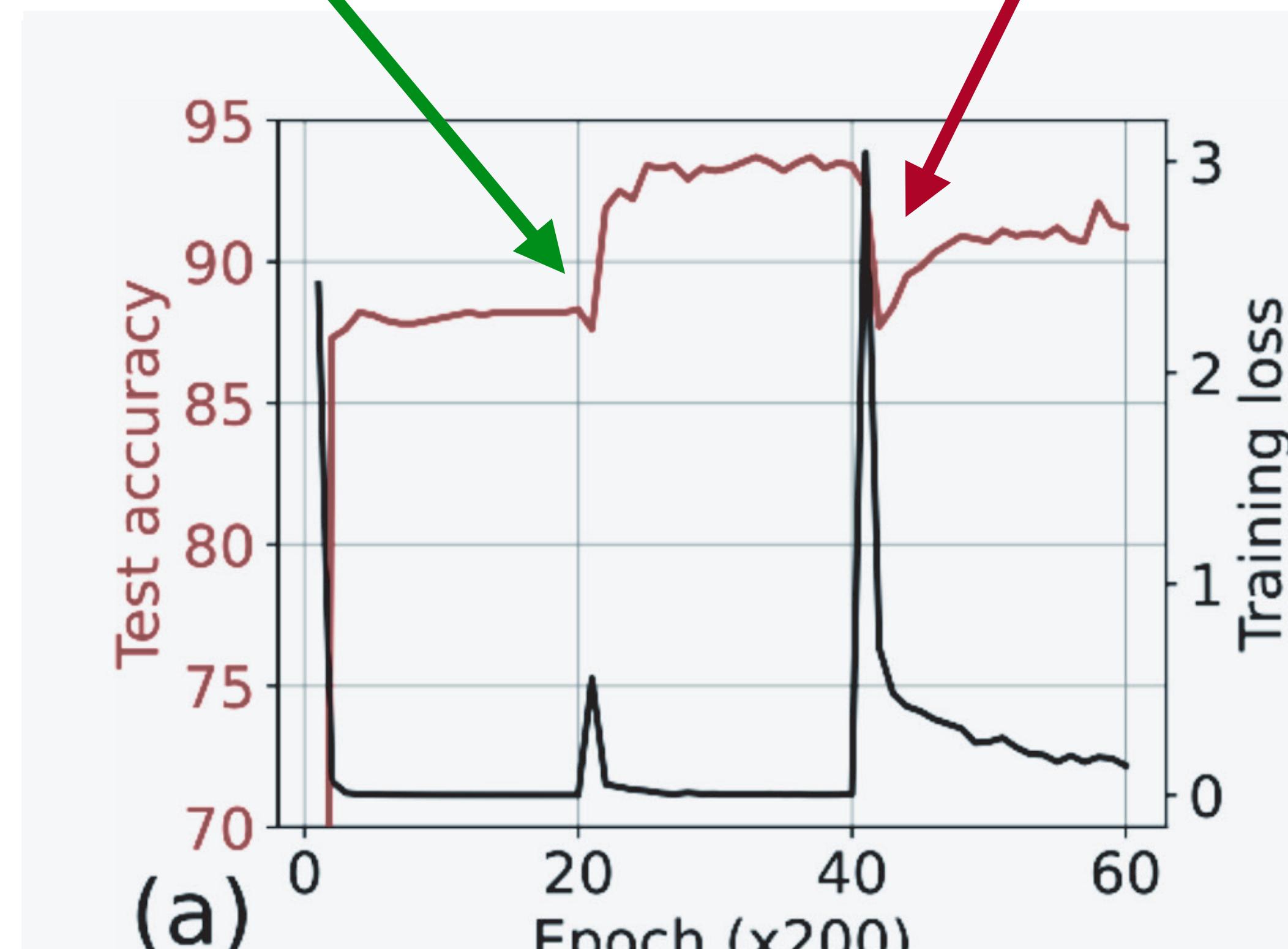


ADD 5% DROPOUT

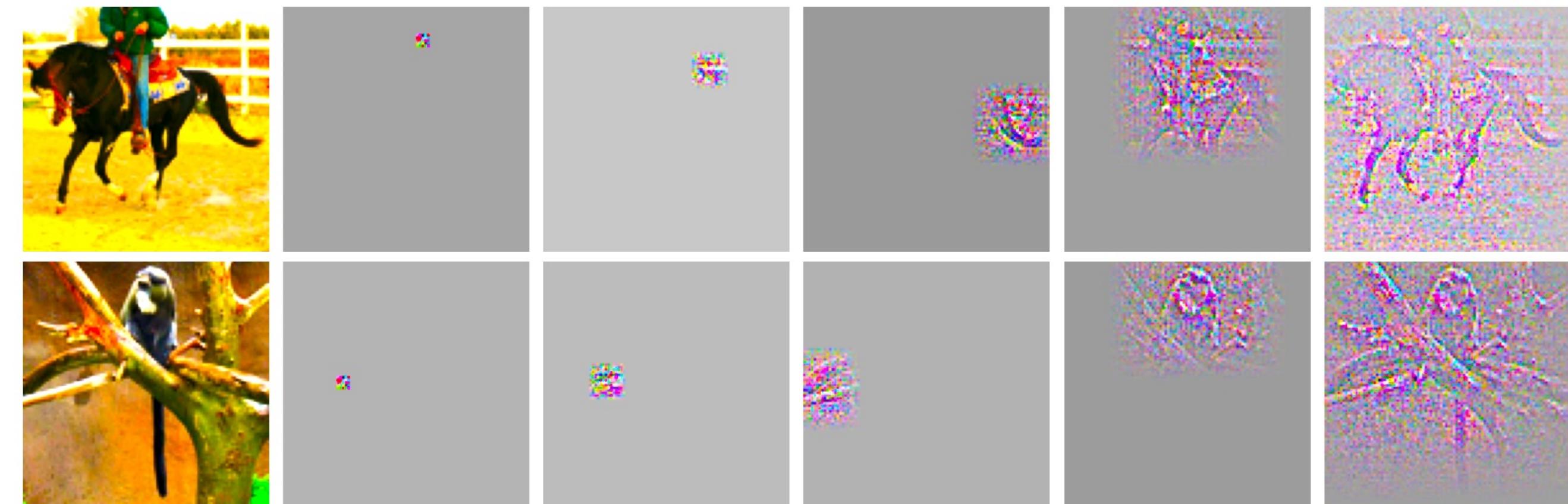


ADD 5% DROPOUT

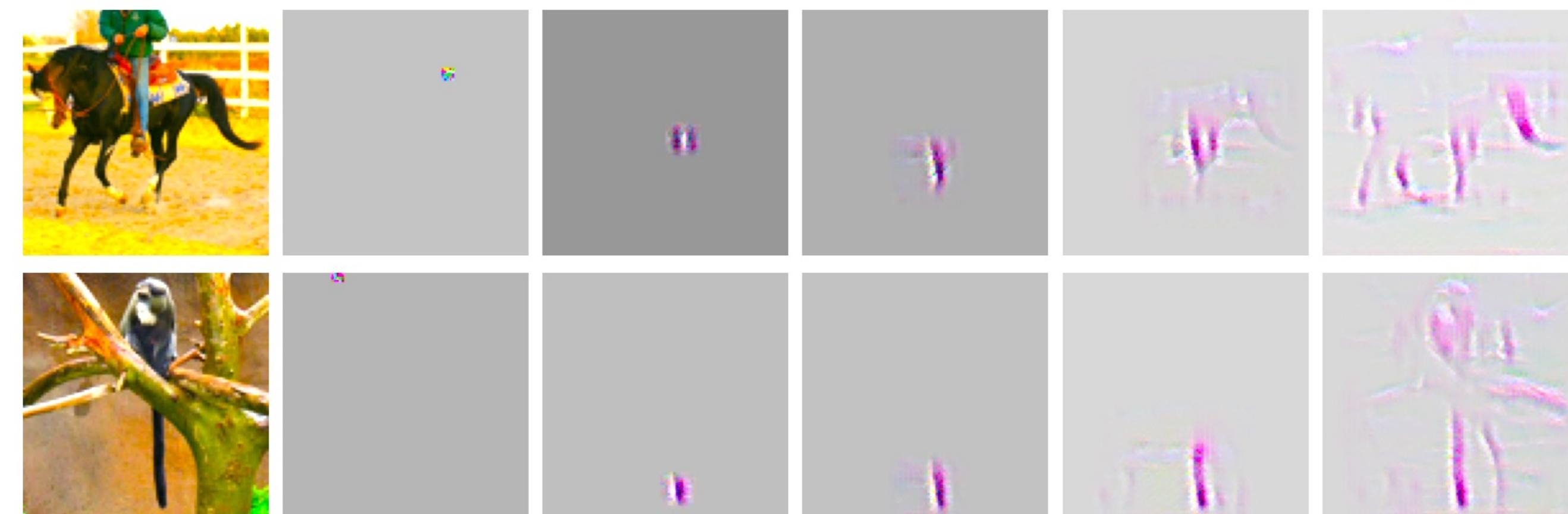
ADD 30% DROPOUT



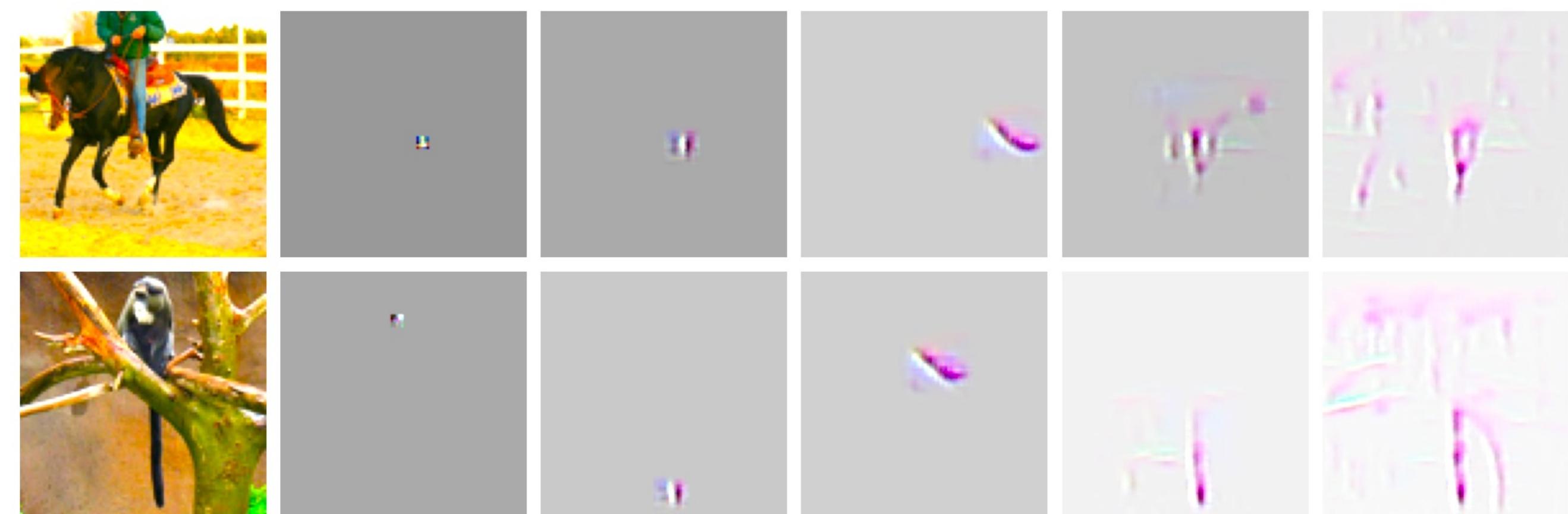
CONCAVE



LINEAR



CONVEX



ONGOING ATTEMPTS AT A FIRST-  
PRINCIPLES THEORY

# DMFT

$$\mathbf{h}_m^1 = \frac{1}{\sqrt{D}} \mathbf{W}^0 \mathbf{x}_m,$$

$$\mathbf{h}_m^{L+1} = \frac{1}{\sqrt{N}} \mathbf{w}_L \cdot \phi(\mathbf{h}_m^L),$$

FEATURE KERNEL

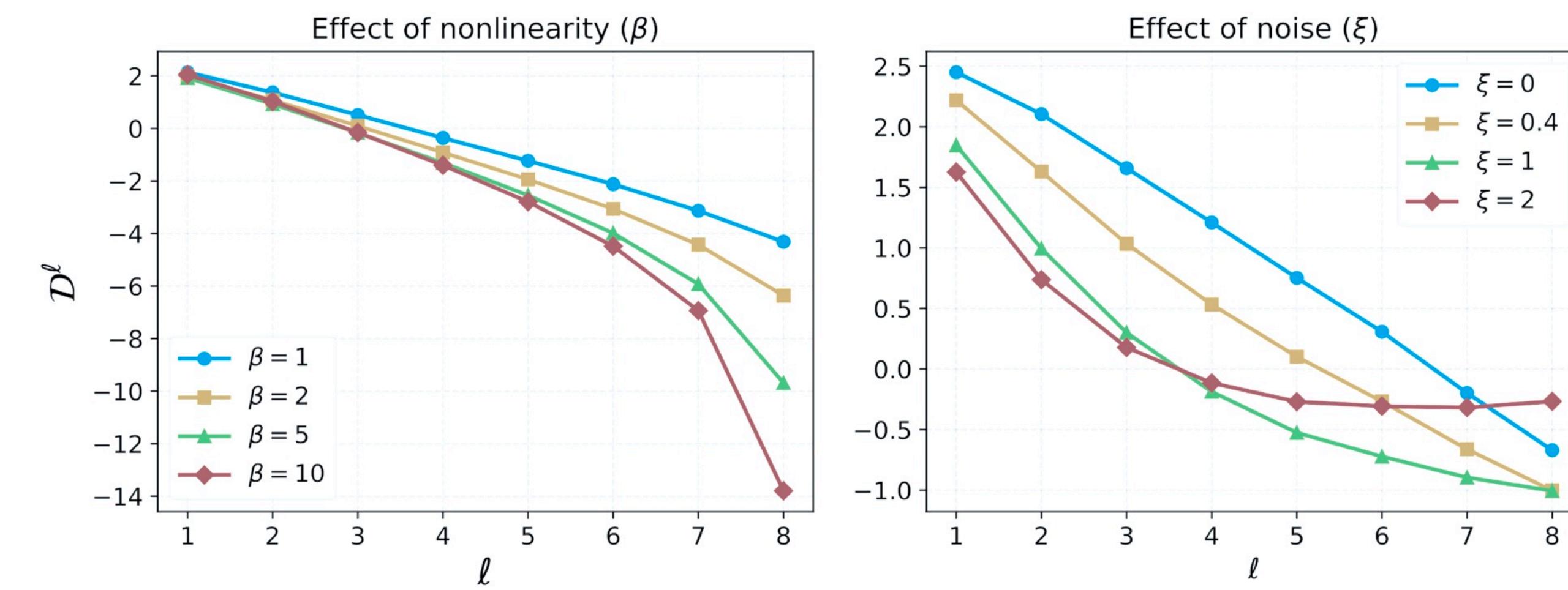
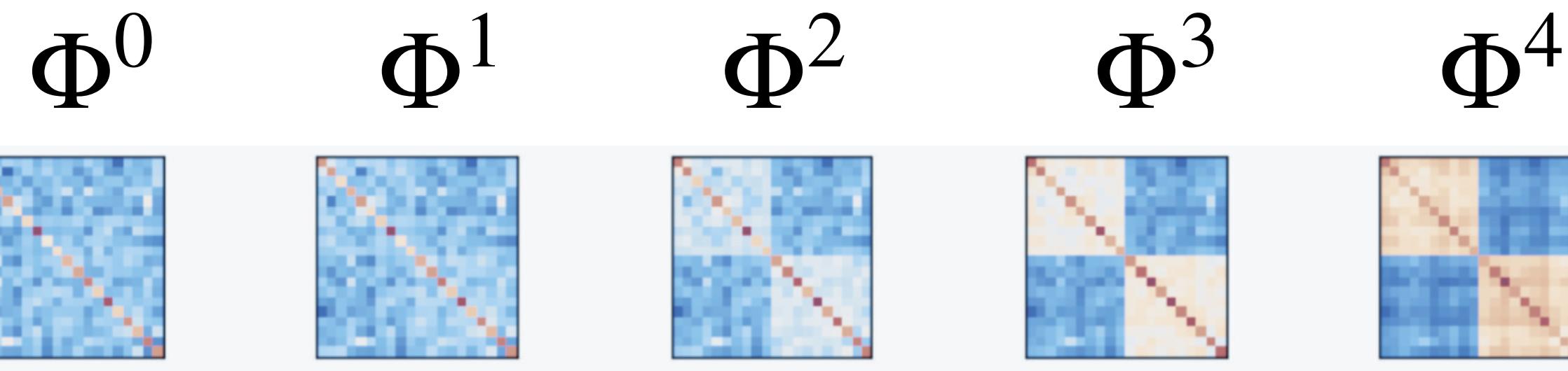
$$\Phi_{mp}^\ell(t, s) = \frac{1}{N} \left\langle \phi(\mathbf{h}_m^\ell(t)), \phi(\mathbf{h}_p^\ell(s)) \right\rangle$$

$$\mathbf{h}_m^{\ell+1} = \frac{1}{\sqrt{N}} \mathbf{W}^\ell \phi(\mathbf{h}_m^\ell)$$

$$f_m = \frac{1}{\gamma_0 \sqrt{N}} \mathbf{h}_m^{L+1}$$

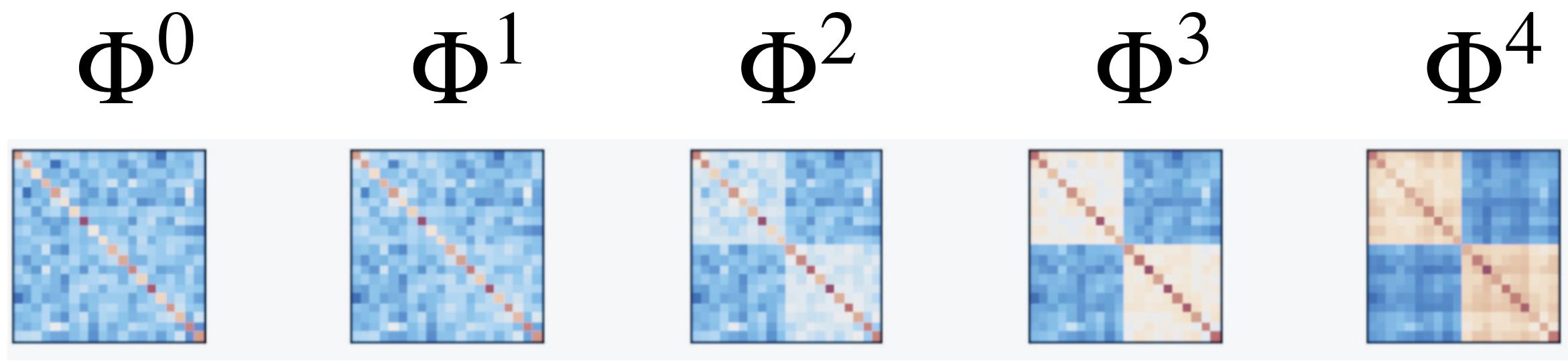
GRADIENT KERNEL

$$G_{mp}^\ell(t, s) = \frac{1}{N} \left\langle \mathbf{g}_m^\ell(t), \mathbf{g}_p^\ell(s) \right\rangle$$

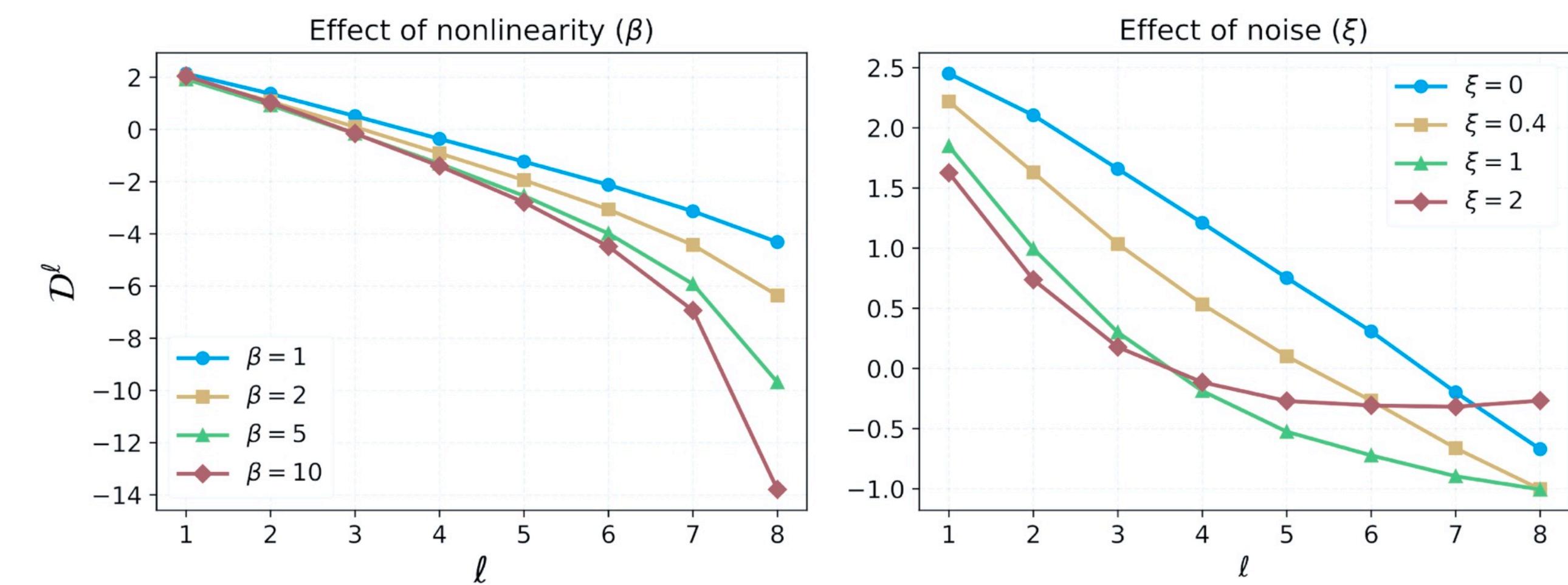


# WHITE DATA, BINARY CLASSIFICATION, LINEAR NET, VERY SMALL INIT

---



$$\mathcal{D}_\ell = \log \left( \frac{M-2}{H_y^\ell} \right) \quad H_y^\ell := \mathbf{y}^\top \Phi^\ell \mathbf{y} / M$$



$$H_y^\ell = (a(t))^\ell$$

$$\begin{aligned} \frac{da}{dt} &= \gamma_0 \Delta a^L, \\ \frac{d\Delta}{dt} &= - (L+1) \Delta a^{2L} \end{aligned}$$

# REPARAMETERIZATION

JACOT ET AL. 2022; “FEATURE LEARNING IN  $L^2$  REGULARIZED DNNs”

$$\hat{Z}_0 := X$$

$$Z_\ell = W_\ell \hat{Z}_{\ell-1}$$

$$\hat{Z}_\ell = \sigma(Z_\ell)$$

# REPARAMETERIZATION

JACOT ET AL. 2022; “FEATURE LEARNING IN  $L^2$  REGULARIZED DNNs”

$$\hat{Z}_0 := X \quad Z_\ell = W_\ell \hat{Z}_{\ell-1} \quad \hat{Z}_\ell = \sigma(Z_\ell)$$

$$W_\ell = Z_\ell \hat{Z}_{\ell-1}^+ + W_\ell^\perp \quad W_\ell^\perp \hat{Z}_{\ell-1} = 0$$

# REPARAMETERIZATION

JACOT ET AL. 2022; “FEATURE LEARNING IN  $L^2$  REGULARIZED DNNs”

$$\hat{Z}_0 := X \quad Z_\ell = W_\ell \hat{Z}_{\ell-1} \quad \hat{Z}_\ell = \sigma(Z_\ell)$$

$$W_\ell = Z_\ell \hat{Z}_{\ell-1}^+ + W_\ell^\perp \quad W_\ell^\perp \hat{Z}_{\ell-1} = 0$$

$$\|W_\ell\|_F^2 = \|Z_\ell \hat{Z}_{\ell-1}^+ + W_\ell^\perp\|_F^2 + \|W_\ell^\perp\|^2$$

# REPARAMETERIZATION

JACOT ET AL. 2022; “FEATURE LEARNING IN  $L^2$  REGULARIZED DNNs”

$$\hat{Z}_0 := X \quad Z_\ell = W_\ell \hat{Z}_{\ell-1} \quad \hat{Z}_\ell = \sigma(Z_\ell)$$

$$W_\ell = Z_\ell \hat{Z}_{\ell-1}^+ + W_\ell^\perp \quad W_\ell^\perp \hat{Z}_{\ell-1} = 0$$

$$\|W_\ell\|_F^2 = \|Z_\ell \hat{Z}_{\ell-1}^+ + W_\ell^\perp\|_F^2 + \|W_\ell^\perp\|^2 \quad \Rightarrow \quad W_\ell^\perp = 0$$

## ORIGINAL LOSS

$$\mathcal{L}(W_1, \dots, W_L) = C(Z_L) + \frac{\lambda}{2} \sum_{\ell=1}^L \|W_\ell\|_F^2$$

## ORIGINAL LOSS

$$C(Z_L) = \|Z_L - Y\|_F^2$$

$$\mathcal{L}(W_1, \dots, W_L) = C(Z_L) + \frac{\lambda}{2} \sum_{\ell=1}^L \|W_\ell\|_F^2$$

## REPARAM 1

$$\mathcal{L}(Z_1, \dots, Z_L) = C(Z_L) + \frac{\lambda}{2} \sum_{\ell=1}^L \|Z_\ell \hat{Z}_{\ell-1}^+\|_F^2$$

## ORIGINAL LOSS

$$C(\mathbf{Z}_L) = \|\mathbf{Z}_L - \mathbf{Y}\|_F^2$$

$$\mathcal{L}(W_1, \dots, W_L) = C(\mathbf{Z}_L) + \frac{\lambda}{2} \sum_{\ell=1}^L \|W_\ell\|_F^2$$

## REPARAM 1

$$\mathcal{L}(\mathbf{Z}_1, \dots, \mathbf{Z}_L) = C(\mathbf{Z}_L) + \frac{\lambda}{2} \sum_{\ell=1}^L \|\mathbf{Z}_\ell \hat{\mathbf{Z}}_{\ell-1}^+\|_F^2$$

$$K_\ell = \mathbf{Z}_\ell^T \mathbf{Z}_\ell$$

$$\hat{K}_\ell = \hat{\mathbf{Z}}_\ell^T \hat{\mathbf{Z}}_\ell$$

$$C(\mathbf{Z}_L) = \|\mathbf{Z}_L - \mathbf{Y}\|_F^2$$

ORIGINAL LOSS

$$\mathcal{L}(W_1, \dots, W_L) = C(\mathbf{Z}_L) + \frac{\lambda}{2} \sum_{\ell=1}^L \|W_\ell\|_F^2$$

REPARAM 1

$$\mathcal{L}(\mathbf{Z}_1, \dots, \mathbf{Z}_L) = C(\mathbf{Z}_L) + \frac{\lambda}{2} \sum_{\ell=1}^L \|\mathbf{Z}_\ell \hat{\mathbf{Z}}_{\ell-1}^+\|_F^2$$

$$K_\ell = \mathbf{Z}_\ell^T \mathbf{Z}_\ell$$
$$\hat{K}_\ell = \hat{\mathbf{Z}}_\ell^T \hat{\mathbf{Z}}_\ell$$

REPARAM 2

$$\mathcal{L}(\mathcal{K}, \mathbf{Z}_L) = C(\mathbf{Z}_L) + \frac{\lambda}{2} \sum_{\ell=1}^L \|K_\ell \hat{K}_{\ell-1}^+\|_F^2$$

# MACROSCOPIC DYNAMICS

WORK IN PROGRESS! POSSIBLE MISTAKES!

$$\varphi_\ell = \text{Tr} [K_\ell]$$

# MACROSCOPIC DYNAMICS

WORK IN PROGRESS! POSSIBLE MISTAKES!

$$\varphi_\ell = \text{Tr} [K_\ell]$$

$$\hat{\varphi}_\ell^+ = \text{Tr} [ \widehat{K}_\ell^+ ]$$

# MACROSCOPIC DYNAMICS

WORK IN PROGRESS! POSSIBLE MISTAKES!

$$\varphi_\ell = \text{Tr} [K_\ell]$$

$$\hat{\varphi}_\ell^+ = \text{Tr} \left[ \widehat{K}_\ell^+ \right]$$

$$\hat{\psi}_\ell = \text{Tr} [K_\ell P_{\ell-1}]$$

# MACROSCOPIC DYNAMICS

WORK IN PROGRESS! POSSIBLE MISTAKES!

$$\hat{Z}_{\ell-1}^+ \hat{Z}_{\ell-1}$$

$$\varphi_\ell = \text{Tr} [K_\ell]$$

$$\hat{\varphi}_\ell^+ = \text{Tr} [ \widehat{K}_\ell^+ ]$$

$$\hat{\psi}_\ell = \text{Tr} [ K_\ell P_{\ell-1} ]$$

# MACROSCOPIC DYNAMICS

WORK IN PROGRESS! POSSIBLE MISTAKES!

$$\hat{Z}_{\ell-1}^+ \hat{Z}_{\ell-1}$$

$$\varphi_\ell = \text{Tr} [K_\ell]$$

$$\hat{\varphi}_\ell^+ = \text{Tr} [ \widehat{K}_\ell^+ ]$$

$$\hat{\psi}_\ell = \text{Tr} [ K_\ell P_{\ell-1} ]$$

$$\omega = \frac{1}{N} \langle Z_L, Y \rangle$$

# MACROSCOPIC DYNAMICS

WORK IN PROGRESS! POSSIBLE MISTAKES!

$$\hat{Z}_{\ell-1}^+ \hat{Z}_{\ell-1}$$

$$\varphi_\ell = \text{Tr} [K_\ell]$$

$$\hat{\varphi}_\ell^+ = \text{Tr} [ \widehat{K}_\ell^+ ]$$

$$\hat{\psi}_\ell = \text{Tr} [ K_\ell P_{\ell-1} ]$$

$$\omega = \frac{1}{N} \langle Z_L, Y \rangle$$

---

$$\frac{d\varphi_\ell}{dt} = -\psi_\ell + \psi_{\ell+1}$$

$$\frac{d\psi_\ell}{dt} = -a_\ell \psi_\ell \hat{\varphi}_{\ell-1}^+ + b_\ell \psi_{\ell+1} \hat{\varphi}_{\ell-1}^+ + c_\ell \psi_{\ell+1} \hat{\varphi}_\ell^+ - d_\ell \frac{\psi_\ell^2 \hat{\varphi}_\ell^+}{\varphi_{\ell-1}}$$

$$\lambda = \frac{1}{2}$$

$$\frac{d\hat{\varphi}_\ell^+}{dt} = \frac{1}{2} e_\ell \hat{\varphi}_\ell^+ \hat{\varphi}_{\ell-1}^+ - \frac{1}{2} f_\ell \frac{\hat{\varphi}_\ell^+ \psi_{\ell+1}}{\varphi_\ell}$$

# LINEAR NETWORK

$$\varphi_\ell = \psi_\ell = \hat{\varphi}_\ell^+$$

---

$$\dot{\varphi}_\ell = \frac{\varphi_{\ell+1}}{\varphi_\ell} - \frac{\varphi_\ell}{\varphi_{\ell-1}}$$

$$\dot{\varphi}_L = 2\omega - 2\varphi_L - \frac{\varphi_L}{\varphi_{L-1}}$$

$$\dot{\omega} = 1 - \omega - \frac{\omega}{\varphi_{L-1}}$$

# LINEAR NETWORK

$$\varphi_\ell = \psi_\ell = \hat{\varphi}_\ell^+$$

$$\dot{\varphi}_\ell = \frac{\varphi_{\ell+1}}{\varphi_\ell} - \frac{\varphi_\ell}{\varphi_{\ell-1}}$$

$$\dot{\varphi}_L = 2\omega - 2\varphi_L - \frac{\varphi_L}{\varphi_{L-1}} \implies \log \varphi_{\ell+1} - \log \varphi_\ell = \log \varphi_\ell - \log \varphi_{\ell-1}$$

$$\dot{\omega} = 1 - \omega - \frac{\omega}{\varphi_{L-1}}$$

A STATIONARY POINT

# LINEAR NETWORK

$$\varphi_\ell = \psi_\ell = \hat{\varphi}_\ell^+$$

$$\dot{\varphi}_\ell = \frac{\varphi_{\ell+1}}{\varphi_\ell} - \frac{\varphi_\ell}{\varphi_{\ell-1}}$$

A STATIONARY POINT

$$\dot{\varphi}_L = 2\omega - 2\varphi_L - \frac{\varphi_L}{\varphi_{L-1}} \implies \log \varphi_{\ell+1} - \log \varphi_\ell = \log \varphi_\ell - \log \varphi_{\ell-1}$$

$$\dot{\omega} = 1 - \omega - \frac{\omega}{\varphi_{L-1}}$$

# NONLINEAR NETWORK (NO NOISE)

$$\log \varphi_{\ell+1} - \log \hat{\varphi}_\ell^+ = \log \varphi_\ell - \log \hat{\varphi}_{\ell-1}^+$$

# ENDNOTES

# SPRINGS & BLOCKS

Bulletin of the Seismological Society of America. Vol. 57, No. 3, pp. 341–371. June, 1967

## MODEL AND THEORETICAL SEISMICITY

By R. BURRIDGE AND L. KNOPOFF

### ABSTRACT

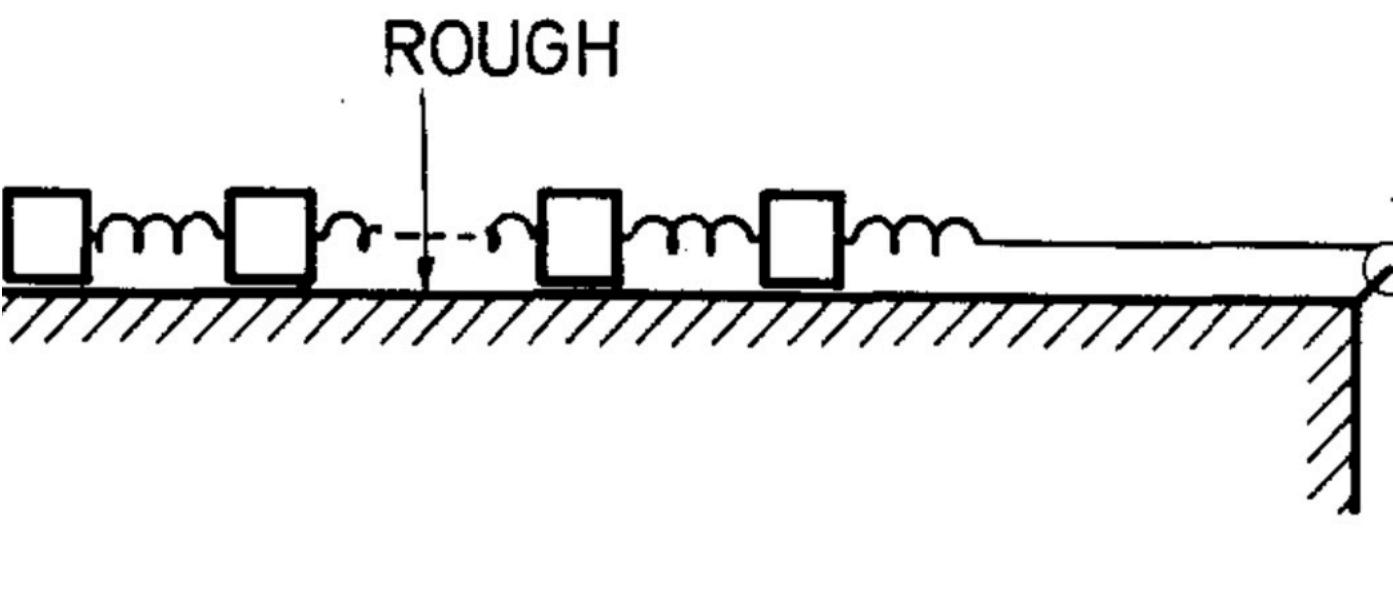


Fig. 3. Schematic diagram of the laboratory model.

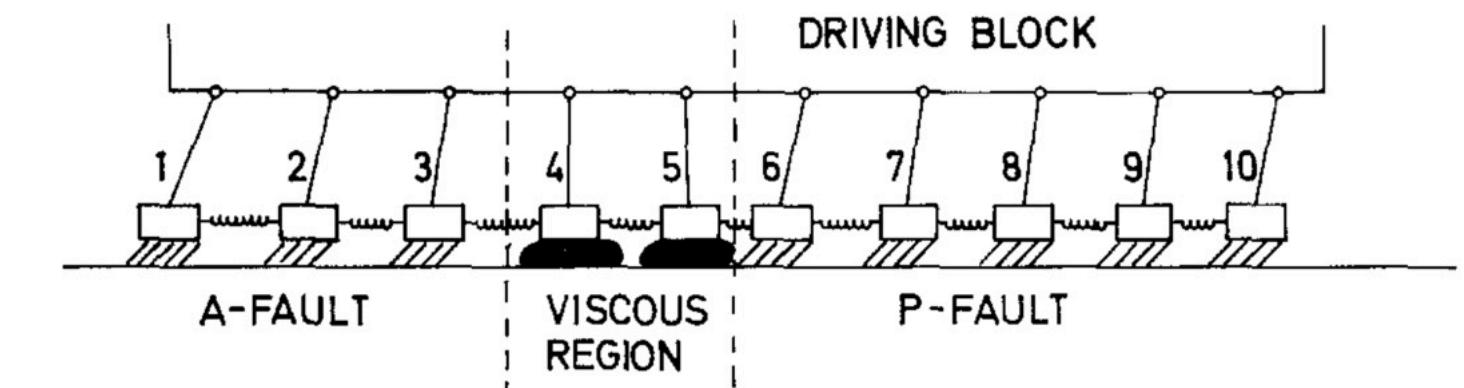
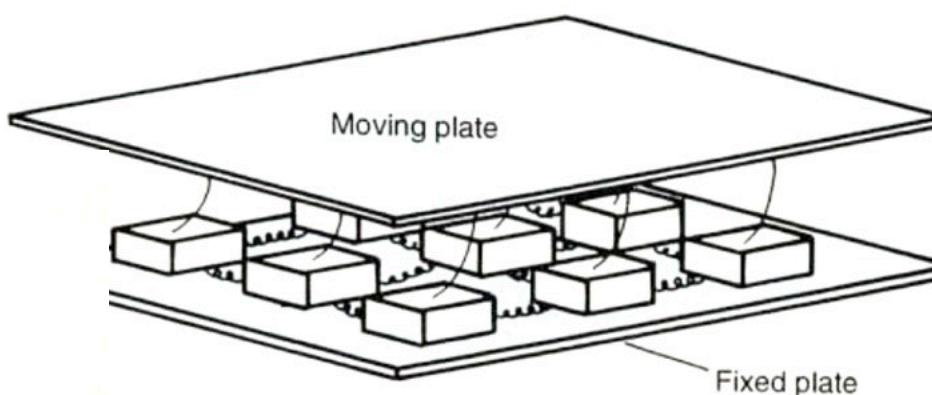


FIG. 15. Schematic diagram of the numerical model.

# NEURONS, DYNAMICS AND COMPUTATION

Brains have long been regarded as biological computers.  
But how do these collections of neurons perform computations?

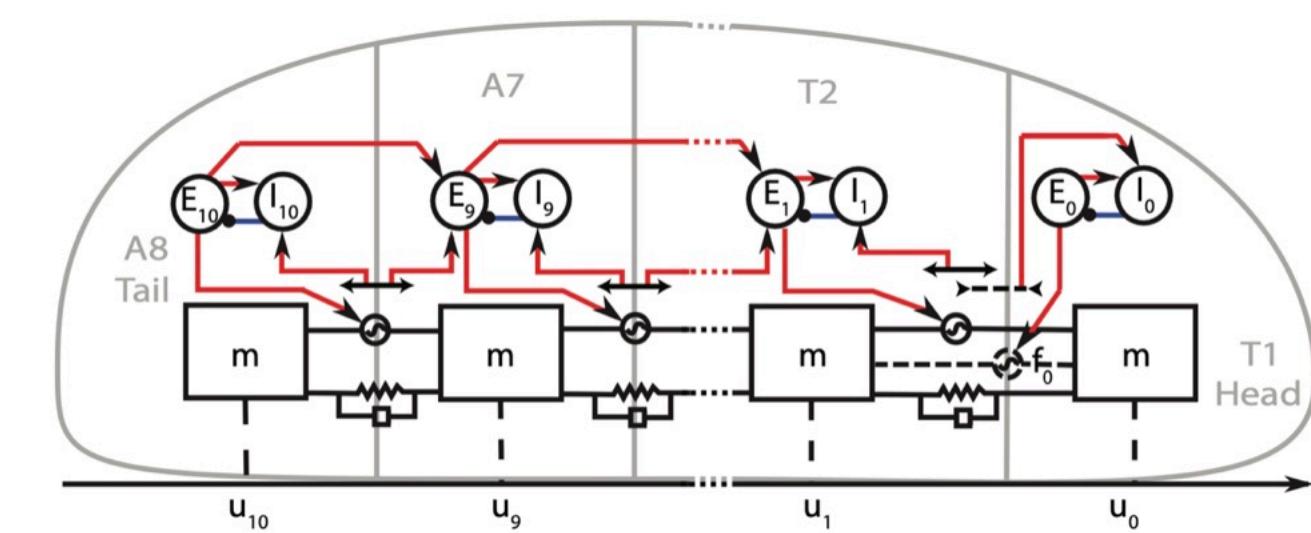
John J. Hopfield

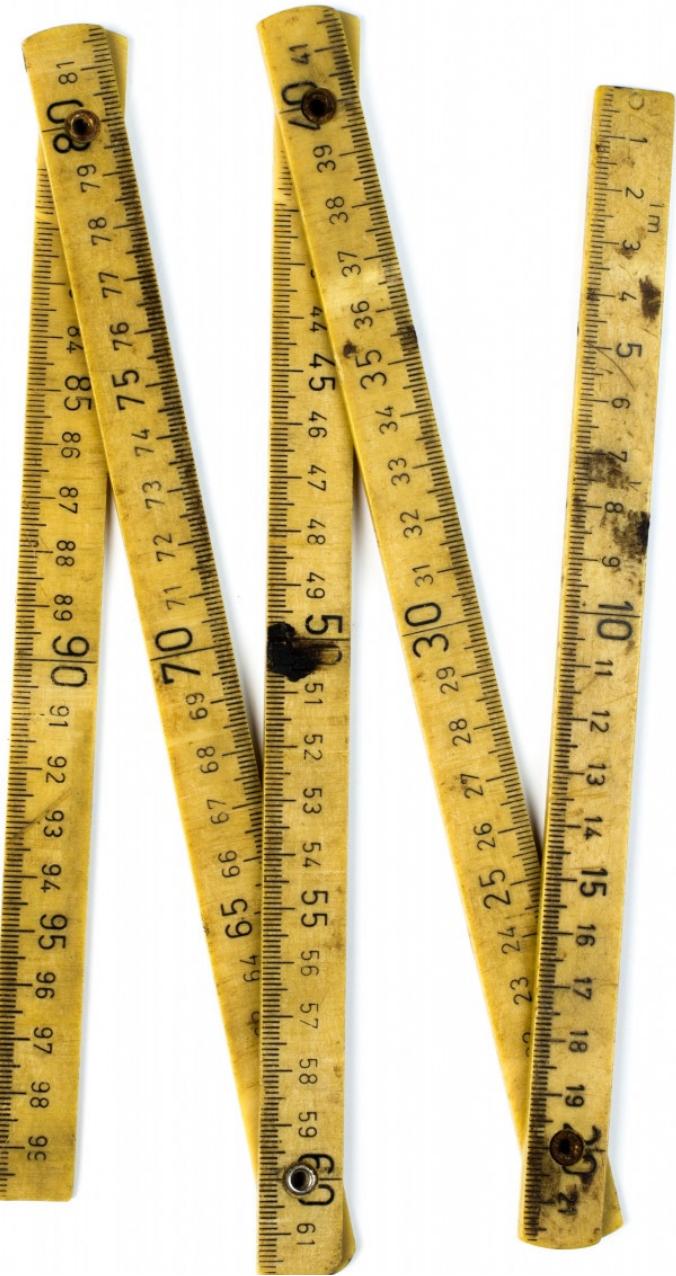
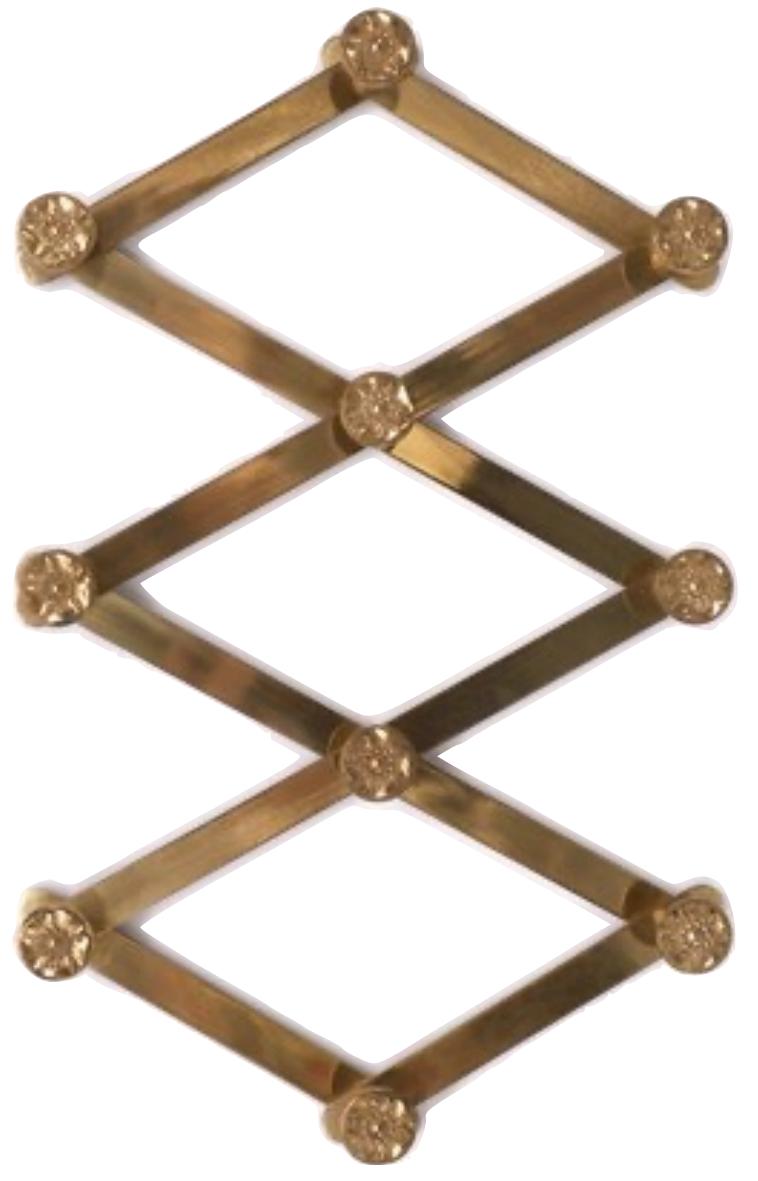


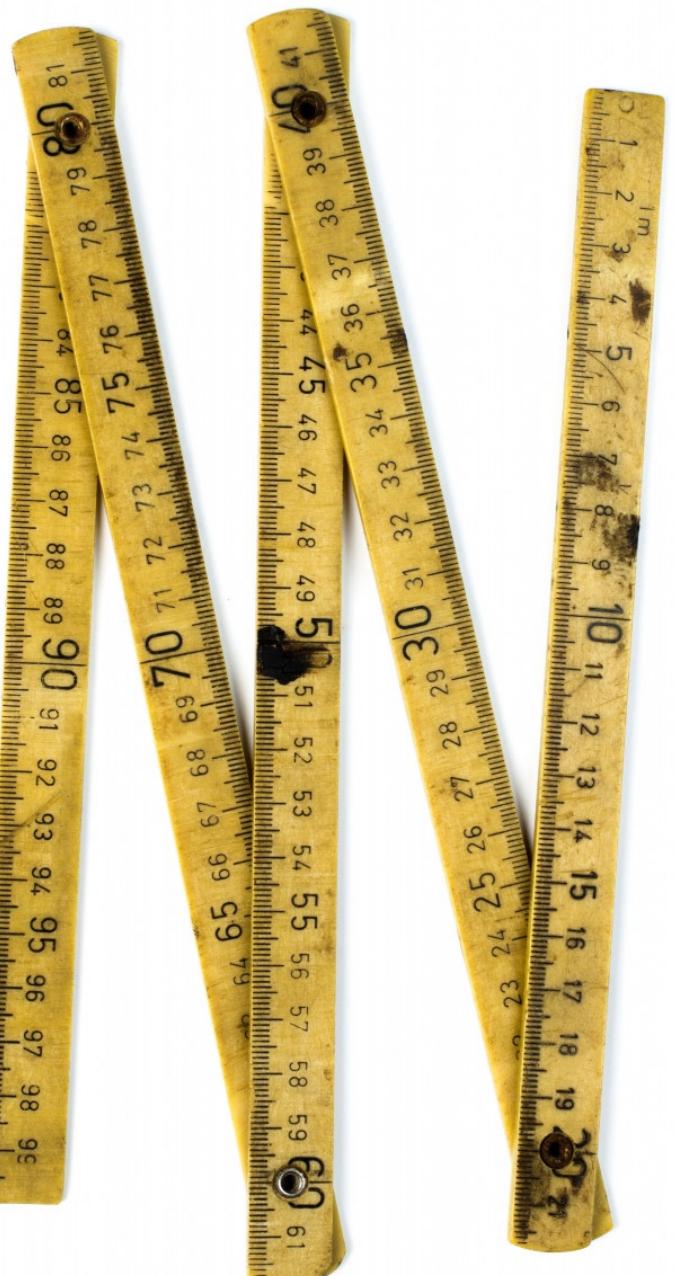
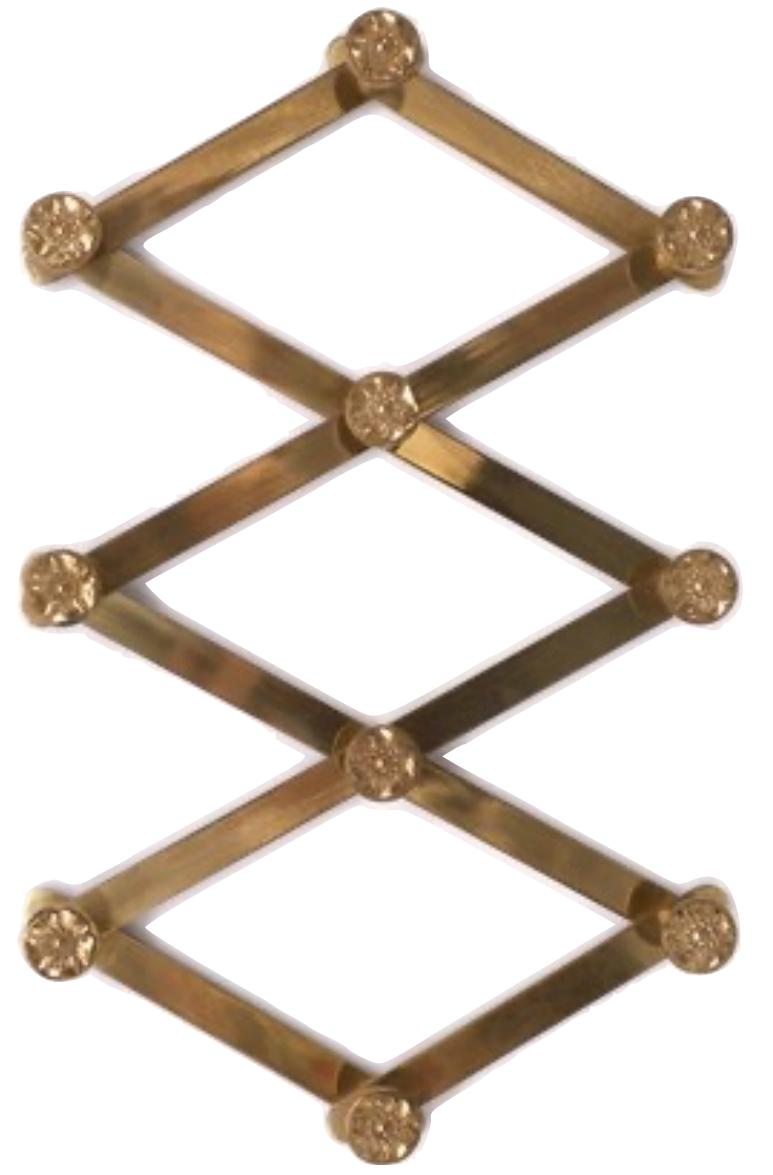
Spring-block model of earthquakes is somewhat analogous to models of sheets of neurons that account for individual neuron spikes. Springs connecting the blocks correspond to synaptic connections between neurons, while the slipping of a block corresponds to the firing of a neuron. Consequently the types of complex behavior seen in the earthquake models can be expected to recur in neurobiology. (Adapted from ref. 17.) **Figure 7**

## Integrative neuromechanics of crawling in *D. melanogaster* larvae

Cengiz Pehlevan<sup>1,2</sup>, Paolo Paoletti<sup>3</sup>, L Mahadevan<sup>4,5,6,7,8\*</sup>







## CONNECTIONS WITH THE WORK OF GERARD & AUKOSH

autonomous dynamics? (somehow)

we insist on understanding what happens over layers

## CONNECTIONS WITH THE WORK OF GERARD & AUKOSH

autonomous dynamics? (somehow)

we insist on understanding what happens over layers

## NOISE IN DYNAMICAL SYSTEMS

linearization / stabilization understood in signal processing and control theory

high dimensional setting is exciting, with focus on cascades (signal flow)

## CONNECTIONS WITH THE WORK OF GERARD & AUKOSH

autonomous dynamics? (somehow)

we insist on understanding what happens over layers

## NOISE IN DYNAMICAL SYSTEMS

linearization / stabilization understood in signal processing and control theory

high dimensional setting is exciting, with focus on cascades (signal flow)

## BILDTHEORIE

you can give a sense of training dynamics to the uninitiated!

# HERTZ, BOLTZMANN, . . .

It is the ubiquitous task of science to explain the more complex in terms of the simpler; or, if preferred, to represent [*anschaulich darstellen*] the complex by means of clear pictures [*bilder*] borrowed from the sphere of simpler phenomena.

# HERTZ, BOLTZMANN, . . .

It is the ubiquitous task of science to explain the more complex in terms of the simpler; or, if preferred, to represent [**anschaulich darstellen**] the complex by means of clear pictures [**bilder**] borrowed from the sphere of simpler phenomena.



PHYSICAL REVIEW LETTERS 134, 257301 (2025)

Editors' Suggestion

## Spring-Block Theory of Feature Learning in Deep Neural Networks

Cheng Shi<sup>1</sup>, Liming Pan<sup>2</sup>, and Ivan Dokmanić<sup>1,3,\*</sup>

<sup>1</sup>Departement Mathematik und Informatik, University of Basel, Spiegelgasse 1, 4051 Basel, Switzerland

<sup>2</sup>School of Cyber Science and Technology, University of Science and Technology of China, 230026 Hefei, China

<sup>3</sup>Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, 306 North Wright Street, Urbana, Illinois 61801, USA