

# A study on quantifying effective training of DLDMD

Joseph A.G. Diaz

Master of Science in Applied Mathematics  
with a Concentration in Dynamical Systems,  
San Diego State University

August 2, 2022

# Introduction

In the study of dynamical systems a central problem is how to derive models from measured data to facilitate the prediction of future states. The data-driven method Dynamic Mode Decomposition and its extensions offer a compelling avenue in the problem of prediction from time-series data.

The marriage of these methods with Machine Learning and Neural Networks allows for leveraging the power of these tools in the space.

If standard metrics of model training are unavailable, what other tools can we use to analyze the performance of a Machine Learning algorithm?

# Introduction – Koopmanism

- Given  $\{\mathbf{y}_j\}_{j=1}^{N^T+1}$ , we seek a model such that

$$\frac{d\mathbf{y}}{dt} = f(\mathbf{y}(t)), \quad \mathbf{y}(0) = \mathbf{x} \in \mathcal{M} \subseteq \mathbb{R}^{N_s} \quad (1)$$

- By [2], Denote  $\varphi(t; \mathbf{x}) = \mathbf{y}(t)$  to be the flow map from  $\mathbf{x}$  and  $g : \mathcal{M} \rightarrow \mathbb{C}$ , be a square integrable observable then,  $\exists \mathcal{K}$  such that

$$\mathcal{K}^t g(\mathbf{x}) = g(\varphi(t; \mathbf{x})), \quad (2)$$

- Finding the eigen-values  $\{\lambda_\ell\}$  and eigen-functions  $\{\phi_\ell\}$  of  $\mathcal{K}$  yields

$$\mathcal{K}^t \phi_\ell = \exp(t\lambda_\ell) \phi_\ell \implies g(\mathbf{x}) = \sum_{\ell \in \mathbb{N}} a_\ell \phi_\ell(\mathbf{x}), \quad (3)$$

# Introduction – Time advancement

- From here advancing the dynamics to time  $t$  is equivalent to writing

$$\mathcal{K}^t g(\mathbf{x}) = \sum_{\ell \in \mathbb{N}} a_\ell \exp(t \lambda_\ell) \phi_\ell(\mathbf{x}) \quad (4)$$

- If we suppose

$$g(\mathbf{x}) = \sum_{\ell=1}^{N_O} a_\ell \psi_\ell(\mathbf{x}) \quad (5)$$

with basis  $\{\psi_\ell\}_{\ell=1}^{N_O}$  of a subspace, then applying  $\mathcal{K}$  for discrete time implies that

$$\mathcal{K}^{\delta t} g(\mathbf{x}) = \sum_{\ell=1}^{N_O} a_\ell \exp(\delta t \lambda_\ell) \psi_\ell(\mathbf{x}) = \sum_{\ell=1}^{N_O} \psi_\ell(\mathbf{x}) (\mathbf{K}_O^T \mathbf{a})_\ell + r(\mathbf{x}; \mathbf{K}_O) \quad (6)$$

# Introduction – Observables

- The  $\mathbf{K}_O$  is a one-step mapping from each data point to the next

$$\mathbf{K}_O = \operatorname{argmin}_K \|\Psi_+ - K\Psi_-||_F^2 \quad (7)$$

where

$$\Psi_- = (\Psi_1 \ \Psi_2 \ \cdots \ \Psi_{N_T}), \quad \Psi_+ = (\Psi_2 \ \Psi_3 \ \cdots \ \Psi_{N_T+1}) \quad (8)$$

are observables of the time series  $\{\mathbf{y}_j\}$ :  $\Psi_j = (\Psi_1(\mathbf{y}_j), \ \Psi_2(\mathbf{y}_j), \ \cdots, \ \Psi_{N_O}(\mathbf{y}_j))$

- Practically,  $\mathbf{K}_O$  is found using a singular value decomposition (SVD). The eigen-spectrum of  $\mathbf{K}_O$  moves us forward

$$\mathbf{K}_O = \mathbf{V}\mathbf{T}\mathbf{V}^{-1}, \quad (9)$$

with  $\lambda_\ell = \ln((\mathbf{T})_{\ell\ell})/\delta t$ .

# Introduction – DMD

- With this the dynamics can be approximated and advanced as

$$y(t; \mathbf{x}) \approx \mathbf{K}_m \exp(t\Lambda) \mathbf{V}^{-1} \Psi(\mathbf{x}) \quad (10)$$

where  $\Psi$  is the representation of the initial condition in terms of the observables,  $\mathbf{K}_m$  is the  $N_S \times N_O$  matrix whose columns are the Koopman modes and  $\Lambda$  is the diagonal matrix whose elements are  $\lambda_\ell = (\Lambda)_{\ell\ell}$ .

- This is the basis for Dynamic Mode Decomposition (DMD) when  $N_O = N_S$  and  $\Psi_i(\mathbf{y}_j) = (\mathbf{y}_j)_i$ ; the ultimate takeaway is that if an optimal set of observables is found, the error from DMD can be reduced.

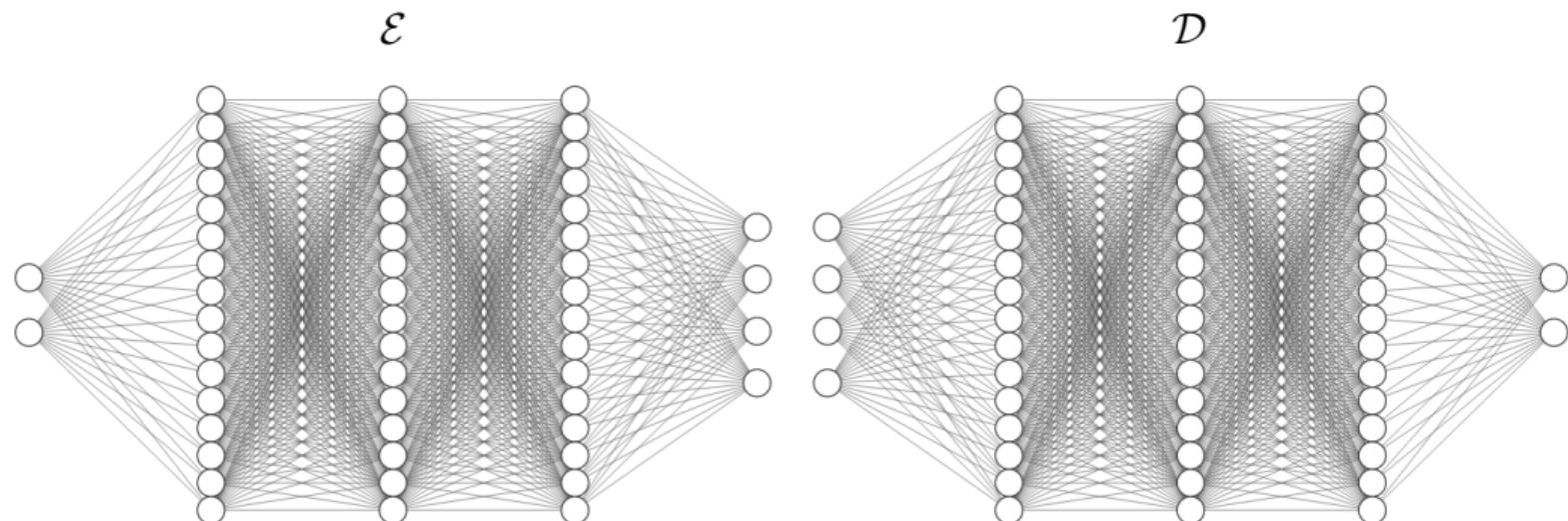
- The key innovation of Lago [1] is to use a neural network to come up with the collection of observables on  $\{\mathbf{y}_j\}$  that allow for the best prediction of future system states, we call this method Deep Learning Enhanced DMD (DLDMD). This is implemented by defining an encoder  $\mathcal{E} : \mathbb{R}^{N_s} \rightarrow \mathbb{R}^{N_o}$  and decoder  $\mathcal{D} : \mathbb{R}^{N_o} \rightarrow \mathbb{R}^{N_s}$  composed of dense layers such that

$$(\mathcal{D} \circ \mathcal{E})(\mathbf{x}) = \mathbf{x} \quad (11)$$

- We choose  $N_o \geq N_s$  and an appropriate loss function so that  $\mathcal{E}$  and  $\mathcal{D}$  give a richer space of observables, called the latent space, for Extended DMD (EDMD) to use when advancing the dynamics. The implementation of NNs for this purpose requires a method of tuning to allow  $\mathcal{E}$  and  $\mathcal{D}$  to learn the best representations possible.

# Introduction – The Network Architecture

Example of DLDMD network with  $N_S = 2$ ,  $N_O = 4$ , and  $N_L = 3$  where every hidden layer has 16 neurons. The layers in the Figure are labeled, sequentially, left to right: Enc in, Enc 0, Enc 1, Enc 2, Enc out, Dec in, Dec 0, Dec 1, Dec 2, Dec out.



# Introduction – The Dilemma

- For dense layers, passing a vector of data  $x \in \mathbb{R}^d$  through a layer  $L$  can be written as

$$L(x; \mathbf{A}, \sigma, \mathbf{b}) = \sigma(\mathbf{A}x + \mathbf{b}) \quad (12)$$

for some matrix  $\mathbf{A} \in \mathbb{R}^{N_O \times d}$ , vector  $\mathbf{b} \in \mathbb{R}^{N_O}$  and activation function  $\sigma : \mathbb{R}^{N_O} \rightarrow \mathbb{R}^{N_O}$ .

- Given a training procedure  $P$ , we can characterize the evolution with

$$Q_{n+1} = P(Q_n) \quad (13)$$

In the interest of examining convergences, one might consider the following limit

$$\lim_{n \rightarrow \infty} \|Q_{n+1} - Q_n\|_2 \quad (14)$$

for some two-norm on the space that  $Q_n$  inhabits.

- Convergence is meaningful, but this might not actually tell us much of anything else.

# Kullback-Leibler Divergence – Entropy

- For a continuous distribution with density function  $f(x)$ , the entropy is given by

$$h[f] = \mathbb{E}[-\log(f(x))] = - \int_X f(x) \log(f(x)) dx \quad (15)$$

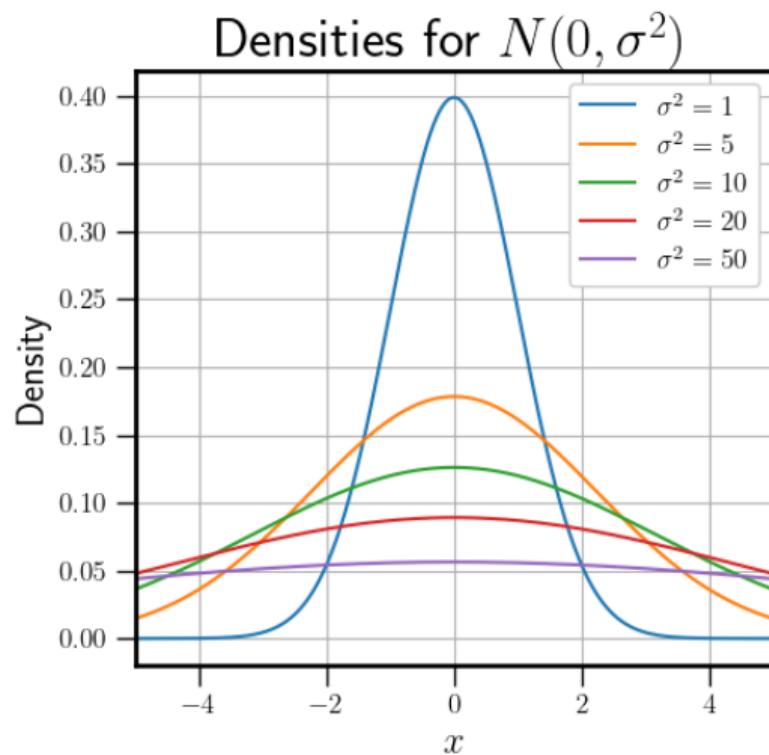
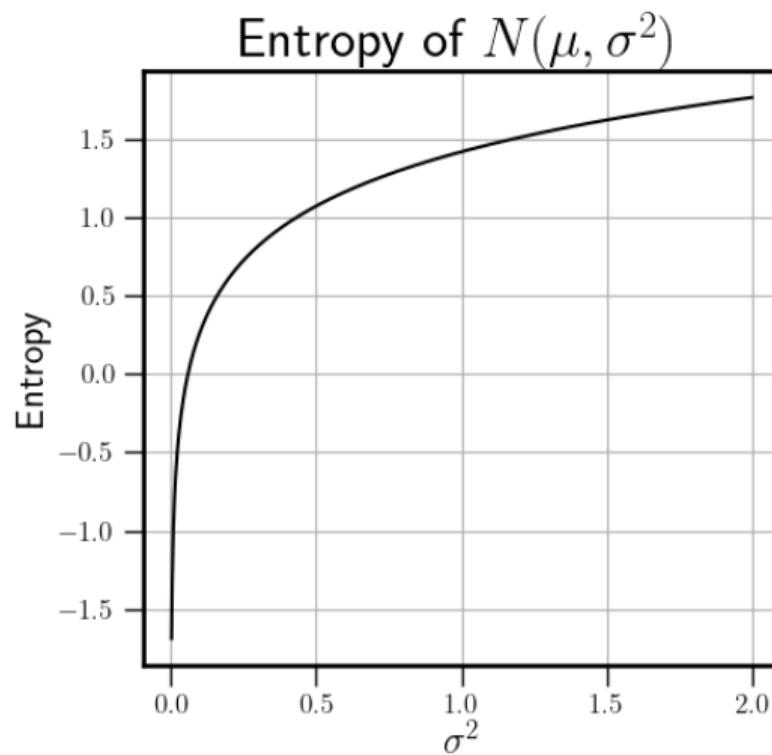
- For a normal probability distribution, we have the probability density function

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right) \quad (16)$$

and the entropy integral for this evaluates to

$$h[f] = \frac{1}{2} (\log(2\pi\sigma^2) + 1) \quad (17)$$

# Kullback-Leibler Divergence – Entropy example



# Kullback-Leibler Divergence – Basic Definitions

- The Kullback-Leibler Divergence (KLD) is a statistical distance between a pair of probability distributions which measures how different a distribution  $P$  is from a reference distribution  $Q$ . If  $P$  and  $Q$  are continuous probability distributions defined on  $X$  with probability density functions  $p$  and  $q$ , the KLD formula is

$$D_{KL}(P \parallel Q) = \int_X p(x) \log(p(x)/q(x)) \ dx \quad (18)$$

- The reason we can consider this a “distance” is that the KLD is non-negative,  $D_{KL}(P \parallel Q) \geq 0$ .

# Kullback-Leibler Divergence – Statistical Distance

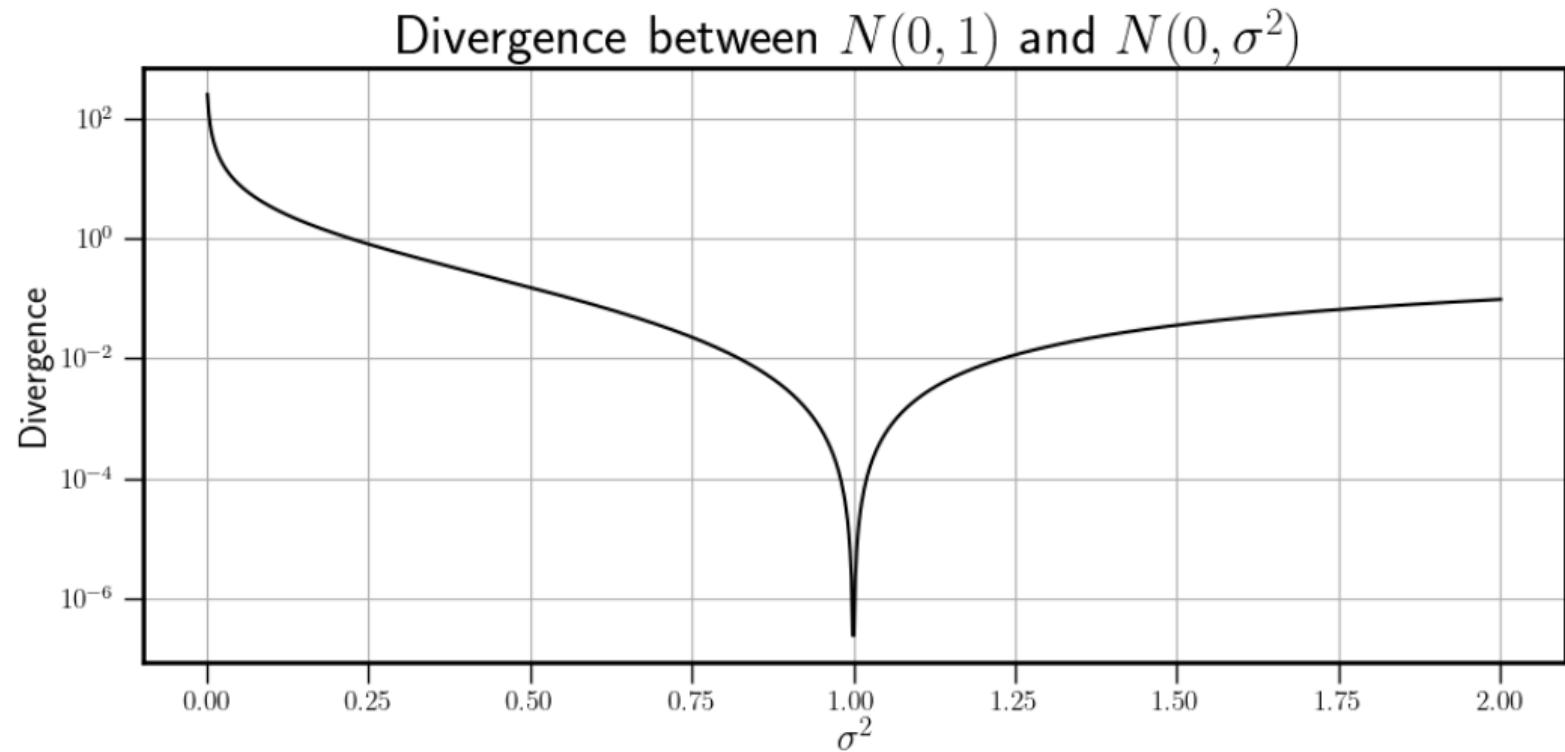
- For example, consider the random variables  $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$  with the appropriate density functions; the divergence between the distributions is

$$D_{KL}(p \parallel q) = \frac{1}{2} \left( \ln \left( \frac{\sigma_2^2}{\sigma_1^2} \right) + \frac{\sigma_1^2}{\sigma_2^2} + \frac{(\mu_1 - \mu_2)^2}{\sigma_2^2} - 1 \right), \quad (19)$$

- Letting  $\mu_1 = \mu_2$ ,  $\sigma_1^2 = 1$ , and allowing  $\sigma_2^2$  to vary, we can write

$$D_{KL}(p \parallel q) = \frac{1}{2} \left( \ln \sigma_2^2 + \frac{1}{\sigma_2^2} - 1 \right), \quad (20)$$

# Kullback-Leibler Divergence – normal distribution example



# Kullback-Leibler Divergence – Kernel Density Estimation

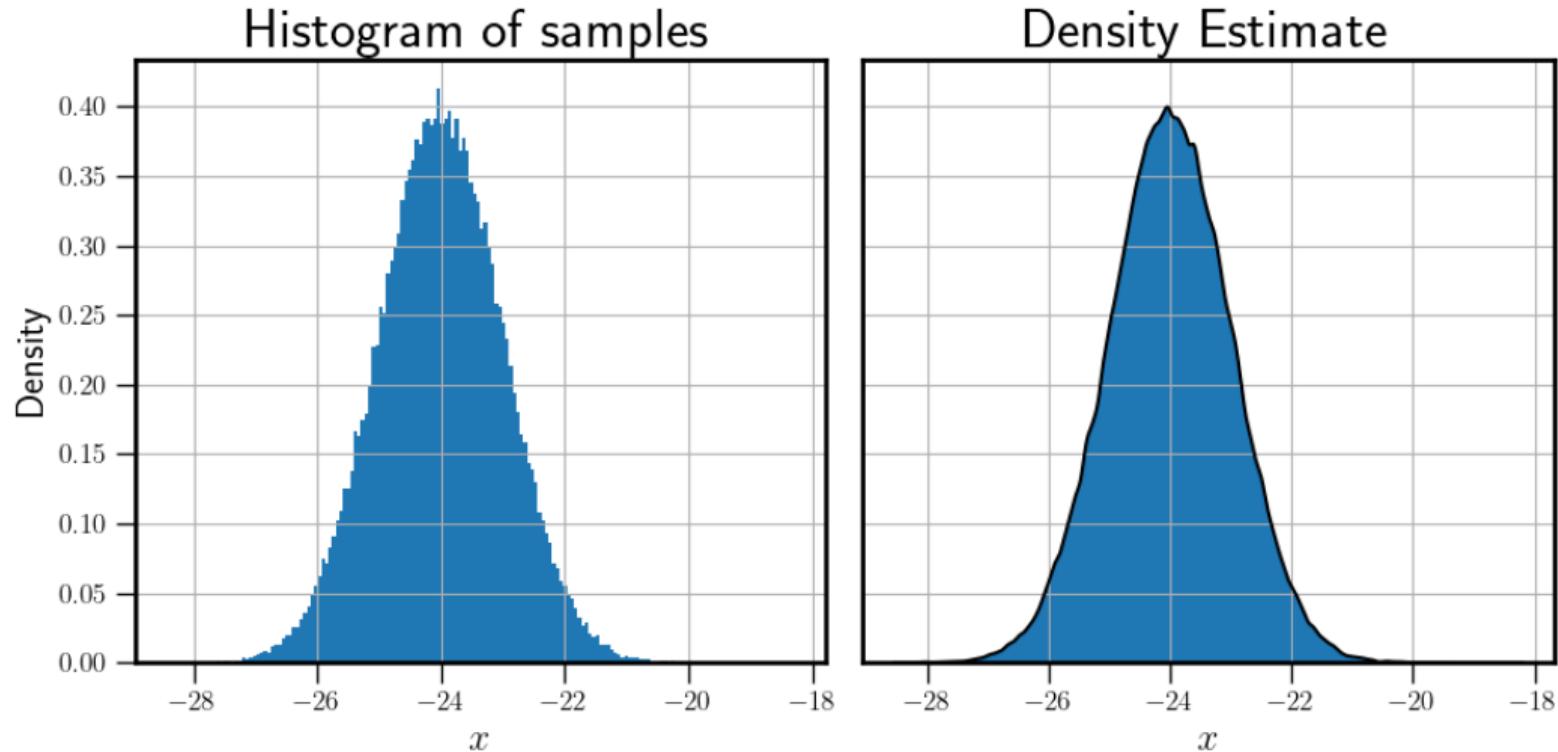
- Kernel Density Estimation (KDE) is a method to approximate a probability density  $f$  from measured data in the form of a histogram. The *kernel density estimator* for  $f$  given by

$$\hat{f}(x; K, h) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right), \quad (21)$$

where  $K$  is the chosen kernel,  $h$  is what's called the bandwidth parameter, and  $\{x_i\}_{i=1}^n$  is the data that generates our histogram.

- The choice of kernel does not provide much statistical significance; but the choice of the bandwidth parameter  $h$  is very crucial for finding a density estimate that approximates the underlying density function appropriately and can be thought of as an analogue of the bin width for the affiliated histogram.

# Kullback-Leibler Divergence – KDE example



# Kullback-Leibler Divergence – Measuring Entropy Flow

- For a converging machine, for each layer we should have that

$$\mathbf{W}_{I,\mathcal{E}}^{(n)}, \mathbf{W}_{I,\mathcal{D}}^{(n)} \rightarrow \mathbf{W}_{I,\mathcal{E}}^*, \mathbf{W}_{I,\mathcal{D}}^* \text{ as } n \rightarrow \infty \quad (22)$$

Consequently, we can characterize any set of weights via the formulas

$$\mathbf{W}_{I,\mathcal{E}}^{(n)} = \mathbf{W}_{I,\mathcal{E}}^* + \mathbf{W}_{I,\mathcal{E}}^{(n),f}, \quad \mathbf{W}_{I,\mathcal{D}}^{(n)} = \mathbf{W}_{I,\mathcal{D}}^* + \mathbf{W}_{I,\mathcal{D}}^{(n),f}, \quad (23)$$

where  $\mathbf{W}_{I,\mathcal{E}}^{(n),f}$  and  $\mathbf{W}_{I,\mathcal{D}}^{(n),f}$  are fluctuations from the steady state. First-order differencing gives us detrended matrices:

$$\delta \mathbf{W}_{I,\mathcal{E}}^{(n)} = \mathbf{W}_{I,\mathcal{E}}^{(n+1)} - \mathbf{W}_{I,\mathcal{E}}^{(n)} = \mathbf{W}_{I,\mathcal{E}}^{(n+1),f} - \mathbf{W}_{I,\mathcal{E}}^{(n),f} \quad (24)$$

- Using KDE we generate an affiliated empirical probability distribution  $p_{I,\mathcal{E}}^{(n)}(w)$ , which we can use to find consecutive divergences with KLD:

$$D = \left\{ D_{KL} \left( p_{I,\mathcal{E}}^{(n+1)} \middle\| p_{I,\mathcal{E}}^{(n)} \right) \right\}_{n=1}^{N_E-2} \quad (25)$$

# Kullback-Leibler Divergence – Implementation details

- After the training of each model for  $N_E$  epochs, the data that we have to play with is

$$\mathbf{W}_{\mathcal{E}} = \left\{ \mathbf{W}_{I,\mathcal{E}}^{(n)} \right\}_{n=1,I=1}^{N_E, N_L}, \mathbf{W}_{\mathcal{D}} = \left\{ \mathbf{W}_{I,\mathcal{D}}^{(n)} \right\}_{n=1,I=1}^{N_E, N_L}, \quad (26)$$

which are the sets of weights of the layers of the encoder and decoder. For each layer  $I$ , the set of detrended matrices are computed

$$\delta \mathbf{W}_{I,\mathcal{E}} = \left\{ \mathbf{W}_{I,\mathcal{E}}^{(n+1),f} - \mathbf{W}_{I,\mathcal{E}}^{(n),f} \right\}_{n=1}^{N_E-1}, \quad (27)$$

these matrices are flattened and their entries are used as the data for KDE.

- The kernel used is the Epanechnikov kernel, defined as

$$K(x) = \begin{cases} \frac{3(1-x^2)}{4}, & |x| < 1 \\ 0, & |x| \geq 1 \end{cases}, \quad (28)$$

which is optimal in the mean-square error sense. The Improved Sheather-Jones algorithm is used for  $h$ .

## Kullback-Leibler Divergence – Implementation details cont.

- From this, we obtain the set of density estimates for each detrended matrix

$$p_{I,\mathcal{E}} = \text{KDE}(\delta \mathbf{W}_{I,\mathcal{E}}) = \left\{ p_{I,\mathcal{E}}^{(n)}(w) \right\}_{n=1}^{N_E-1} \quad (29)$$

- Given these density approximations, we can now compute what we are truly interested in:

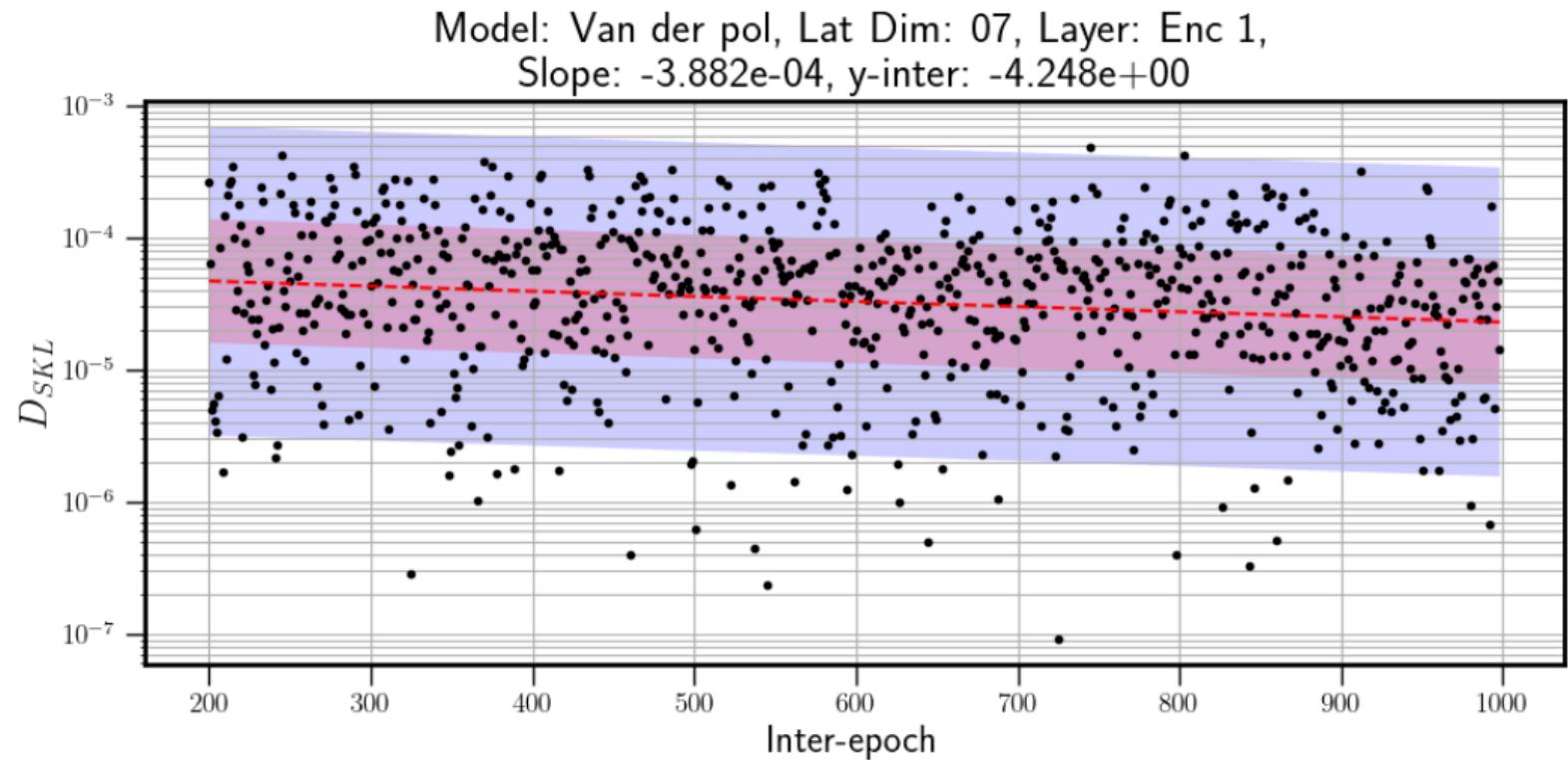
$$D_I = \left\{ D_{SKL} \left( p_{I,\mathcal{E}}^{(n+1)} \middle\| p_{I,\mathcal{E}}^{(n)} \right) \right\}_{n=1}^{N_E-2} \quad (30)$$

Here

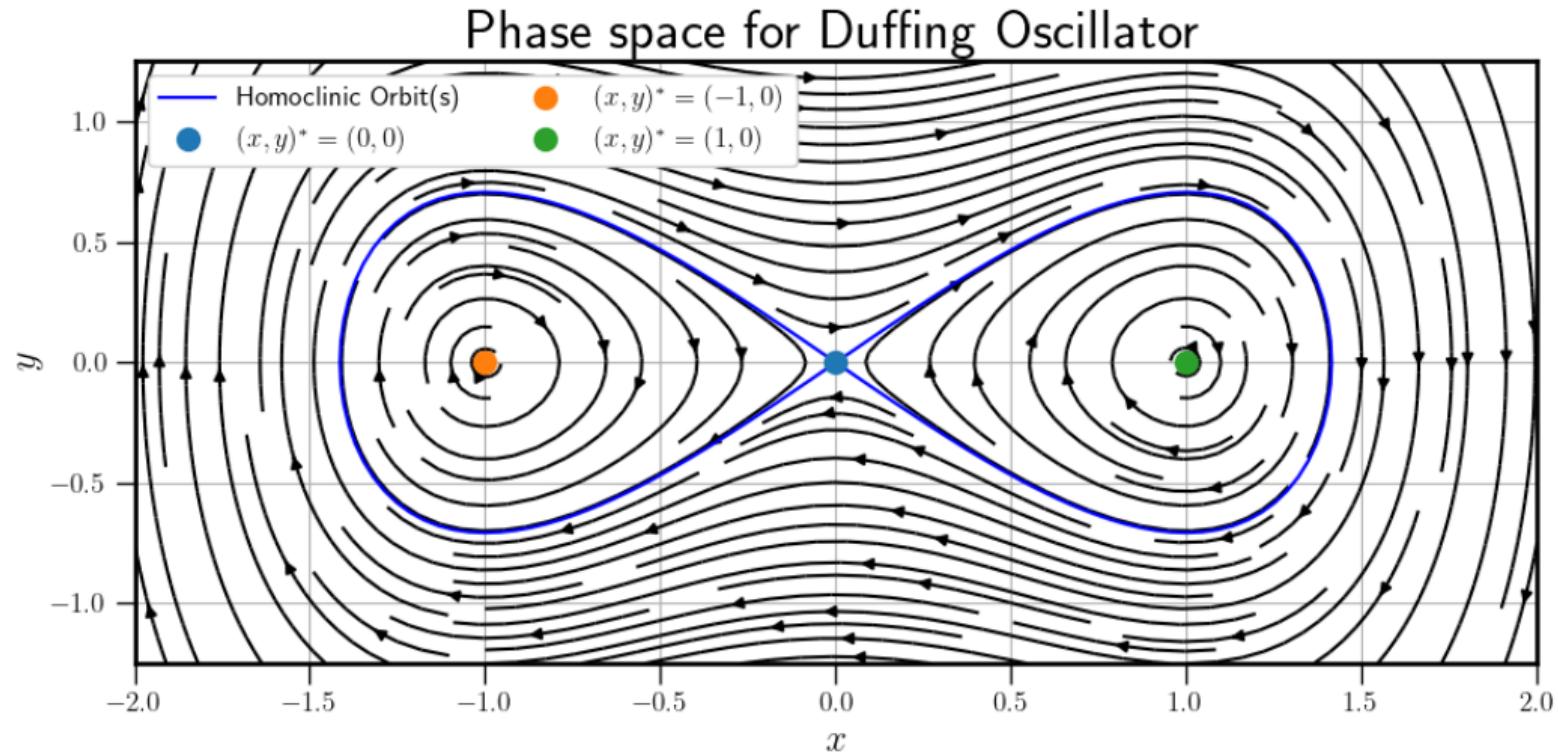
$$D_{SKL}(p, q) = \frac{1}{2} (D_{KL}(p \parallel q) + D_{KL}(q \parallel p)) \quad (31)$$

- We are interested in what these divergences tell us about information flow in the model's weights and seek to identify specific classifiers of good training.

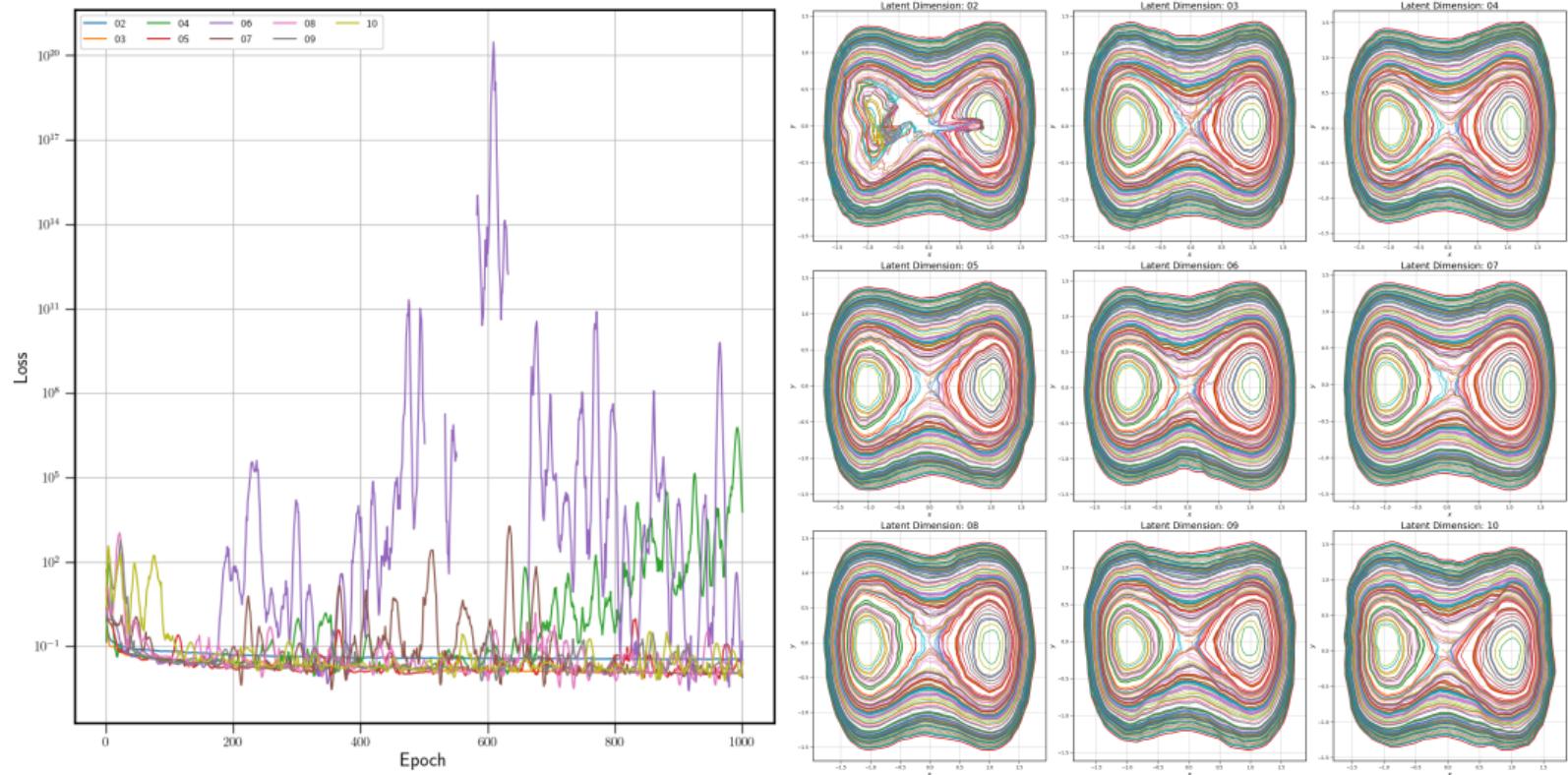
# Kullback-Leibler Divergence – Example fitting



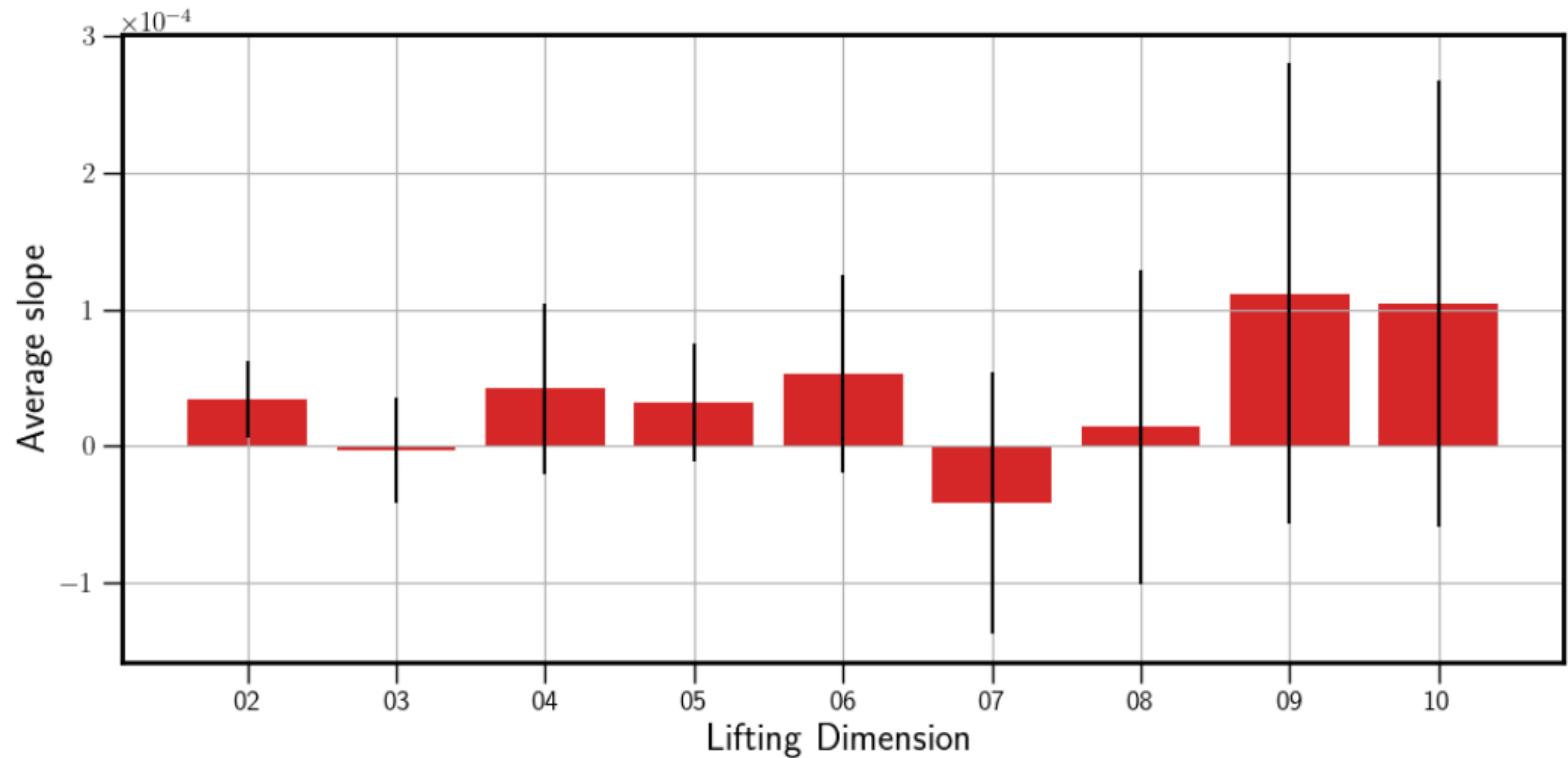
# Results - Duffing



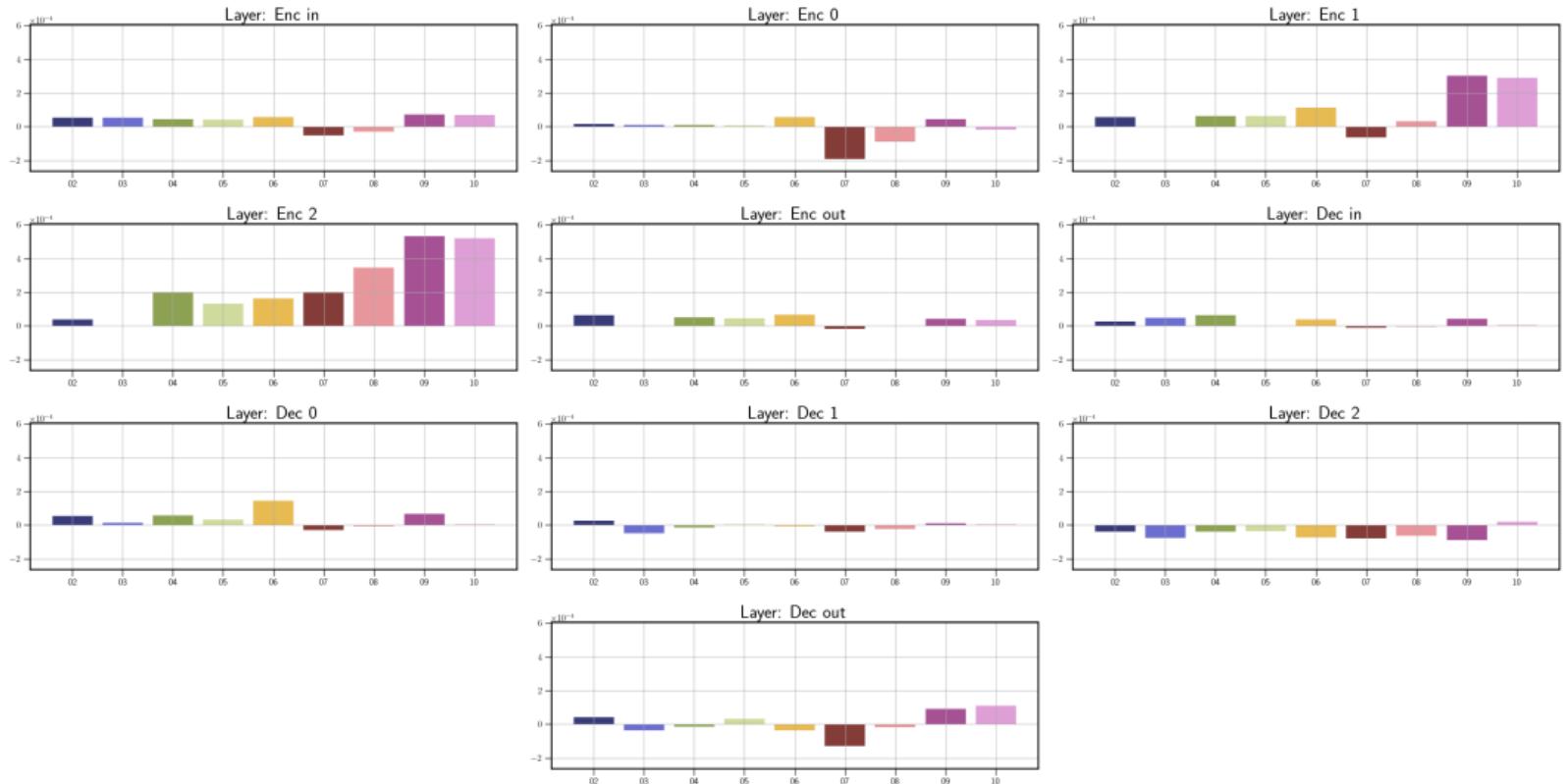
# Results – Duffing loss curves and phase space



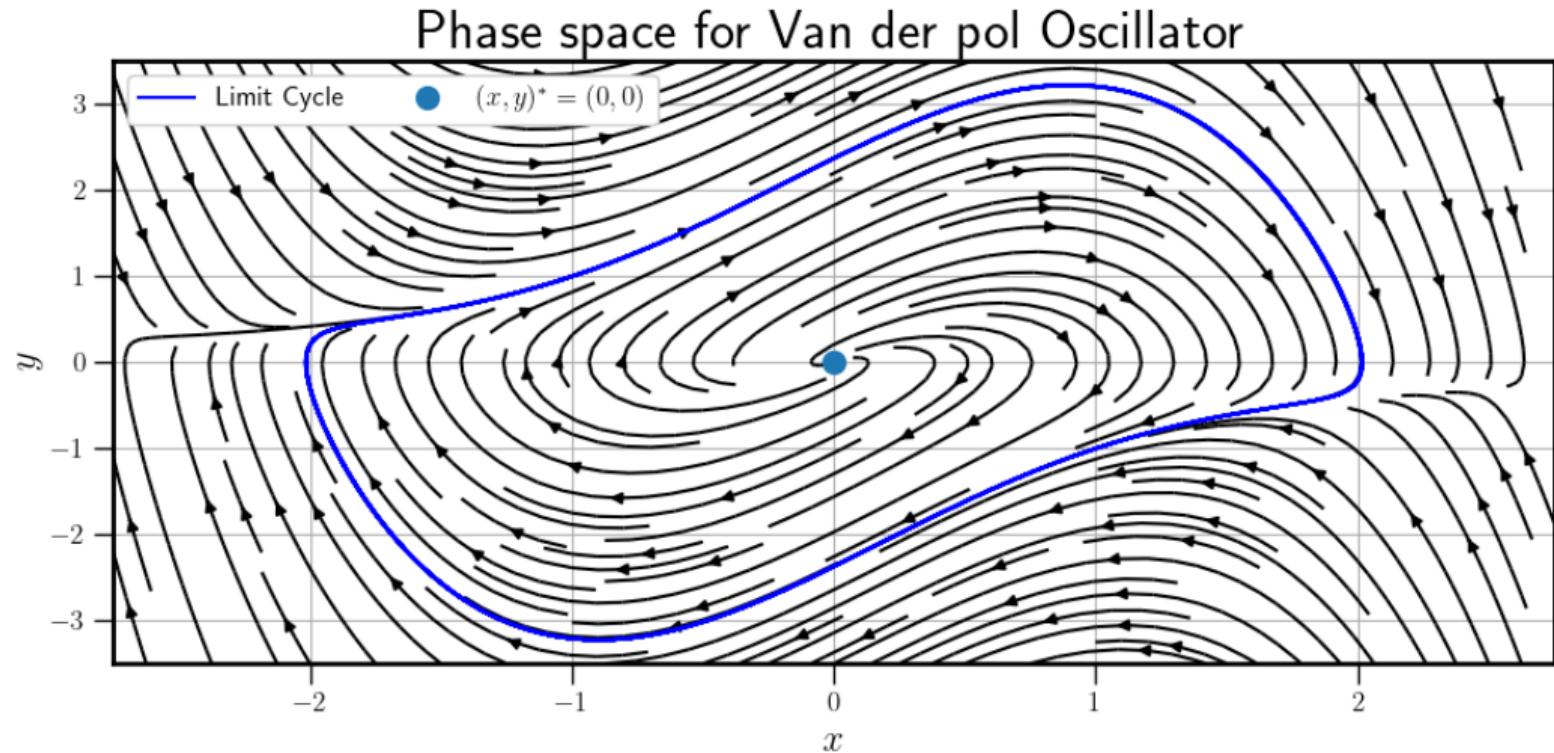
## Results – Duffing linear fit slope averages and variances



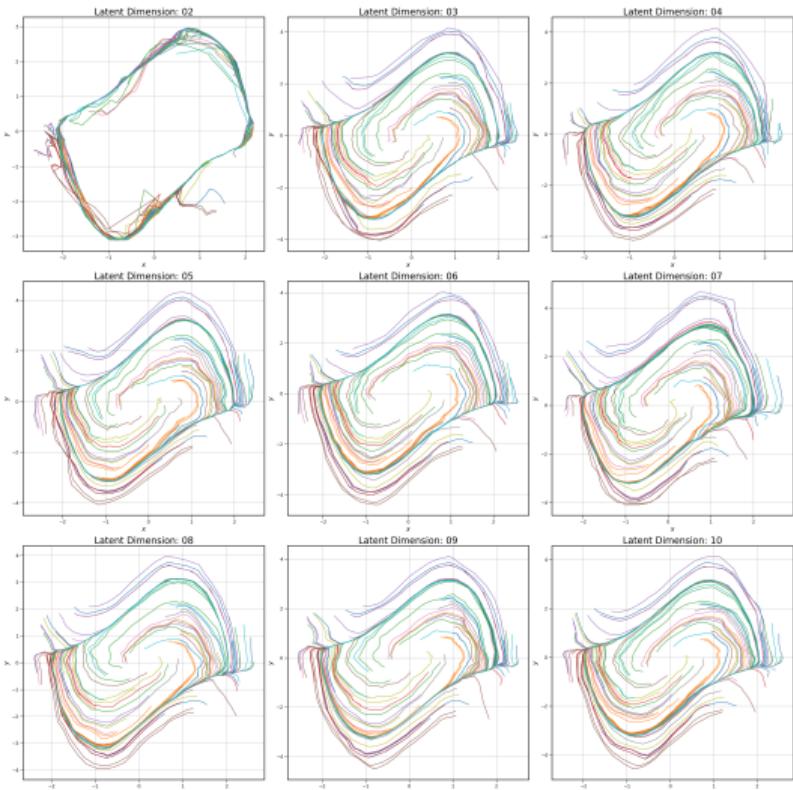
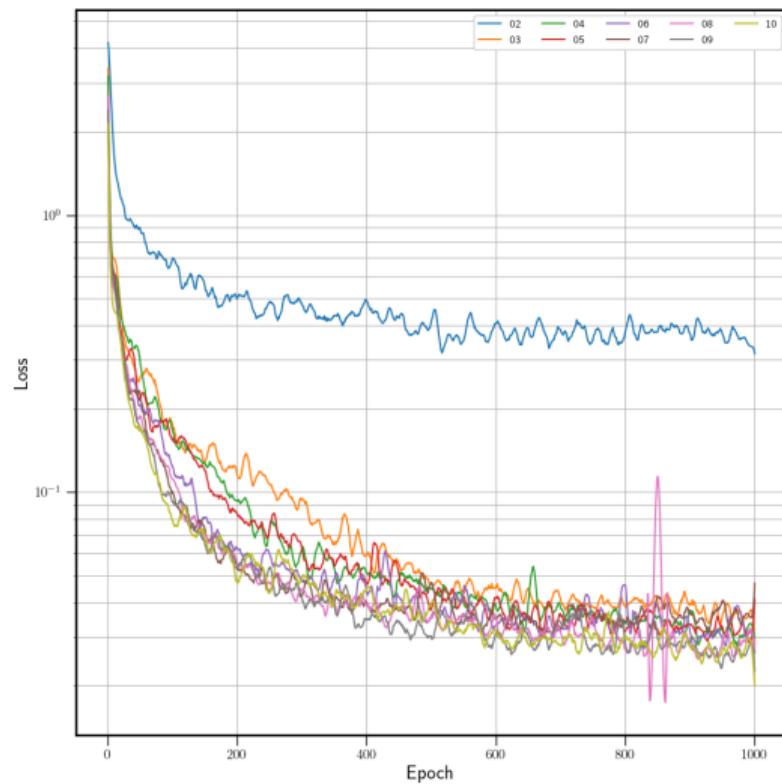
# Results – Duffing linear fit slopes



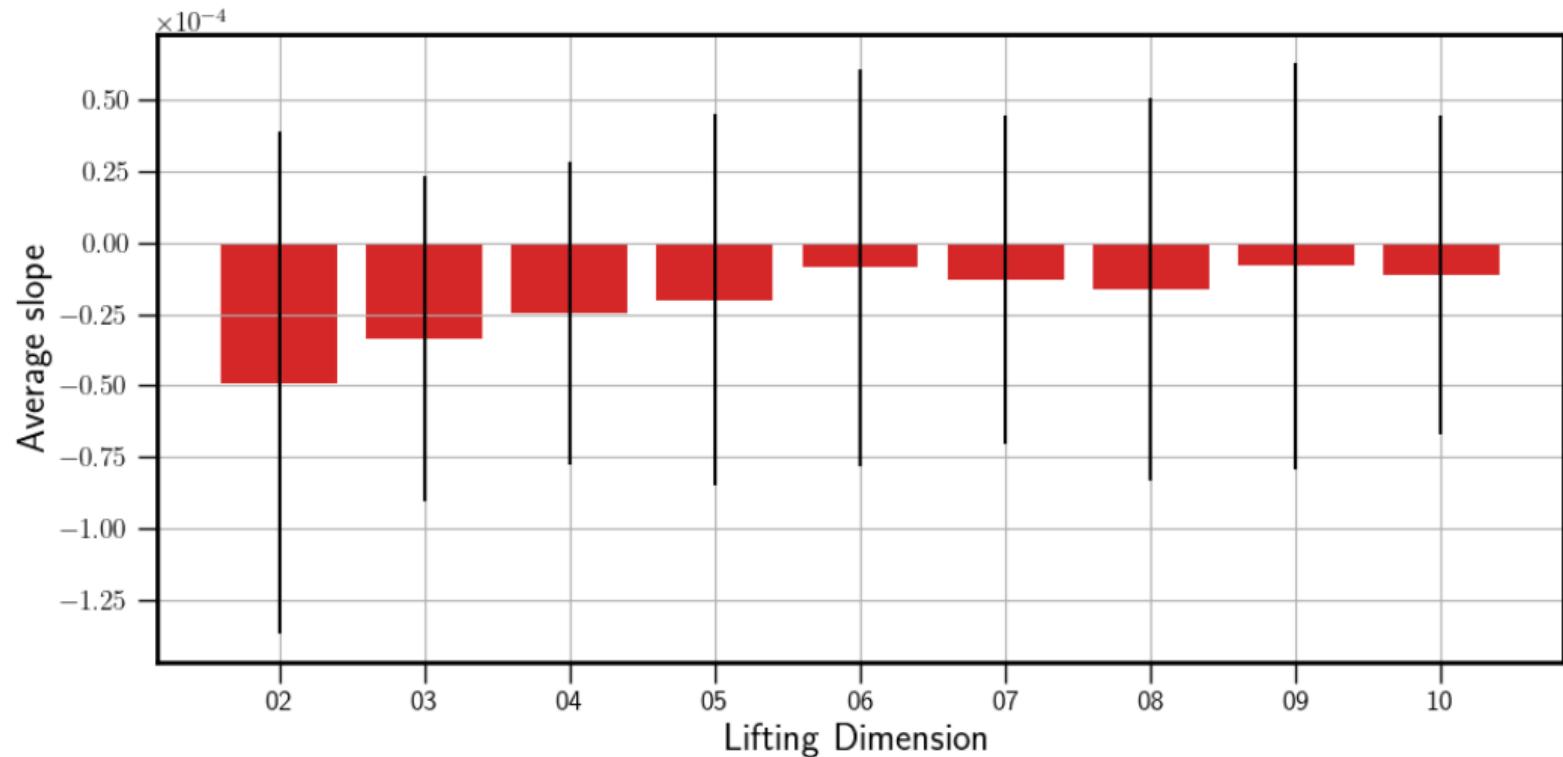
## Results – Van der Pol



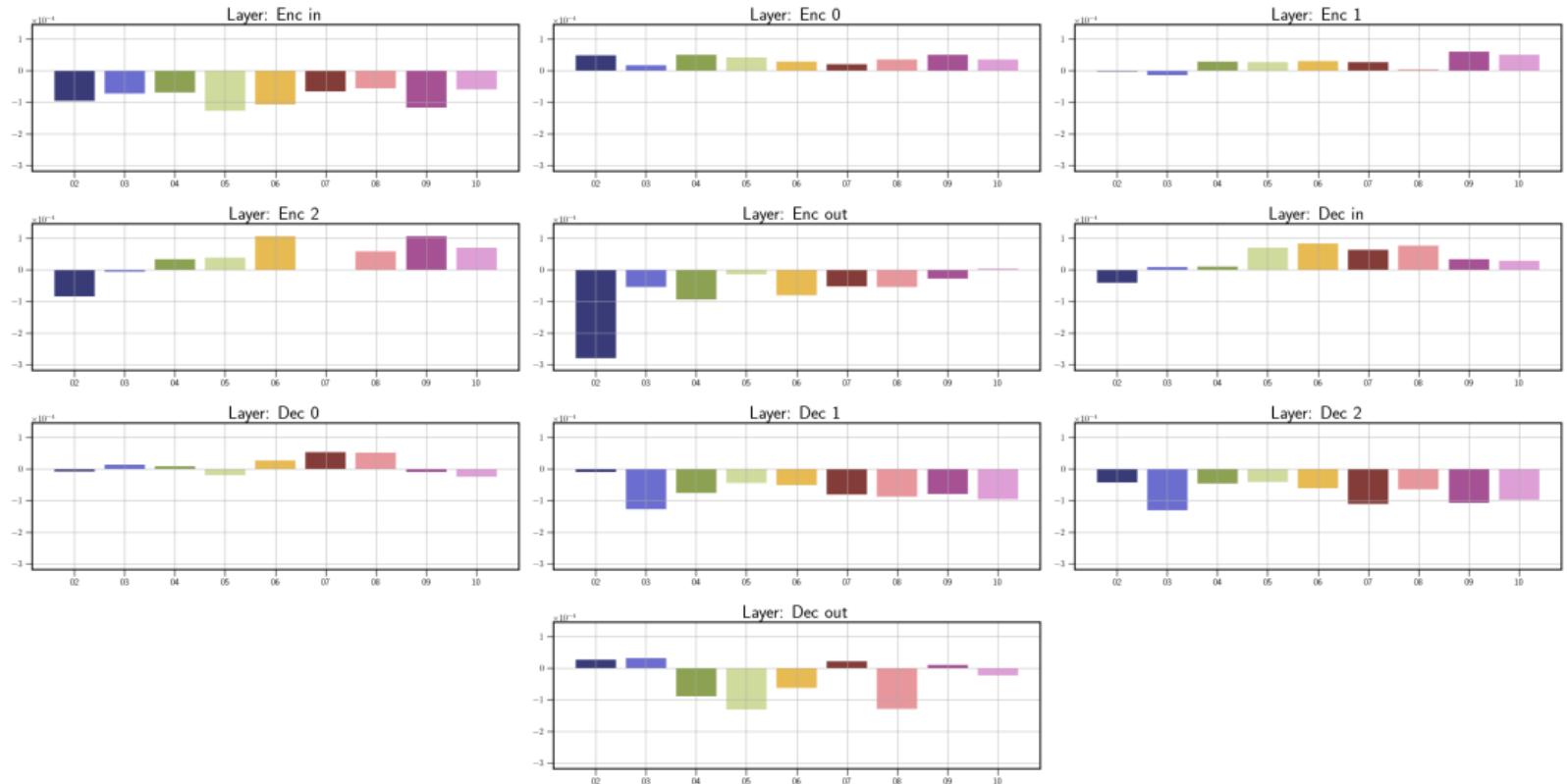
# Results – Van der Pol loss curves and phase space



## Results – Van der Pol linear fit slope averages and variances



# Results – Van der Pol linear fit slopes



# Discussion

- The Good
- Optimal Parameters
- Non-unitary densities
- Statistical issues with edge layers
- Future work

# References

-  D.J. Alford-Lago, C. W. Curtis, A. T. Ihler, and O. Issan.  
Deep Learning Enhanced Dynamic Mode Decomposition.  
*Chaos: An Interdisciplinary Journal of Nonlinear Science*, 32, 2022.
-  B.O. Koopman.  
Hamiltonian systems and transformations in Hilbert space.  
*Proc. Nat. Acad. Sci.*, 17:315–318, 1931.

The End!

Thank you for your time!  
Questions?