# Towards a Phenomenological Understanding of Neural Networks

Samuel Tovey, Sven Krippendorf, Michael Spannowsky

Konstantin Nikolaou, Christian Holm

ML4Jets, DESY, Hamburg -- stovey@icp.uni-stuttgart.de

# Why is this Necessary?

# Why is this Necessary?

| | |
|---|---|
| **Parameters** | 175 billion |
| **Training Time** | Several months |
| **Training Cost** | ~ $4.6 million |

https://openai.com, , [1] (Sterling & Laughlin, 2015), [2] (Thorpe et al., 1996), freepik

# Why is this Necessary?

| | |
|---|---|
| **Parameters** | 175 billion |
| **Training Time** | Several months |
| **Training Cost** | ~ $4.6 million |

We cannot afford to perform hyperparameter searches here.

# Why is this Necessary?

| | |
|---|---|
| **Parameters** | 175 billion |
| **Training Time** | Several months |
| **Training Cost** | ~ $4.6 million |



We cannot afford to perform hyperparameter searches here.

https://openai.com, , [1] (Sterling & Laughlin, 2015), [2] (Thorpe et al., 1996), freepik

# Why is this Necessary?

| Parameters | 175 billion | Neurons | 86 billion |
|---|---|---|---|
| Training Time | Several months | Object recognition time[2] | 150 ms |
| Training Cost | ~ $4.6 million | Energy cost[1] | < 20 W |

We cannot afford to perform hyperparameter searches here.

9/11/2023

https://openai.com, , [1] (Sterling & Laughlin, 2015), [2]  (Thorpe et al., 1996), freepik

# Collective Variables for Neural Networks

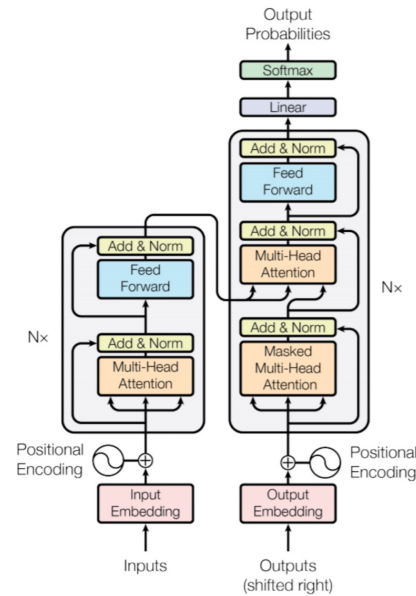Entropy, trace, and more...

# What is a Neural Network? (The NN Zoo)
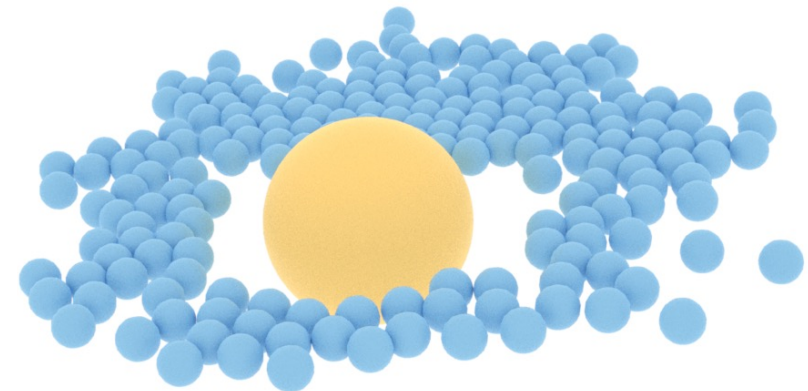
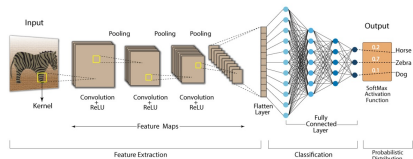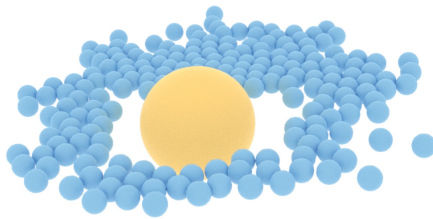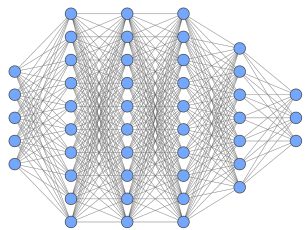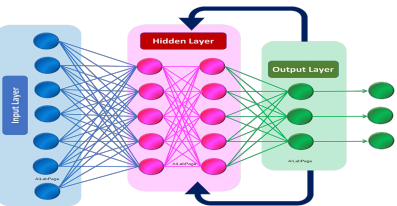## Function fitting in a very high dimensional space.



https://medium.datadriveninvestor.com/recurrent-neural-network-58484977c445



$$h_i = \sigma(\sum_{i \in N_j} c_{ij} W h_j)$$



https://developersbreach.com/convolution-neural-network-deep-learning/

https://builtin.com/artificial-intelligence/transformer-neural-network

# What is a Neural Network? (The NN Zoo)

Function fitting in a very high dimensional space.



https://developersbreach.com/convolution-neural-network-deep-learning/

https://medium.datadriveninvestor.com/recurrent-neural-network-58484977c445

https://builtin.com/artificial-intelligence/transformer-neural-network

$$f_\theta : X \to Y$$

$$\theta = \{\theta_0, \dots, \theta_N\}$$

$$\theta_i' = \theta_i - \eta \cdot \partial_{\theta_i} \mathcal{L}(f(X), Y)$$

# How do they evolve?

Architecture component

$$f'_\theta(X) = f_\theta(x) - \begin{pmatrix} \nabla_\theta f_\theta(x_0) \cdot \nabla_\theta f_\theta(x_0) & \cdots & \nabla_\theta f_\theta(x_0) \cdot \nabla_\theta f_\theta(x_N) \\ \vdots & \ddots & \vdots \\ \nabla_\theta f_\theta(x_N) \cdot \nabla_\theta f_\theta(x_0) & \cdots & \nabla_\theta f_\theta(x_N) \cdot \nabla_\theta f_\theta(x_N) \end{pmatrix} \cdot \nabla_{f_\theta} \mathcal{L}(X)$$
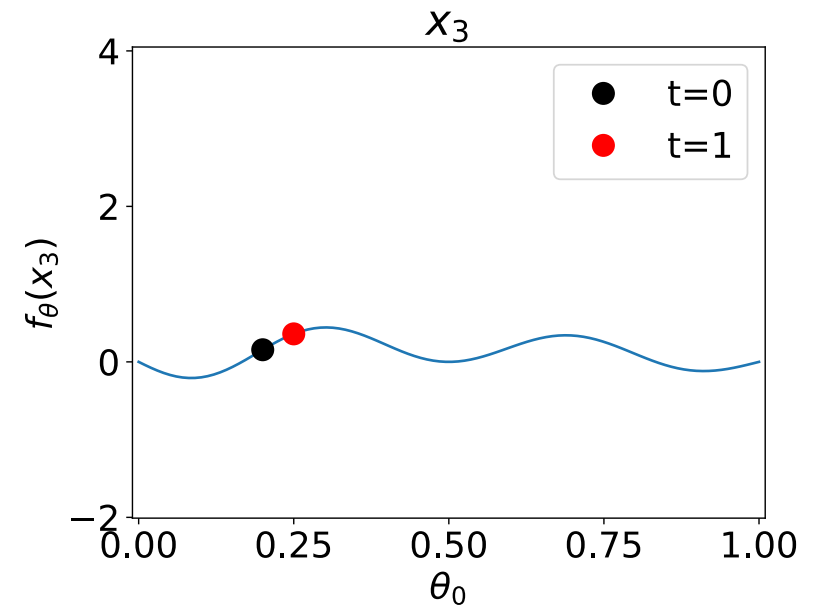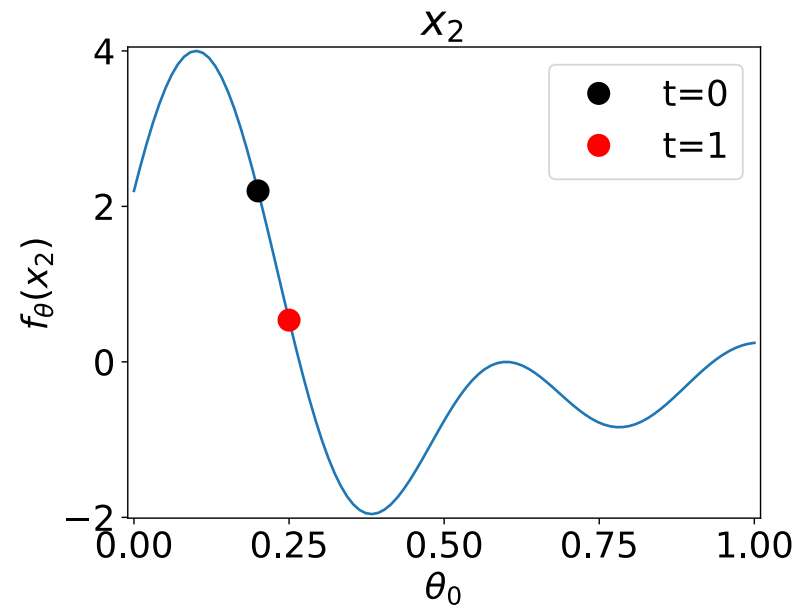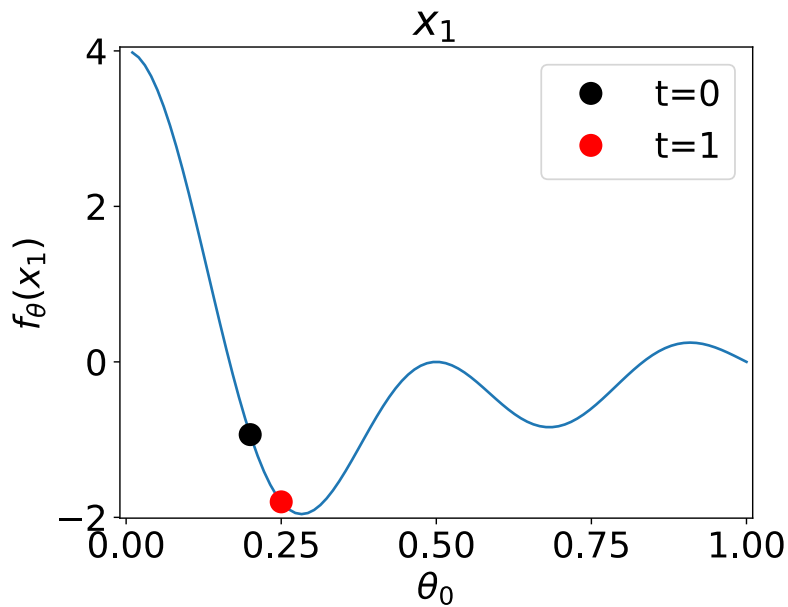
Loss component

$\Theta$

# How do they evolve?

Architecture component

$$f'_\theta(X) = f_\theta(x) - \begin{pmatrix} \nabla_\theta f_\theta(x_0) \cdot \nabla_\theta f_\theta(x_0) & \cdots & \nabla_\theta f_\theta(x_0) \cdot \nabla_\theta f_\theta(x_N) \\ \vdots & \ddots & \vdots \\ \nabla_\theta f_\theta(x_N) \cdot \nabla_\theta f_\theta(x_0) & \cdots & \nabla_\theta f_\theta(x_N) \cdot \nabla_\theta f_\theta(x_N) \end{pmatrix} \cdot \nabla_{f_\theta} \mathcal{L}(X)$$

Loss component

$\Theta$



9/11/2023

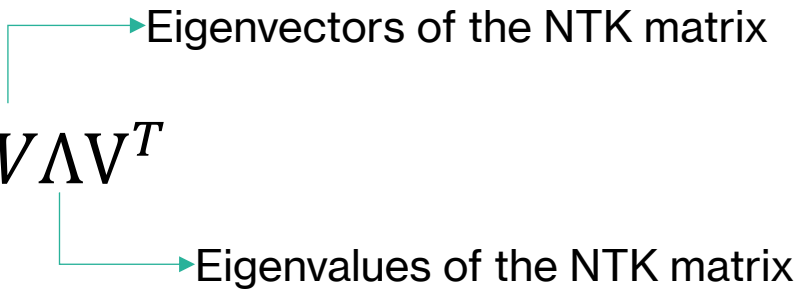# Collective Variables

Eigenvectors of the NTK matrix

$$\Theta = V\Lambda V^T$$

Eigenvalues of the NTK matrix

# Collective Variables

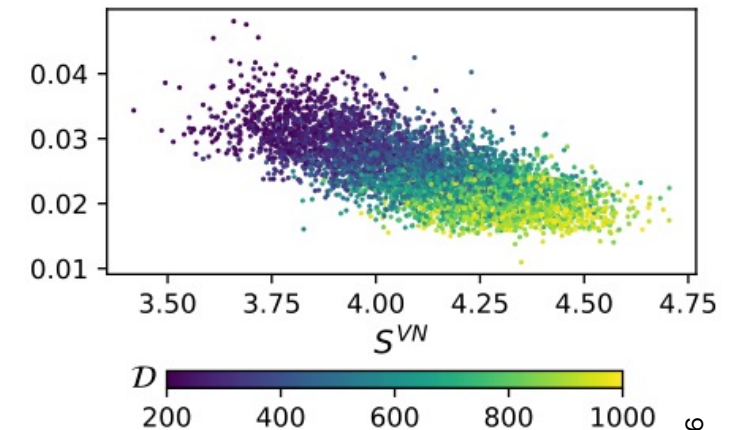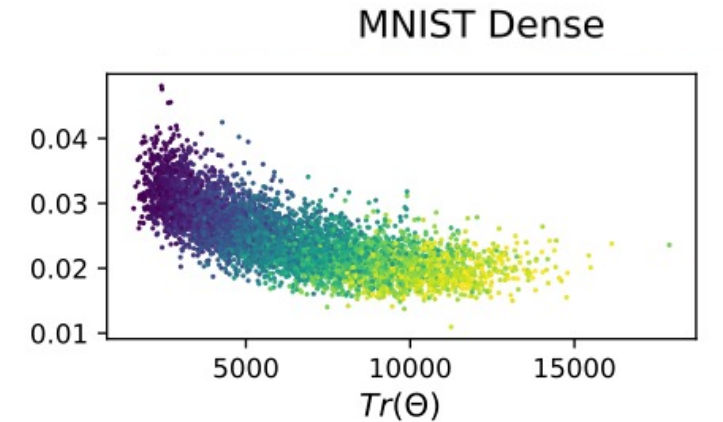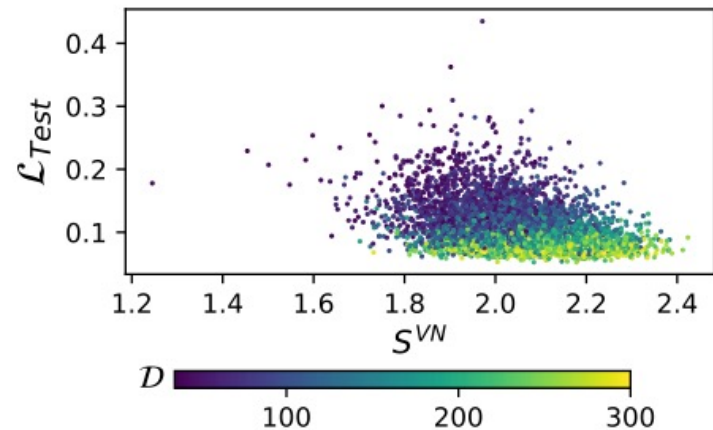Eigenvectors of the NTK matrix

$$\Theta = V \Lambda V^T$$

Eigenvalues of the NTK matrix

$$S = -\sum_i \lambda_i \cdot \ln \lambda_i$$

Measure of correlation in data

$$Tr(\Theta) = \sum_i \Theta_{ii}$$

Weighting of largest step direction

# Collective Variables

Eigenvectors of the NTK matrix
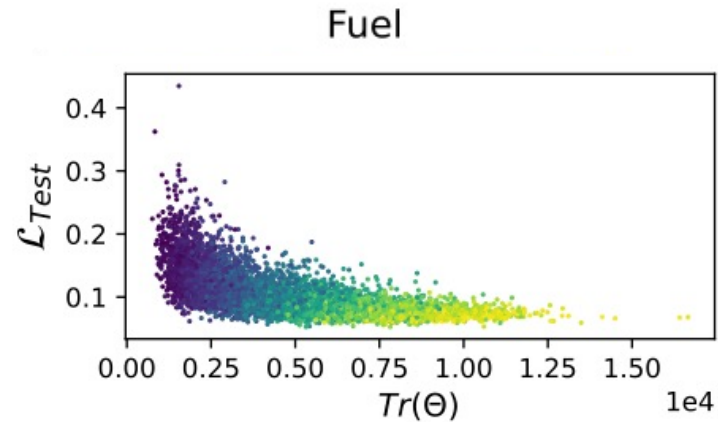
$$\Theta = V \Lambda V^T$$

Eigenvalues of the NTK matrix

$$S = - \sum_i \lambda_i \cdot \ln \lambda_i$$

Measure of correlation in data

$$Tr(\Theta) = \sum_i \Theta_{ii}$$

Weighting of largest step direction

9/11/2023

# Collective Variables

Eigenvectors of the NTK matrix
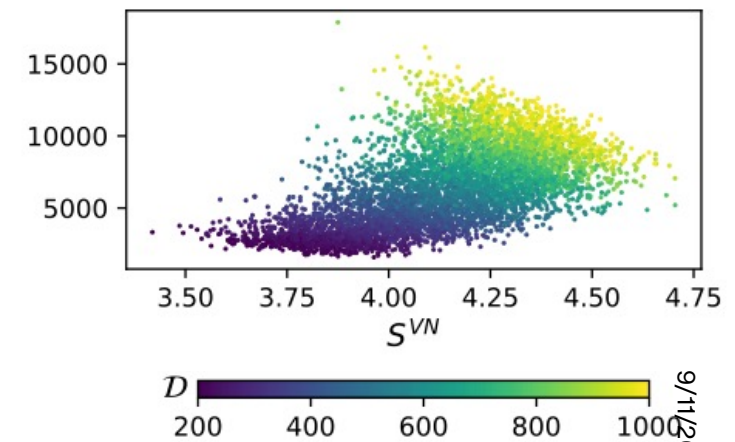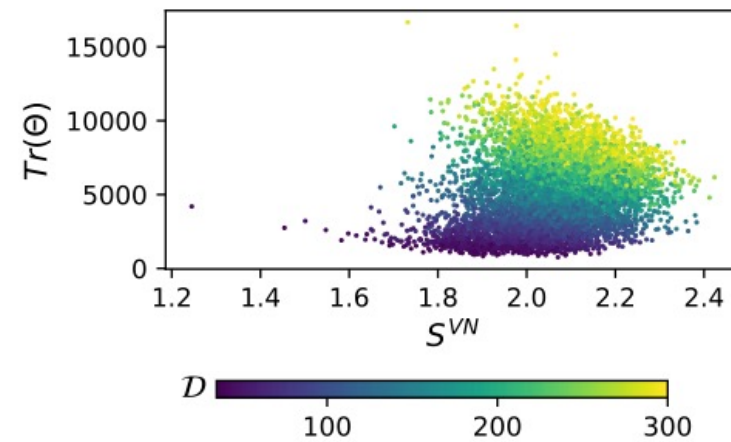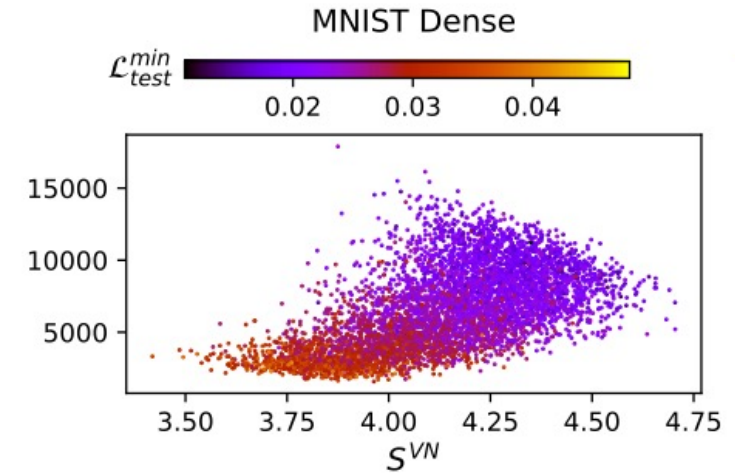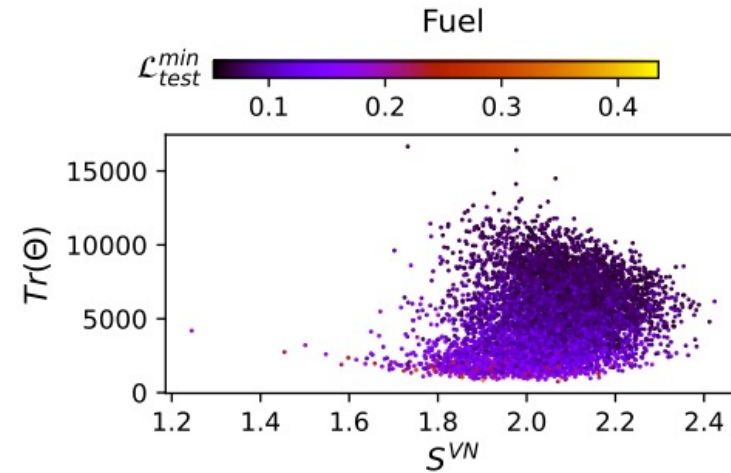
$$\Theta = V \Lambda V^T$$

Eigenvalues of the NTK matrix

$$S = - \sum_i \lambda_i \cdot \ln \lambda_i$$

Measure of correlation in data

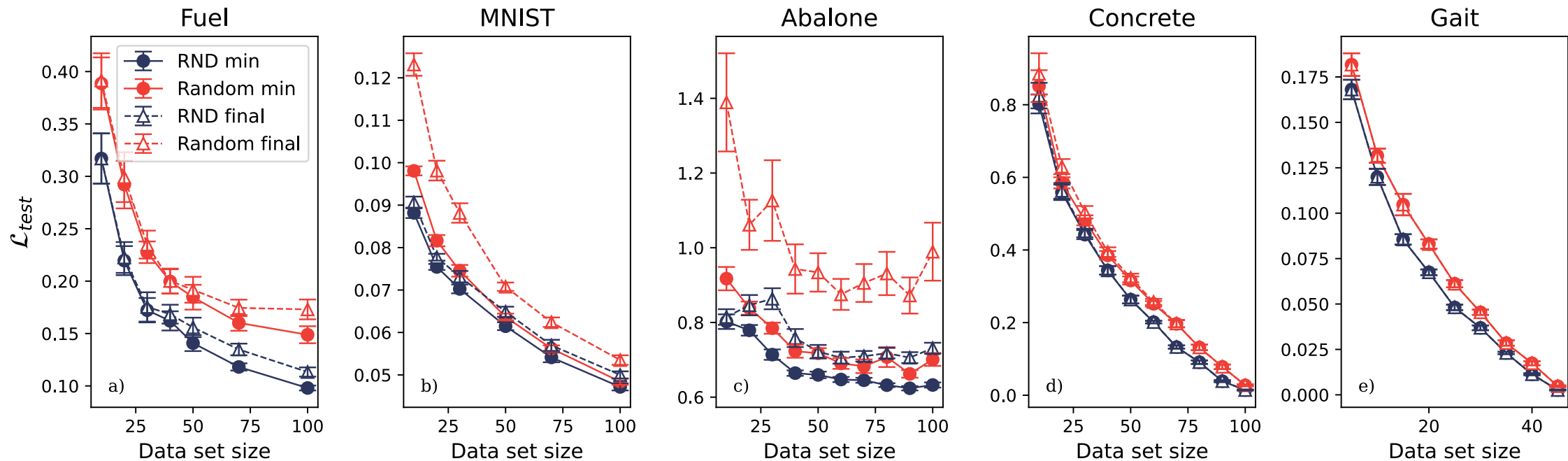$$Tr(\Theta) = \sum_i \Theta_{ii}$$

Weighting of largest step direction

9/11/2023

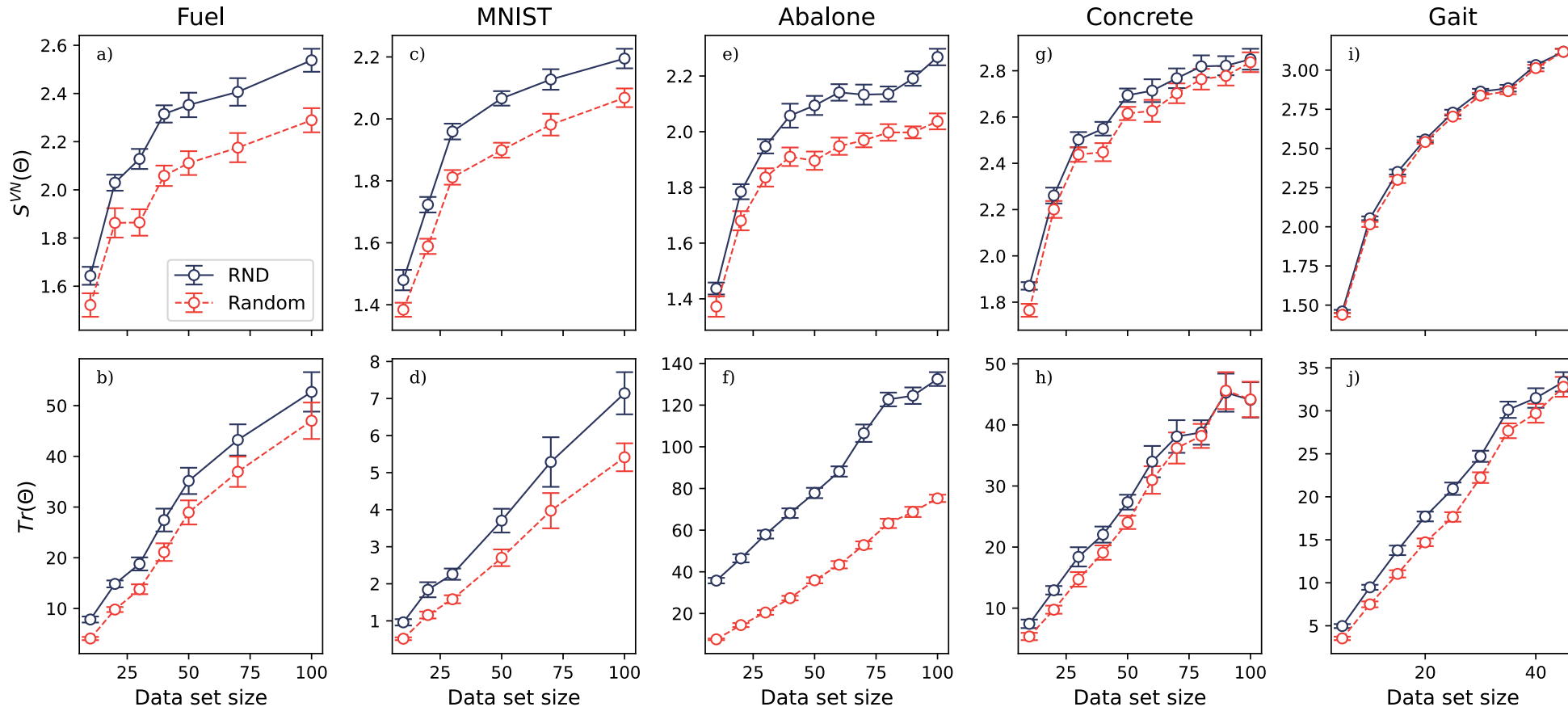| | |
|---|---|
| **Datasets** | 5 |
| **Optimizer** | ADAM(0.001) |
| **Loss Function** | MSE or CE |
| **Architectures** | 5 |
| **Epochs** | 200 |

# Data

Levelling the playing field

# Data Selection in Neural Networks
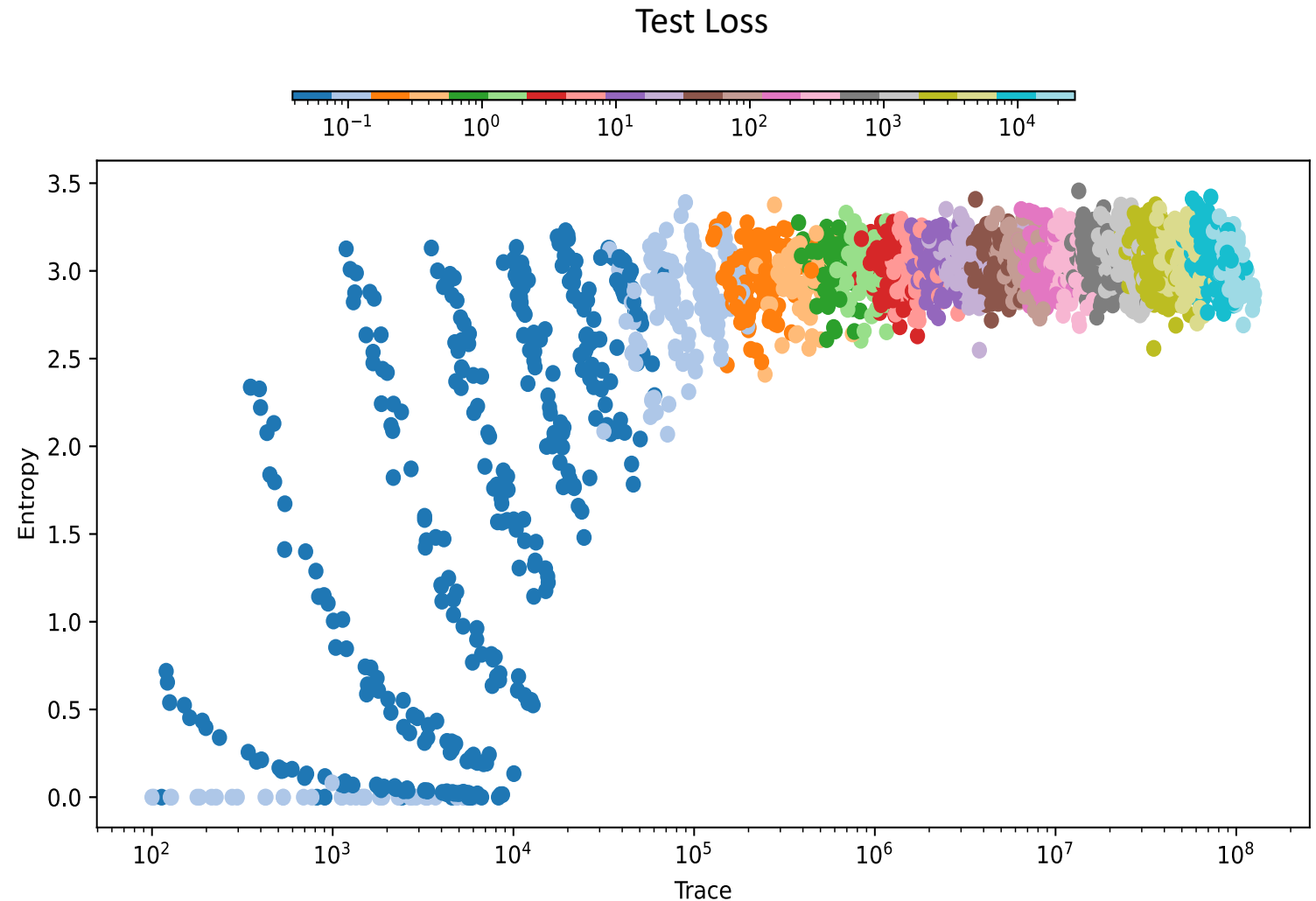


RND selected data-sets outperform random selection

9/11/2023

# What do the collective variables say?
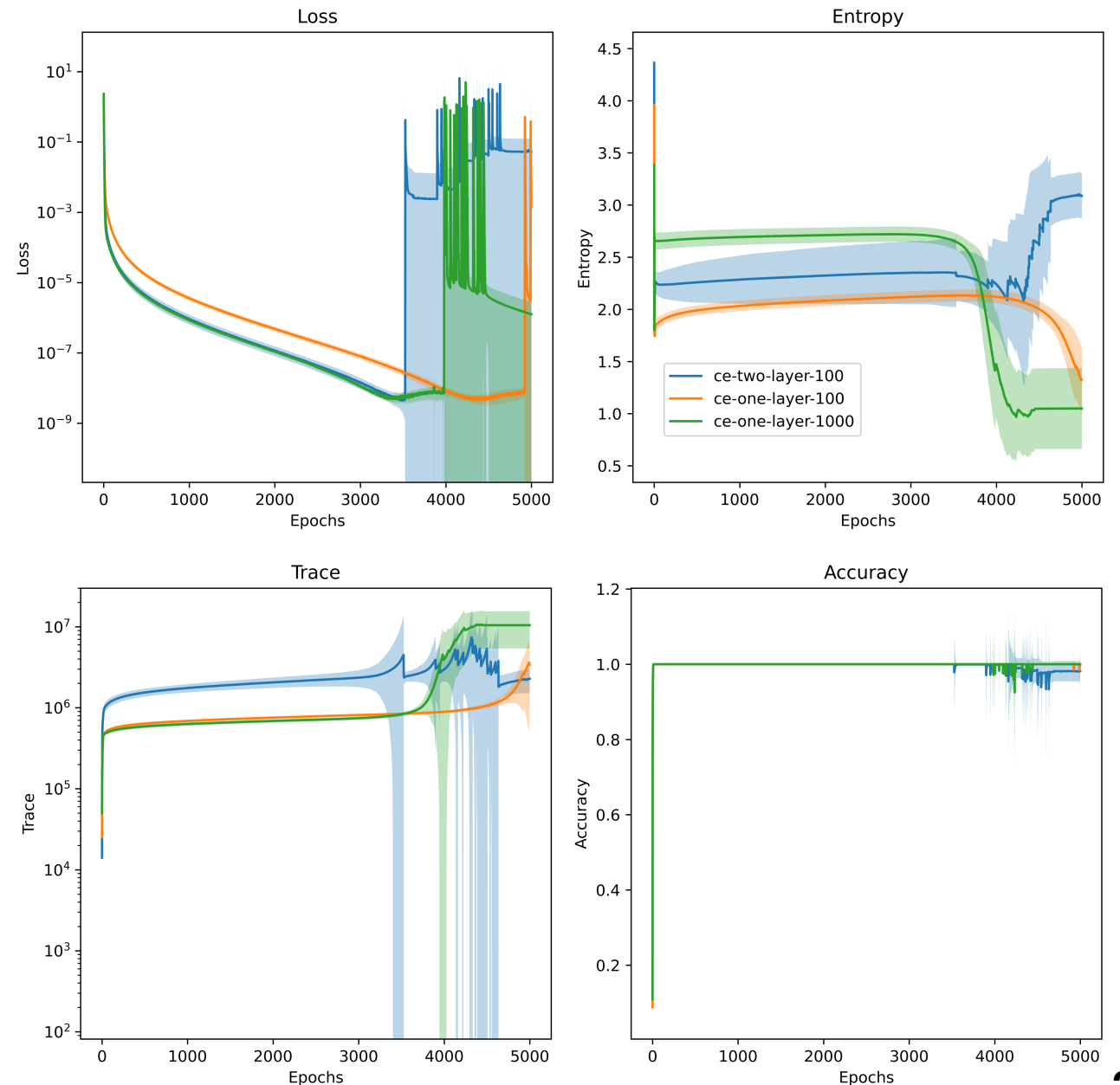


Datasets with larger trace / entropy perform better!

9/11/2023

# Next Steps: Initialization

| | |
|---|---|
| **Dataset** | MNIST (1000) |
| **Optimizer** | ADAM(0.001) |
| **Architecture** | $\mathcal{D}^{128}r\mathcal{D}^{128}r\mathcal{D}^{10}$ |
| **Weight std** | 0.0 – 1.0 |
| **Bias std** | 0.0 - 1.0 |



Test Loss

# Next Steps: Dynamics

- Compute CVs at all epochs

- Search for universal behaviour

- Interesting long-time behaviour

| Dataset | MNIST (1000) |
|---|---|
| Loss Function | Cross-entropy |
| Optimizer | SGD(0.01) |
| Architectures | 3 |



9/11/2023

ML4Jets, DESY, Hamburg -- stovey@icp.uni-stuttgart.de

Towards a phenomenological understanding of neural networks: data

Samuel Tovey[3,4,1] (iD), Sven Krippendorf[3,4,2] (iD), Konstantin Nikolaou[3,1] (iD) and Christian Holm[1] (iD)

📄 Article PDF

zincware/**ZnNL**

Python package to perform random network distillation.

👥 4 Contributors    ⊙ 16 Issues    ☆ 5 Stars    ⑂ 0 Forks

https://github.com/zincware/ZnNL

# Wrapping Up: ZnNL

- Tools of physics can help us understand neural networks
  - Statistical Physics
  - Quantum Mechanics

- We can leverage this understanding
  - Data Selection
  - Initialization
  - Optimization and dynamics

DALL . E 3