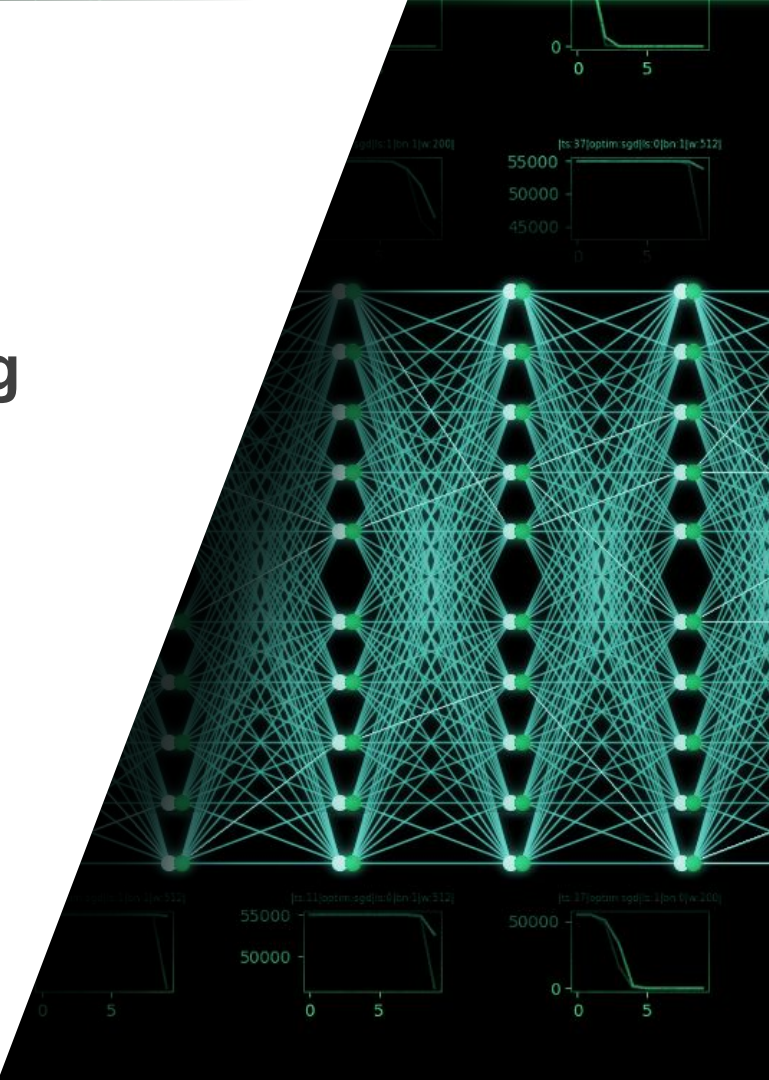# Exploring neural network training dynamics through binary node activations

*Gerbrand Haasbroek & Marelie H Davel*

Multilingual Speech Technologies, North-West University, South Africa
Centre for Artificial Intelligence Research

# Background

- No theoretical framework for DNN training or generalisation

- Generalisation studied using

    - Geometry of loss landscape

    - Stability and robustness

    - Complexity of hypothesis space

    - Margin distributions

NWU®

# Background

- *Sample Set* - Samples that activate a node

- Sample sets formed during the forward pass

- Sample sets refined during the backward pass

- Evolution of sample sets has not been explored

NWU®

# Contributions

- Motivation of the importance of sample sets

- Demonstration of a consistent evolution of sample sets

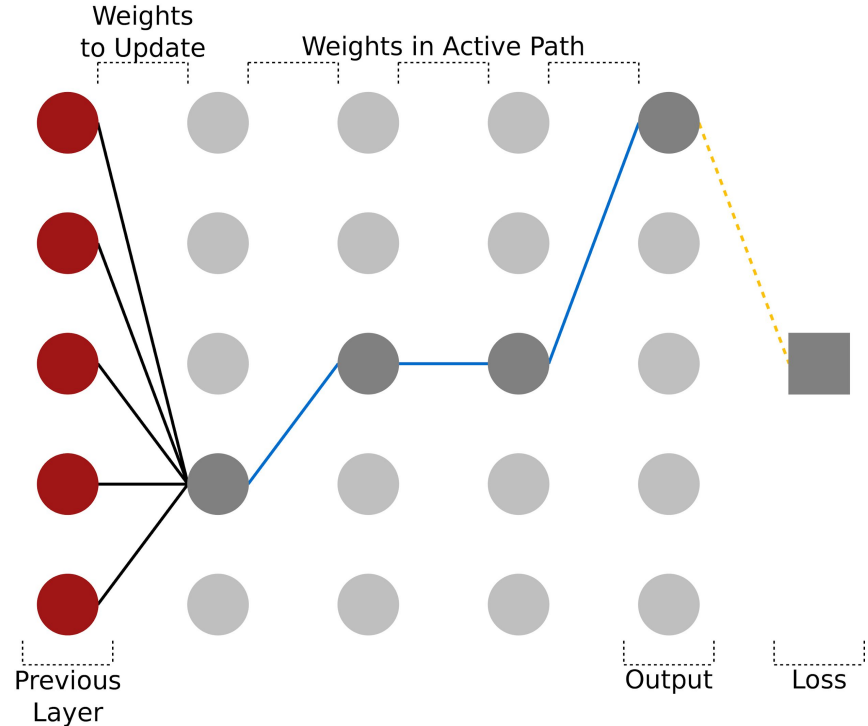- Interpretation of findings in terms of clustering

# Sample Sets

- Sample added to set if

$$\mathbf{z}^{\mathbf{t}}_{l-1} \cdot \mathbf{w}_{l,j} > \eta \sum_{\mathbf{s} \in \hat{S}_{b,l,j}} \left( \mathbf{z}^{\mathbf{t}}_{l-1} \cdot \mathbf{z}^{\mathbf{s}}_{l-1} \right) \phi^{\mathbf{s}}_{l,j}$$

- Sample removed otherwise

- Node-supported cost

$$\phi^{\mathbf{s}}_{l,j} = \sum_{p \in P^{\mathbf{s}}_j} \left( \lambda^{\mathbf{s}}_p \prod_{k=1}^{N-l} w_{p_k} \right)$$



Weights to Update

Weights in Active Path

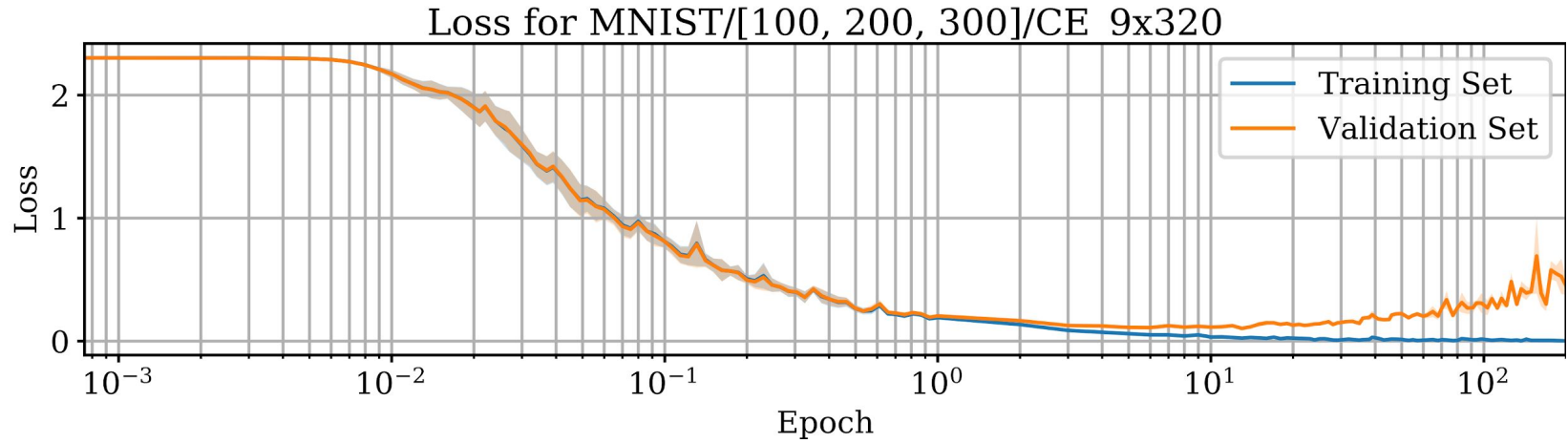Previous Layer

Output

Loss

# Experiments

- Train fully connected networks:

    - Number of layers: 2, 5, 9

    - Number of nodes per layer: 20, 80, 320

    - Activation functions: ReLU

    - Loss functions: mean squared error, cross entropy

    - Adam optimiser

    - MNIST dataset
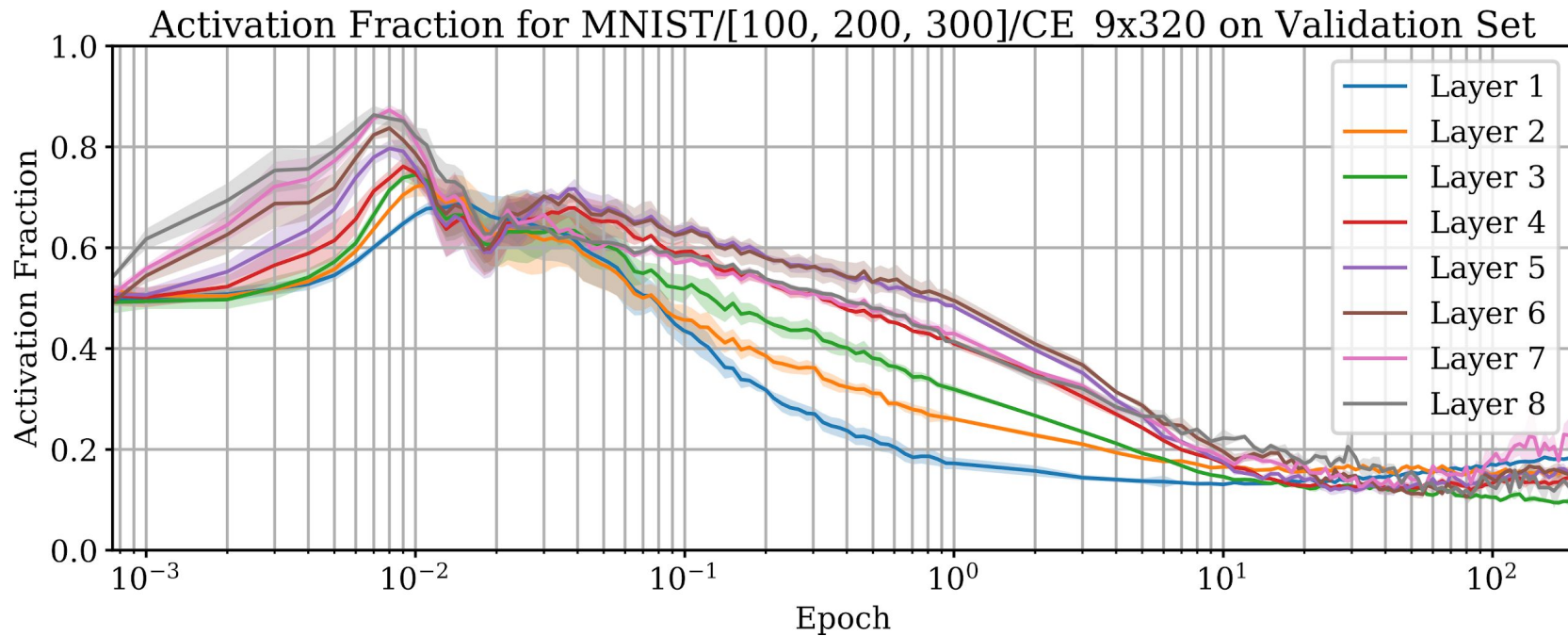
    - Epochs: 200

NWU®

# Measurements

- *Activation fraction* - Fraction of all samples in sample set

- *Predictability* - Average amount of information about binary activation given class

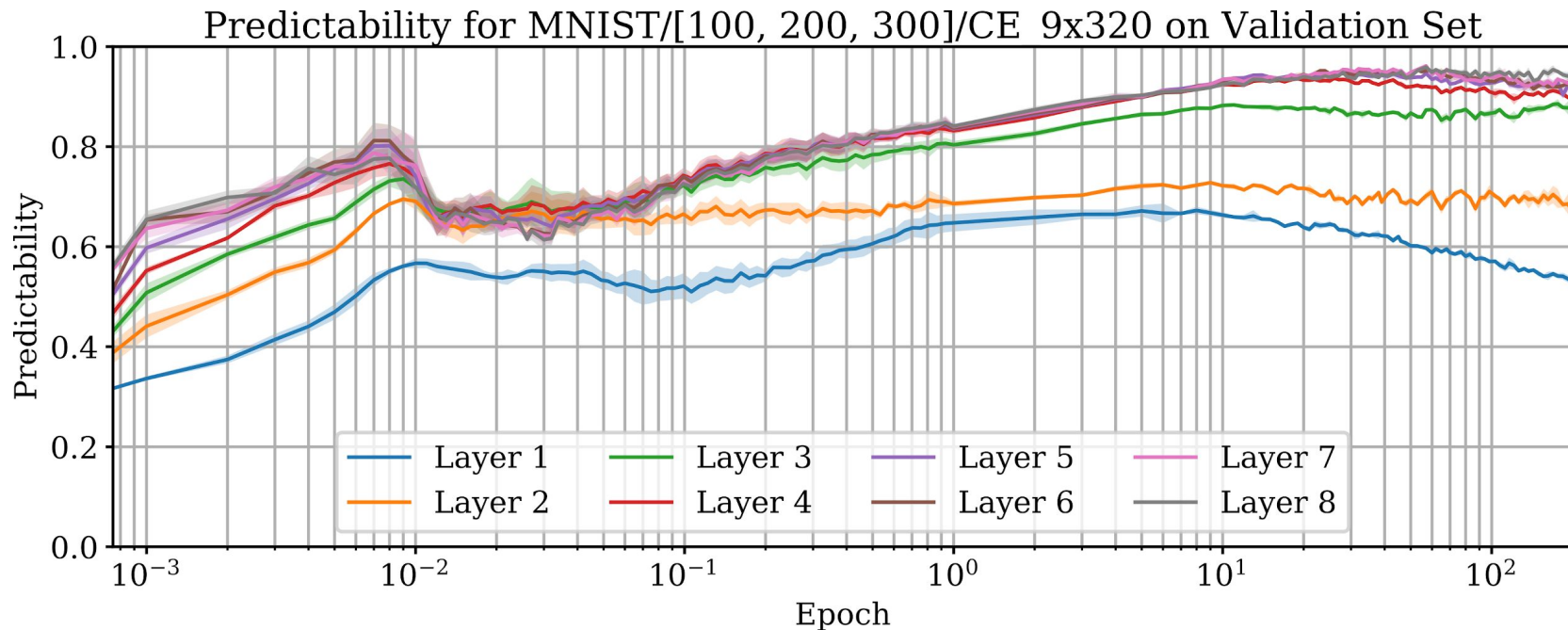- *Informativeness* - Average amount of information about class given binary activation

NWU®

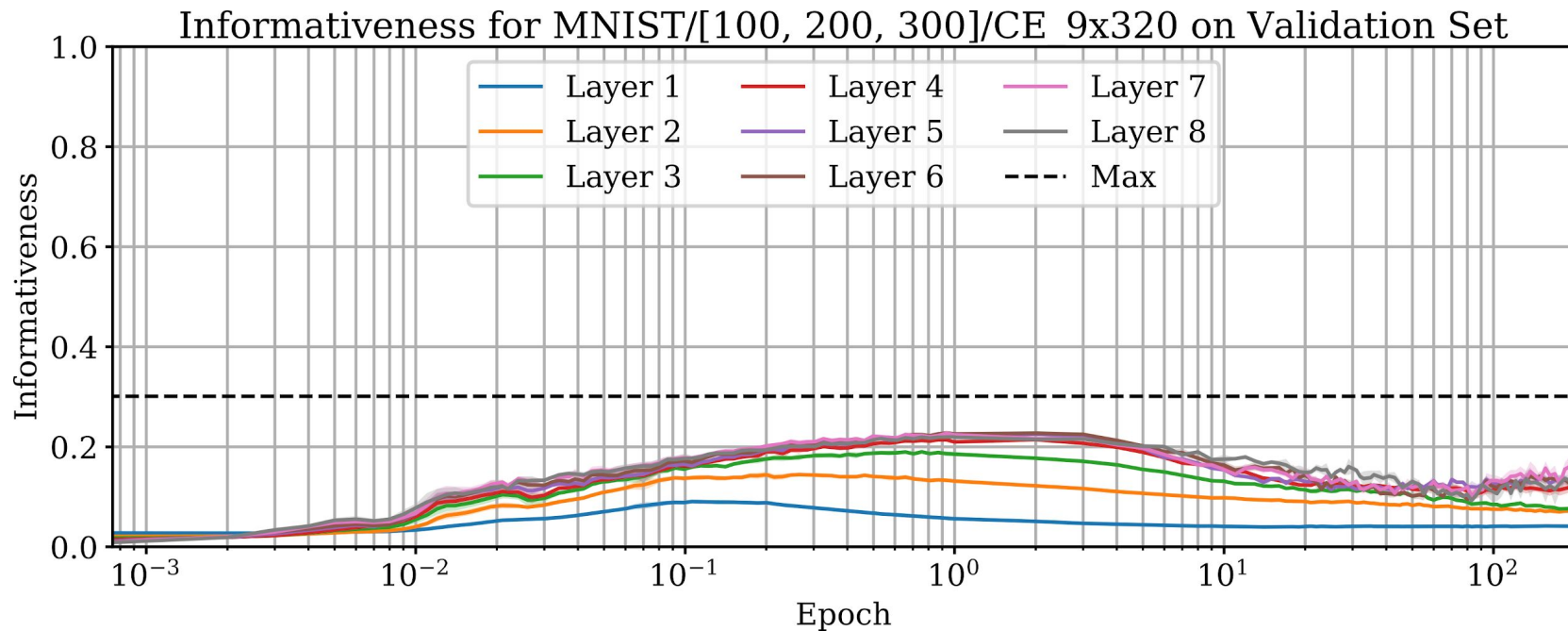# Results



Loss for MNIST/[100, 200, 300]/CE 9x320

Activation Fraction for MNIST/[100, 200, 300]/CE 9x320 on Validation Set

# Results



Predictability for MNIST/[100, 200, 300]/CE 9x320 on Validation Set

# Results



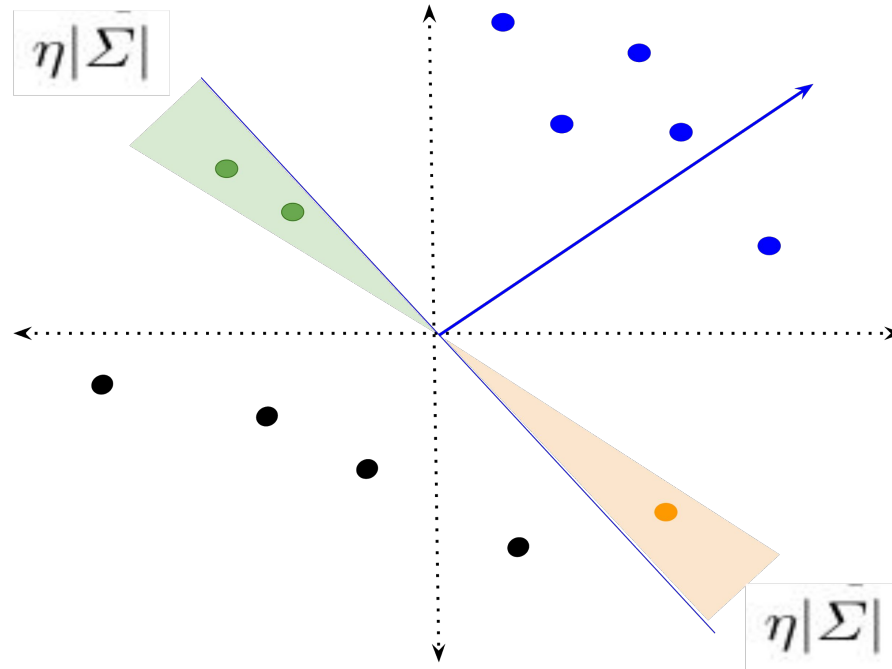Informativeness for MNIST/[100, 200, 300]/CE 9x320 on Validation Set

# Interpretation

- Initial loss magnitude is large

- Sample sets grow until loss starts decreasing

- Sample sets shrink as loss decreases

- Samples with zero loss do not contribute to weight updates

- Training finds clusters of samples that cause coherent behaviour

NWU ®

# Questions

NWU®

# Appendix

- Predictability

$$p_{l,i} \approx 1 - \frac{1}{\log{(2)}} H\left(\hat{Z}_{l,i} \mid Y\right)$$

- Informativeness

$$u_{l,i} \approx 1 - \frac{1}{\log{(|C|)}} H\left(Y \mid \hat{Z}_{l,i}\right)$$

NWU®