

# Flattening the Path to Generalization in LLMs

Ting Han and Prof. Dr. Michael Kamp

Lamarr Institute, TU Dortmund and Institute of Artificial Intelligence in Medicine (IKIM)

Partner institutions:









Institutionally funded by:



Bundesministerium für Forschung, Technologie und Raumfahrt Ministerium für Kultur und Wissenschaft des Landes Nordrhein-Westfalen



# AI/LLMs are EVERYWHERE!





Machine Learning









AlphaFold, developed by DeepMind, has solved the 50-year-old protein folding problem with remarkable accuracy.



ChatGPT reached 1 million users in 5 days.



Machine Learning









AlphaFold, developed by DeepMind, has solved the 50-year-old protein folding problem with remarkable accuracy.



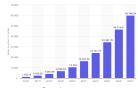
ChatGPT reached 1 million users in 5 days.



Al-powered personalization increases sales by 10-15%.



Machine Learning



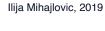
In 2022, the global AI market was valued at over \$100 billion and is projected to reach \$190 billion by 2025.

Statista, 2024











AlphaFold, developed by DeepMind, has solved the 50-year-old protein folding problem with remarkable accuracy.



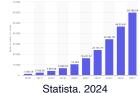
ChatGPT reached 1 million users in 5 days.



Al-powered personalization increases sales by 10-15%.



**Machine Learning** 



In 2022, the global AI market was valued at over \$100 billion and is projected to reach \$190 billion by 2025.











AlphaFold, developed by DeepMind, has solved the 50-year-old protein folding problem with remarkable accuracy.



ChatGPT reached 1 million users in 5 days.

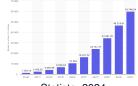


Al-powered personalization increases sales by 10-15%.



Machine Learning





In 2022, the global AI market was valued at over \$100 billion and is projected to reach \$190 billion by 2025.

Statista, 2024











AlphaFold, developed by DeepMind, has solved the 50-year-old protein folding problem with remarkable accuracy.



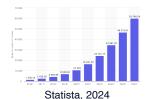
ChatGPT reached 1 million users in 5 days.



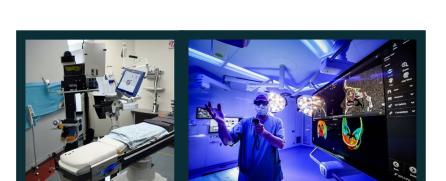
Al-powered personalization increases sales by 10-15%.



Machine Learning



In 2022, the global AI market was valued at over \$100 billion and is projected to reach \$190 billion by 2025.



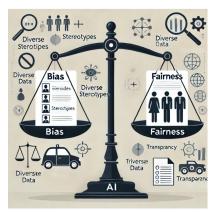


# Are AI/LLMs TRUELY everywhere? Not yet!





#### Bias & Fairness





#### Private data recovery



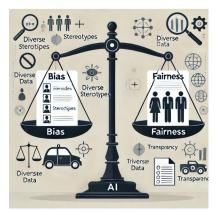
# Are AI/LLMs TRUELY everywhere?



Trustwhothiess & Responsibility



#### Bias & Fairness





#### Private data recovery



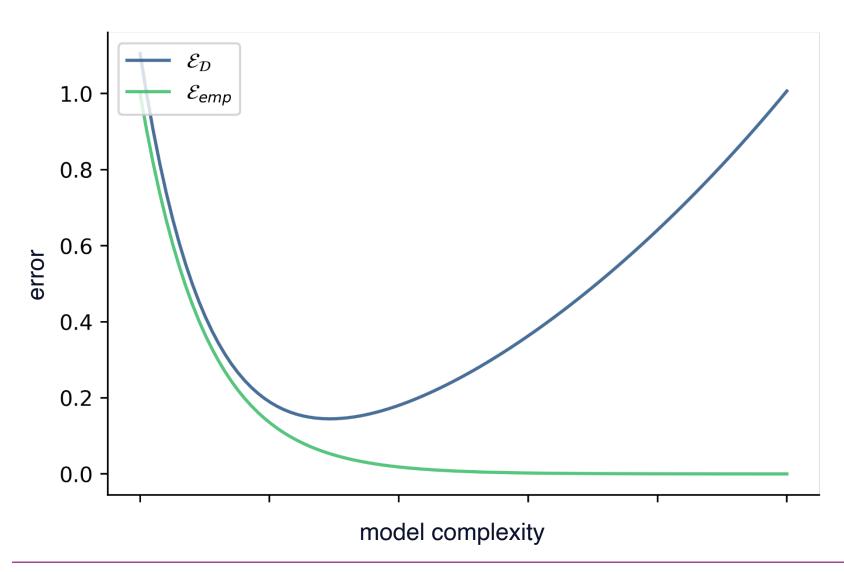




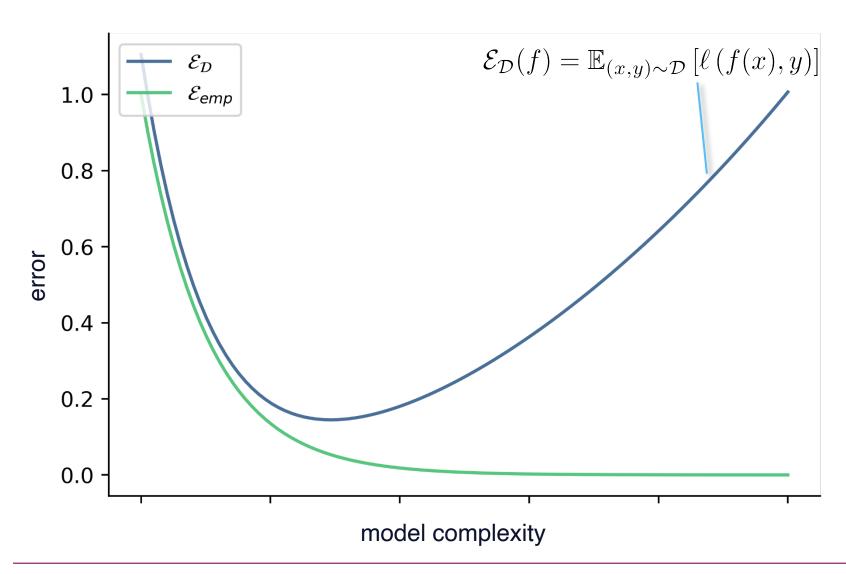


# Generalization

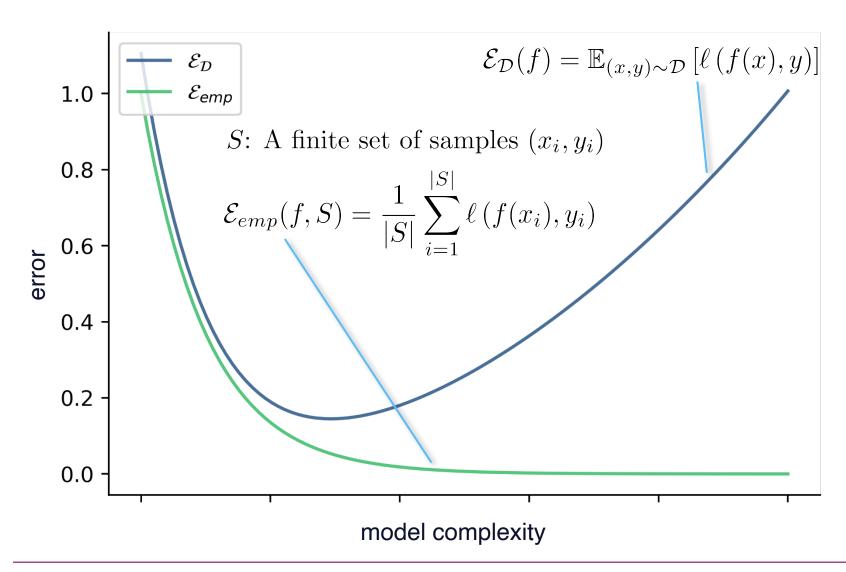




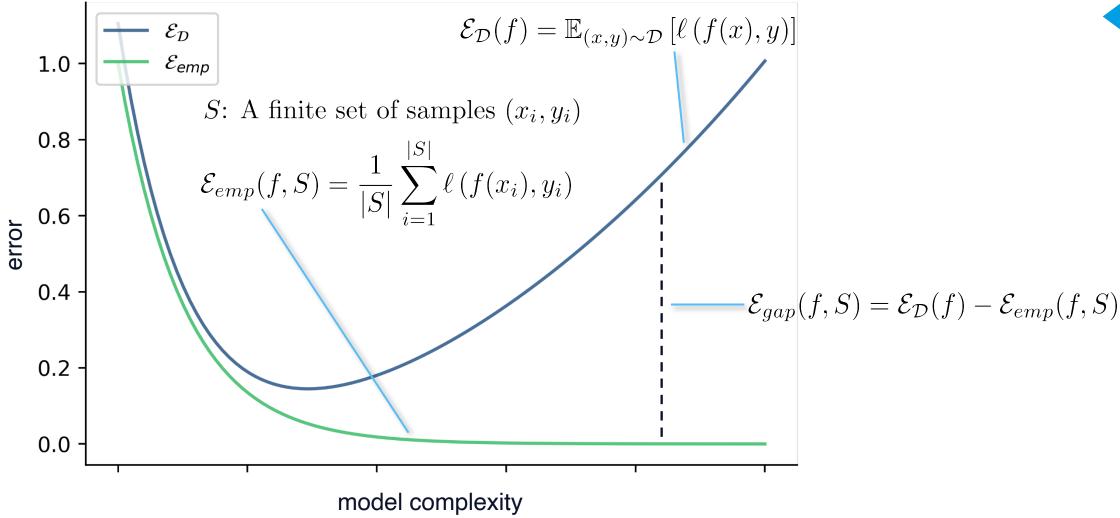














1.0
$$\mathcal{E}_{\mathcal{D}} \qquad \qquad \mathcal{E}_{\mathcal{D}} \qquad \qquad \mathcal{E}_{\mathcal{D}} \left[ \ell \left( f(x), y \right) \right] = \mathbb{E}_{x \sim \mathcal{D}_{x}} \left[ \ell \left( f(x), y \right) \right]$$

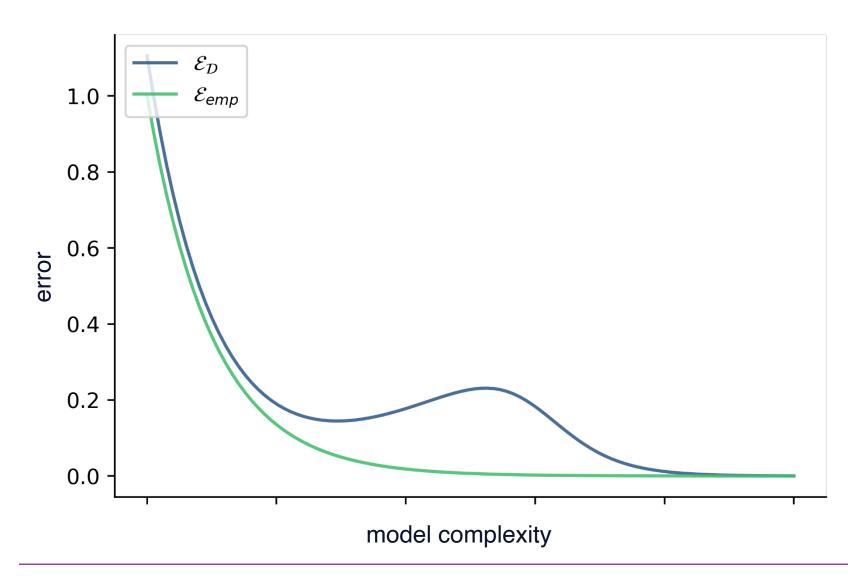
$$S: \text{ A finite set of samples } (x_{i}, y_{i})$$

# statistical learning theory:

Guarantee: with probability : 
$$\delta$$
  $\mathcal{E}_{\mathcal{D}} \leq \mathcal{E}_{emp} + \sqrt{\frac{1}{N} \left( \mathcal{C}(\mathcal{H}) \log \frac{1}{\delta} \right)}$   $\mathcal{C}(\mathcal{H})$ : complexity of model class  $\mathcal{H}$  model complexity

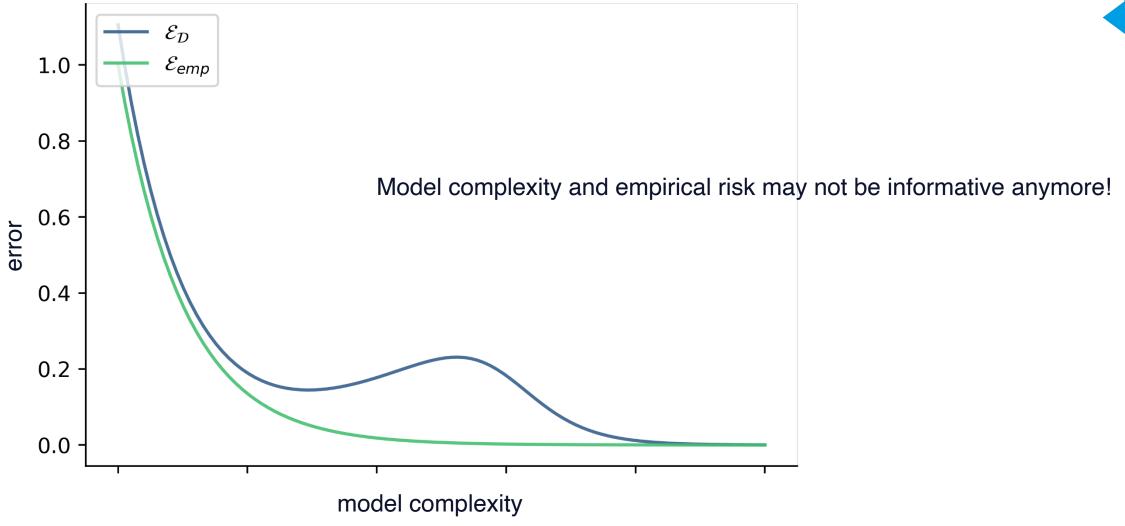
# **Double Descent**





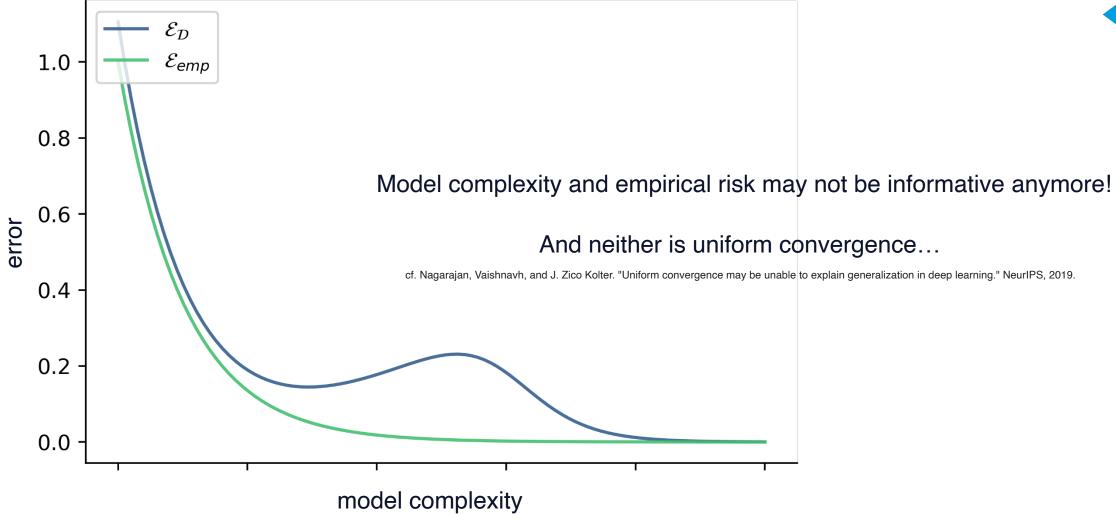
# **Double Descent**





### **Double Descent**







Feature Robustness and Relative Flatness

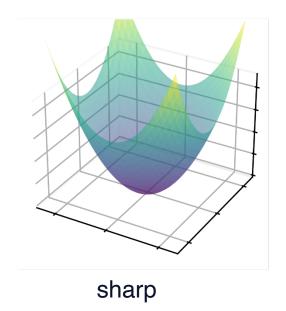
### Flatness of the Loss Curve

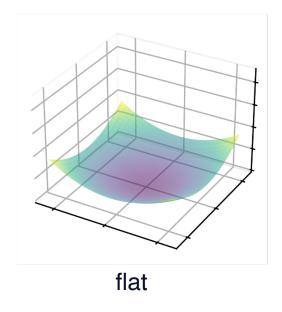


# "Large region in weight space with the property that each weight vector from that region leads to similar small error"

Hochreiter and Schmidhuber. "Simplifying neural nets by discovering flat minima." NIPS, 1995.

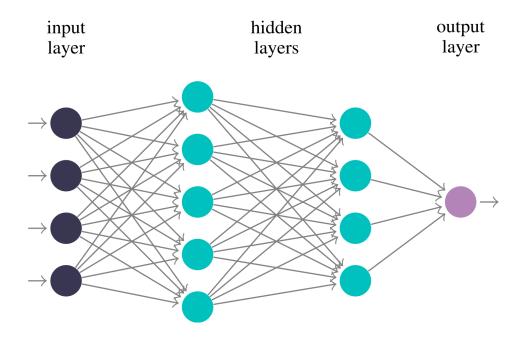
$$\mathcal{E}_{emp}(\mathbf{w} + \delta \mathbf{v}, S) \approx \mathcal{E}_{emp}(\mathbf{w}, S)$$





# Weights and Feature Space

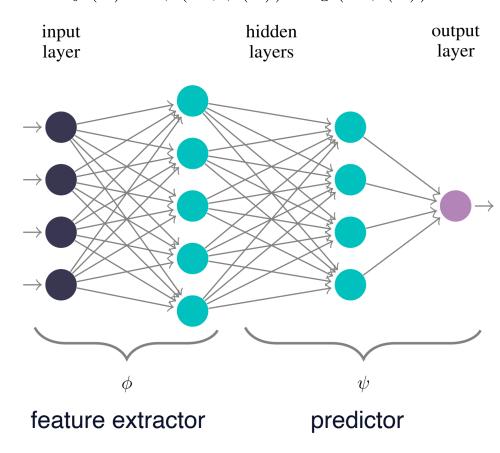




# Weights and Feature Space



$$f(x) = \psi(\mathbf{w}, \phi(x)) = g(\mathbf{w}\phi(x))$$





# key insight:

$$g((\mathbf{w} + \delta \mathbf{v})\phi(x))$$



# key insight:

$$g((\mathbf{w} + \delta \mathbf{v})\phi(x)) = g((\mathbf{w} + \delta \mathbf{w}A)\phi(x))$$

$$\phi(x) \in \mathbb{R}^m, \ A \in \mathbb{R}^{m \times m}$$



# key insight:

$$g((\mathbf{w} + \delta \mathbf{v})\phi(x)) = g((\mathbf{w} + \delta \mathbf{w}A)\phi(x))$$
$$= g(\mathbf{w}(\phi(x) + \delta A\phi(x)))$$
$$\phi(x) \in \mathbb{R}^m, \ A \in \mathbb{R}^{m \times m}$$

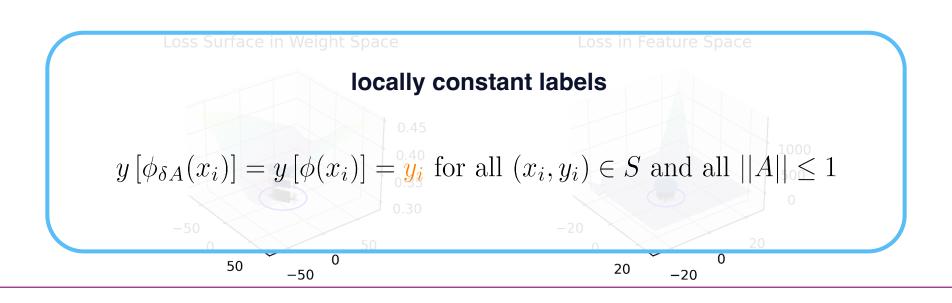


$$f(x) = \psi(\mathbf{w}, \phi(x)) = g(\mathbf{w}\phi(x))$$
$$g((\mathbf{w} + \delta \mathbf{w} A) \phi(x)) = g(\mathbf{w} (\phi(x) + \delta A \phi(x)))$$

$$\mathcal{E}_{emp}(\mathbf{w} + \delta \mathbf{w} A, \phi(S)) - \mathcal{E}_{emp}(\mathbf{w}, \phi(S)) = \frac{1}{|S|} \sum_{i=1}^{|S|} \ell(\psi(\mathbf{w}, \phi_{\delta A}(x_i)), y_i) - \ell(\psi(\mathbf{w}, \phi(x_i)), y_i)$$

flatness

feature robustness





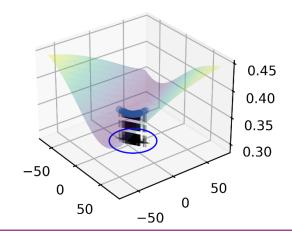
$$f(x) = \psi(\mathbf{w}, \phi(x)) = g(\mathbf{w}\phi(x))$$
$$g((\mathbf{w} + \delta \mathbf{w} A) \phi(x)) = g(\mathbf{w} (\phi(x) + \delta A \phi(x)))$$

$$\mathcal{E}_{emp}(\mathbf{w} + \delta \mathbf{w} A, \phi(S)) - \mathcal{E}_{emp}(\mathbf{w}, \phi(S)) = \frac{1}{|S|} \sum_{i=1}^{|S|} \ell(\psi(\mathbf{w}, \phi_{\delta A}(x_i)), y_i) - \ell(\psi(\mathbf{w}, \phi(x_i)), y_i)$$

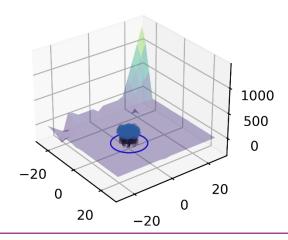
flatness

feature robustness

#### Loss Surface in Weight Space

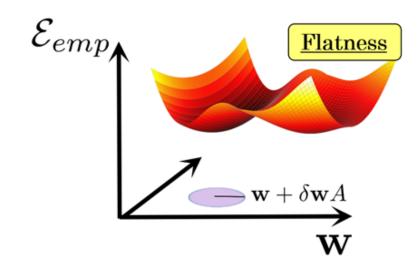


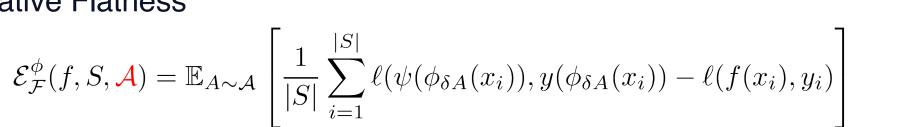
#### Loss in Feature Space



$$\mathcal{E}_{\mathcal{F}}^{\phi}(f, S, \mathcal{A}) = \mathbb{E}_{A \sim \mathcal{A}} \left[ \frac{1}{|S|} \sum_{i=1}^{|S|} \ell(\psi(\phi_{\delta A}(x_i)), y(\phi_{\delta A}(x_i)) - \ell(f(x_i), y_i)) \right]$$

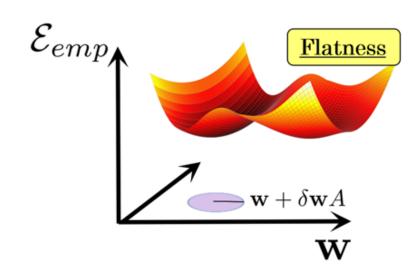


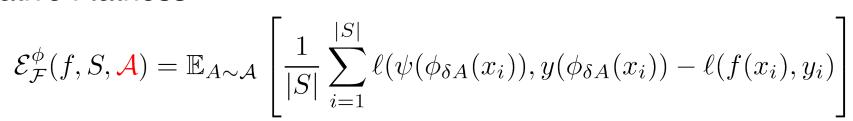






Second order **Taylor decomposition** of **feature robustness** and taking expectation over a sensibly chosen distribution on feature matrices results in a Hessian based flatness measure:

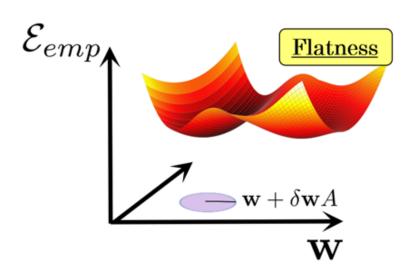






Second order **Taylor decomposition** of **feature robustness** and taking expectation over a sensibly chosen distribution on feature matrices results in a Hessian based flatness measure:

$$H_{s,s'}(\mathbf{w},\phi(S)) = \left[\frac{\partial^2 \mathcal{E}_{emp}(\mathbf{w},\phi(S))}{\partial w_{s,t} \partial w_{s',t'}}\right]_{1 \le t,t' \le m}$$
$$\mathbf{w}_s = (w_{s,t})_t \in \mathbb{R}^{1 \times m}$$



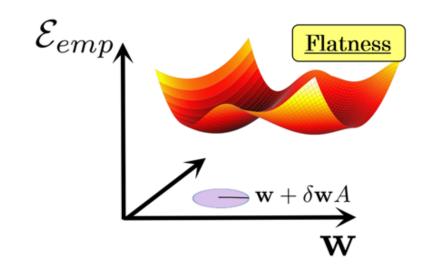
$$\mathcal{E}_{\mathcal{F}}^{\phi}(f, S, \mathcal{A}) = \mathbb{E}_{A \sim \mathcal{A}} \left[ \frac{1}{|S|} \sum_{i=1}^{|S|} \ell(\psi(\phi_{\delta A}(x_i)), y(\phi_{\delta A}(x_i)) - \ell(f(x_i), y_i)) \right]$$



Second order **Taylor decomposition** of **feature robustness** and taking expectation over a sensibly chosen distribution on feature matrices results in a Hessian based flatness measure:

$$H_{s,s'}(\mathbf{w},\phi(S)) = \left[\frac{\partial^2 \mathcal{E}_{emp}(\mathbf{w},\phi(S))}{\partial w_{s,t}\partial w_{s',t'}}\right]_{1 \le t,t' \le m}$$
$$\mathbf{w}_s = (w_{s,t})_t \in \mathbb{R}^{1 \times m}$$

$$K_{Tr}^{\phi}(\mathbf{w}) = \sum_{s,s'=1}^{d} \langle \mathbf{w}_s, \mathbf{w}_{s'} \rangle \cdot Tr(H_{s,s'}(\mathbf{w}, \phi(S)))$$



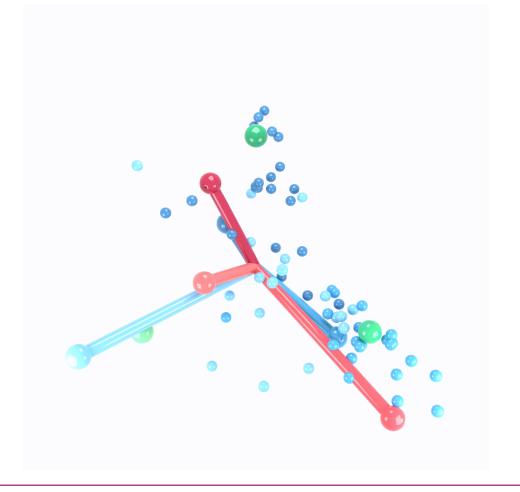




Neural Collapse is a phenomenon at the end of training where last-layer features and classifiers form a simple, symmetric structure.

- NC1: Variability Collapse: Features within each class collapse tightly to their class mean.
  - NC2: Convergence to Simplex ETF
  - NC3: Convergence to self-duality
  - NC4: Nearest Class-Center Simplification

$$NCC := \sum_{c \neq c'} \frac{V_c + V_{c'}}{2\|\mu_c - \mu_{c'}\|^2}$$





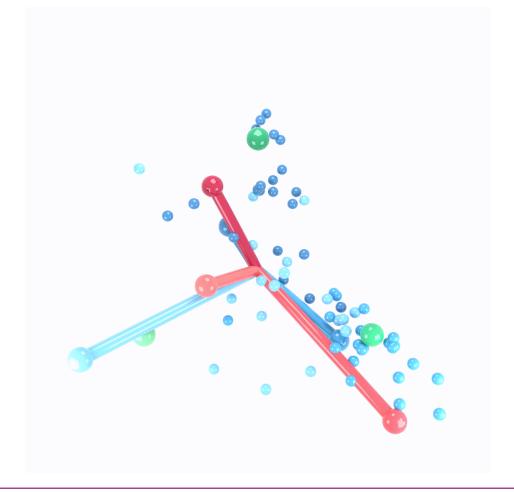
Neural Collapse is a phenomenon at the end of training where last-layer features and classifiers form a simple, symmetric structure.

- NC1: Variability Collapse: Features within each class collapse tightly to their class mean.
  - NC2: Convergence to Simplex ETF
  - NC3: Convergence to self-duality
  - NC4: Nearest Class-Center Simplification

$$NCC := \sum_{c \neq c'} \frac{V_c + V_{c'}}{2\|\mu_c - \mu_{c'}\|^2}$$

### Neural Collapse and Generalization





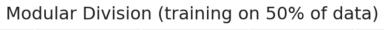


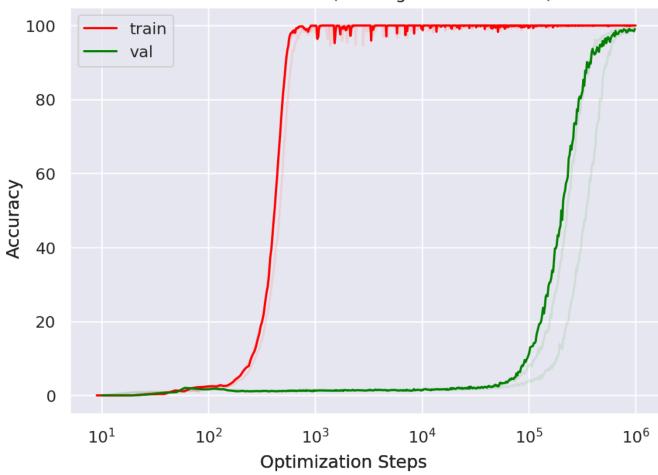
Grokking

# Grokking

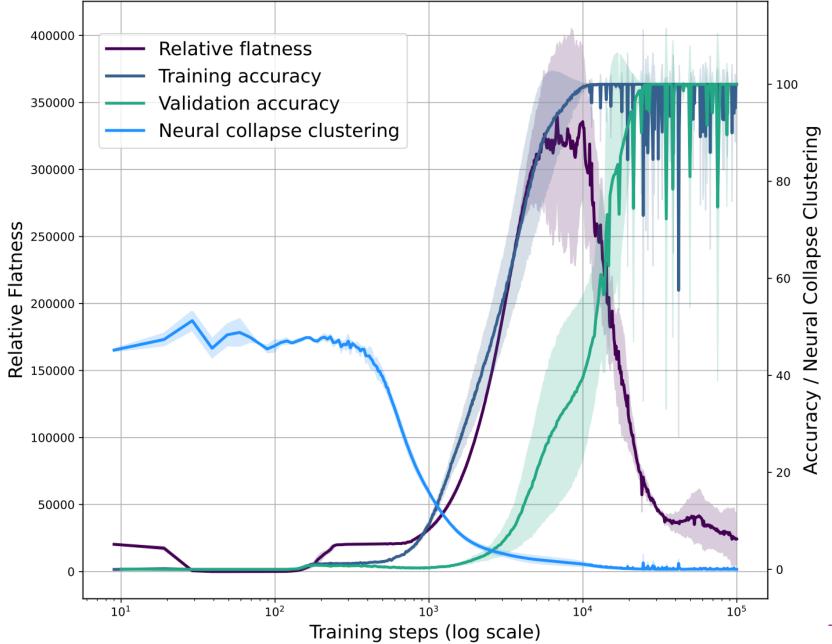










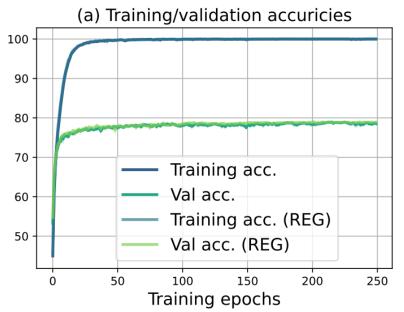


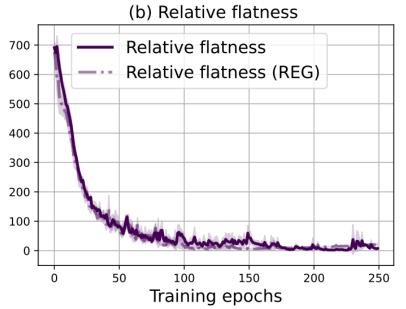


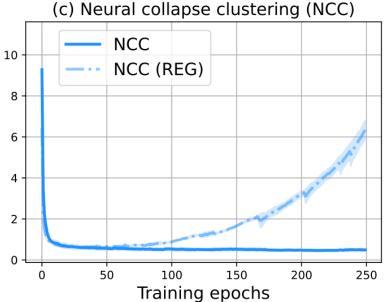


# Causal Interventions





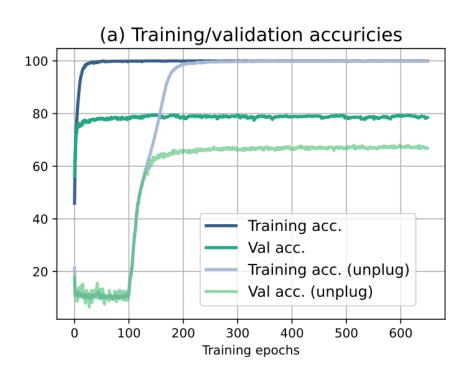


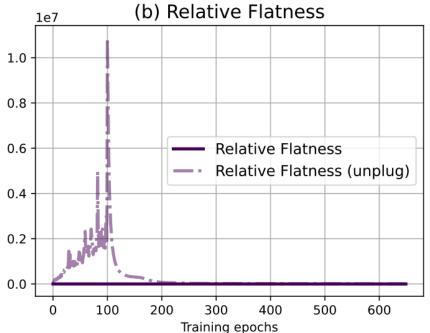


$$\mathcal{L}_{\text{NC\_REG}} = \mathcal{L}_{\text{CE}} - \lambda \cdot \text{NCC}$$

$$\text{NCC} := \sum_{c \neq c'} \frac{V_c + V_{c'}}{2\|\mu_c - \mu_{c'}\|^2}$$







$$\mathcal{L}_{ ext{NC\_RF}} = \mathcal{L}_{ ext{CE}} - \lambda * \kappa_{Tr}^{\phi}(\mathbf{w})$$

$$\kappa_{\mathrm{Tr}}^{\phi}(\mathbf{w}) := \sum_{s,s'=1}^{d} \langle \mathbf{w}_s, \mathbf{w}_{s'} \rangle \cdot \mathrm{Tr}(H_{s,s'}(\mathbf{w}, \phi(S)))$$



# Conclude:

Relative flatness, rather than Neural Collapse, is necessary to Generalization.

### Relative flatness + LLMs



- 1. LLM intervention:
  - Relative Flatness-Regularized Training, Placement, Diagnostic tool...
- 2. Cross lingual Generalization
  - -A, B and C?
- 3. Relative flatness + alignment
  - -D, E and F?



# Thanks and QA?

Partner institutions:









#### Institutionally funded by:



Bundesministerium für Forschung, Technologie und Raumfahrt Ministerium für Kultur und Wissenschaft des Landes Nordrhein-Westfalen

