



上海交通大學

SHANGHAI JIAO TONG UNIVERSITY



Towards the Difficulty for a Deep Neural Network to Learn Concepts of Different Complexities

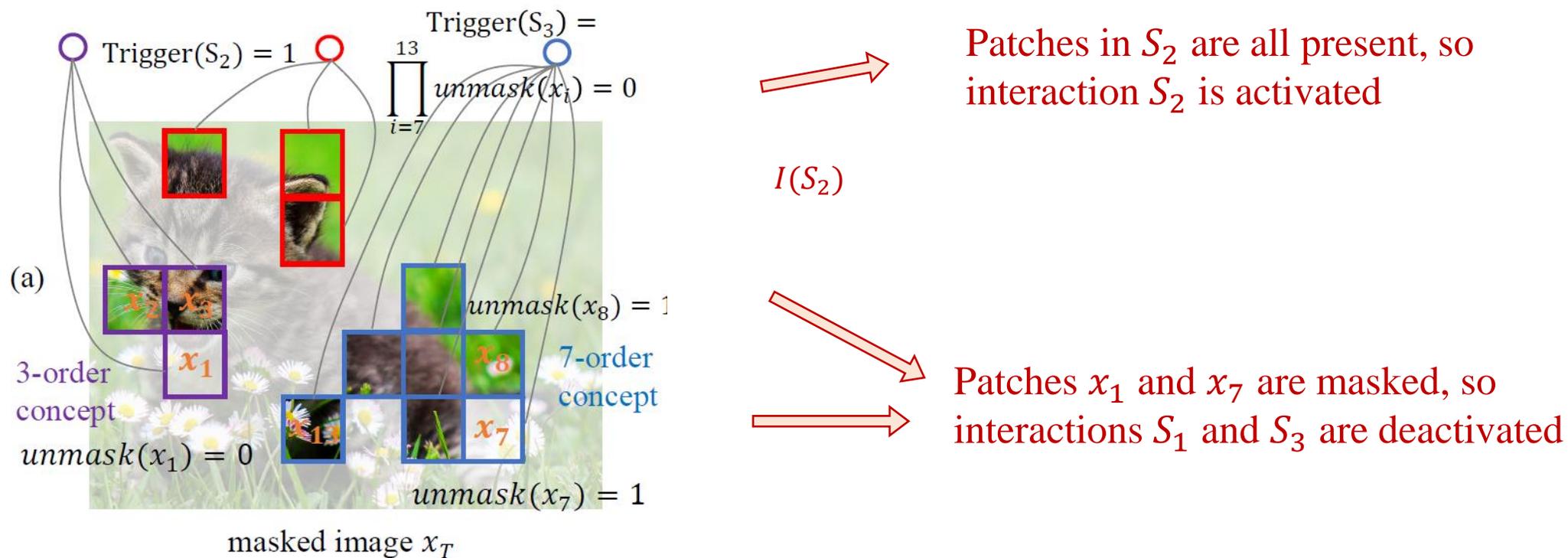
Dongrui Liu^{1*}, Huiqi Deng^{1*}, Xu Cheng¹, Qihan Ren¹, Kangrui Wang², Quanshi Zhang¹

1. Shanghai Jiao Tong University 2. University of Chicago

➤ Preliminary: understanding interactions as AND relationships

[1-3] used the Harsanyi interaction $I(S)$ to study the emergence of concepts in neural networks

Each interaction represents an **AND relationship** between input variables in a set S



[1] Ren et al. Defining and Quantifying the Emergence of Sparse Concepts in DNNs. CVPR, 2023.

[2] Ren et al. Can we faithfully represent absence states to compute shapley values on a DNN?. ICLR, 2022.

[3] Li et al. Does a neural network really encode symbolic concept. ICML, 2023.

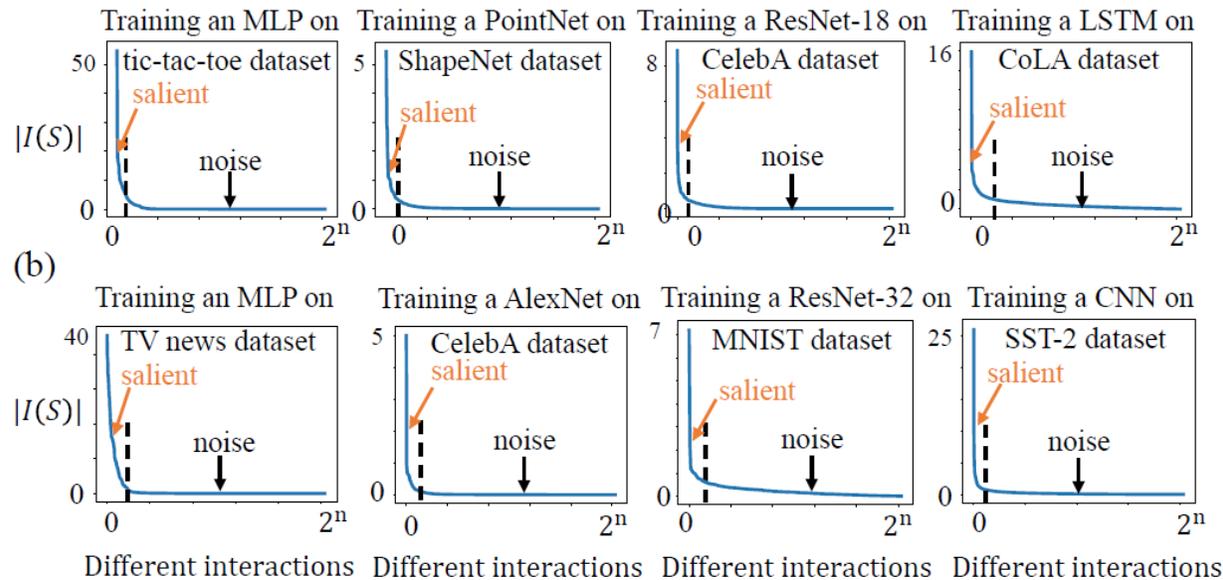


Preliminary: neural networks usually only encode a small number of salient concepts

- **Faithfulness:** The output of a neural network on any arbitrarily masked sample x_S can be disentangled into the sum of the effects of all interactions within set S

$$\forall S \subseteq N, v(x_S) = \sum_{T \subseteq S} I(T|x)$$

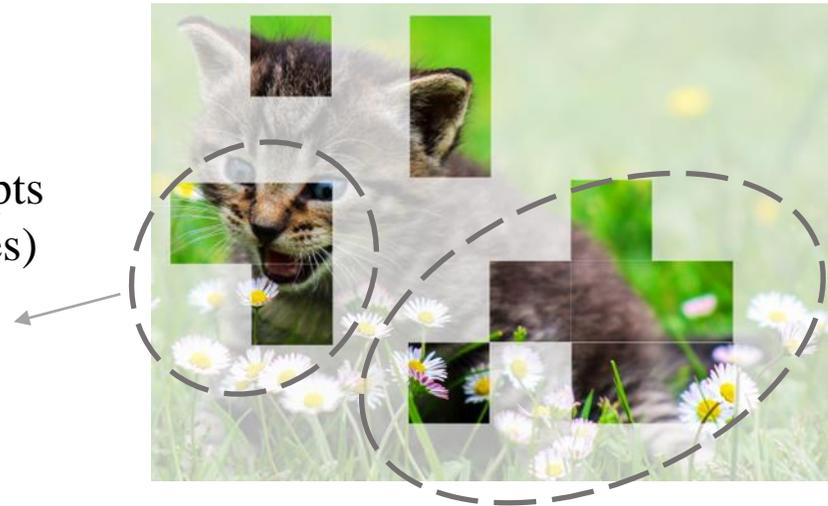
- **Sparsity:** **most** interactions have **near-zero effects** (noisy patterns), but only **a small number** of interactions have **significant effects** (salient interactive concepts).



Overview

We theoretically explain the reason why it is easier for DNNs to learn simple concepts than complex concepts.

Simple (low-order) concepts
(composed of a **few** patches)



Complex (high-order) concepts
(composed of **massive** patches)

- We show that low-order interactive concepts in the data are much **more stable than** high-order concepts, which makes low-order interactive concepts more likely to be encoded.
- We provide new insights into several empirical findings *w.r.t.* the training of DNNs.

Low-order interactive concepts in data are more stable

- Using input perturbations to roughly represent inevitable variations in data

There are some inevitable variations in data. For example, image classification suffers from different variations, such as small shape deformation and small object rotations.

In this paper, we analyze **the impact of data variations on interactive concepts of different complexities**. Such data variations are quite difficult to formulate, so we use a Gaussian perturbation as a rough representation.

➤ Low-order interactive concepts in data are more stable

The Harsanyi interaction effect of each concept can be rewritten.

Theorem 2 (proven in the supplementary material). *Given a neural network v and an arbitrary input sample $x' \in \mathbb{R}^n$, the network output can be decomposed by using the Taylor expansion i.e., $v(x') = \sum_{S \subseteq N} \sum_{\pi \in Q_S} U_{S,\pi} \cdot J(S, \pi | x')$. In this way, according to Equation (4), the Harsanyi interaction effect $I(S|x')$ on the sample x' can be reformulated as follows.*

$$I(S|x') = \sum_{\pi \in Q_S} U_{S,\pi} \cdot J(S, \pi | x'). \quad (4)$$

Here, $J(S, \pi | x') = \prod_{i \in S} \left(\text{sign}(x'_i - b_i) \cdot \frac{x'_i - b_i}{\tau} \right)^{\pi_i}$ denotes a Taylor expansion term of the degree π , where the degree $\pi \in Q_S = \{[\pi_1, \dots, \pi_n] | \forall i \in S, \pi_i \in \mathbb{N}^+; \forall i \notin S, \pi_i = 0\}$ and b_i is the baseline value to mask the input variable x_i . In addition, $U_{S,\pi} = \frac{\tau^m}{\prod_{i=1}^n \pi_i!} \frac{\partial^m v(x'_0)}{\partial x_1^{\pi_1} \dots \partial x_n^{\pi_n}} \cdot \prod_{i \in S} [\text{sign}(x'_i - b_i)]^{\pi_i}$, where x'_0 denotes the sample whose input variables are all masked. $m = \sum_{i=1}^n \pi_i$.

The stability of interactive concepts decreases along with the order of interactive concepts.

Theorem 3. *Let us add a Gaussian perturbation $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ to the input sample x . Let us first consider the case with the lowest degree $\hat{\pi} = [\hat{\pi}_1, \dots, \hat{\pi}_n] \in Q_S$, satisfying that $\forall i \in S, \hat{\pi}_i = 1; \forall i \notin S, \hat{\pi}_i = 0$. The mean and variance of $J(S, \hat{\pi} | x + \epsilon)$ over the Gaussian perturbation ϵ are given as*

$$\mathbb{E}_\epsilon[J(S, \hat{\pi} | x + \epsilon)] = 1, \quad \text{Var}_\epsilon[J(S, \hat{\pi} | x + \epsilon)] = \left(1 + \left(\frac{\sigma}{\tau} \right)^2 \right)^{|S|} - 1. \quad (5)$$

Furthermore, for the more general case with an arbitrary degree $\pi \in Q_S = \{[\pi_1, \dots, \pi_n] | \forall i \in S, \pi_i \in \mathbb{N}^+; \forall i \notin S, \pi_i = 0\}$, the mean and variance of $J(S, \pi | x + \epsilon)$ are computed as

$$\mathbb{E}_\epsilon[J(S, \pi | x + \epsilon)] = \mathbb{E}_\epsilon\left[\prod_{i \in S} \left(1 + \frac{\epsilon_i}{\tau} \right)^{\pi_i}\right], \quad \text{Var}_\epsilon[J(S, \pi | x + \epsilon)] = \text{Var}_\epsilon\left[\prod_{i \in S} \left(1 + \frac{\epsilon_i}{\tau} \right)^{\pi_i}\right]. \quad (6)$$

DNNs mainly learn low-order interactive concepts

The following equation enables us to understand a DNN for the classification task as a pseudo-linear function.

$$v(\mathbf{x}_T) = Y(\mathbf{x}_T) = \sum_{S \in \Omega} U_S \cdot C_S(\mathbf{x}_T)$$

- If an interactive concept S has a stable value (i.e., C_S stably being present/absent) across all samples in the same category, then we consider this concept is discriminative and easy to learn.
- If the concept cannot be consistently present or consistently absent over samples in the same category, then this concept is hard to learn.

DNNs mainly learn low-order interactive concepts

Experiment 1: verifying the claim that high-order interactive concepts are more sensitive to data variations than low-order interactive concepts. Therefore, we use the following two metrics to evaluate the discrimination power of each interactive concept S

Relative consistency of the interactive concept S

$$\beta(S) = \mathbb{E}_c [|\mathbb{E}_{\mathbf{x} \in \mathbf{X}_c} [I(S|\mathbf{x})]| / \text{Std}_{\mathbf{x} \in \mathbf{X}_c} [I(S|\mathbf{x})]]$$

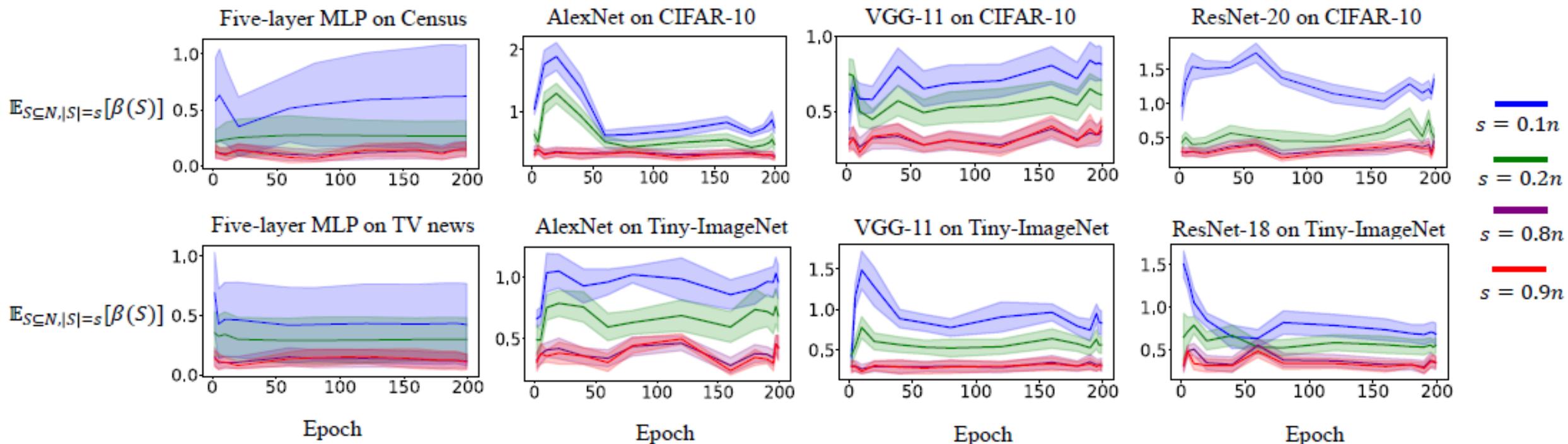
Relative instability of the interactive concept S

$$\kappa(S) = \mathbb{E}_{\mathbf{x} \in \mathbf{X}} [\mathbb{E}_{\epsilon} [||I(S|\mathbf{x} + \epsilon) - I(S|\mathbf{x})||]] / \mathbb{E}_{\mathbf{x} \in \mathbf{X}} [||I(S|\mathbf{x})||]$$

➤ DNNs mainly learn low-order interactive concepts

Experiment 1: verifying the claim that high-order interactive concepts are more sensitive to data variations than low-order interaction.

Consistency of the interactive concept S $\beta(S) = \mathbb{E}_c [|\mathbb{E}_{\mathbf{x} \in \mathcal{X}_c} [I(S|\mathbf{x})]| / \text{Std}_{\mathbf{x} \in \mathcal{X}_c} [I(S|\mathbf{x})]]$

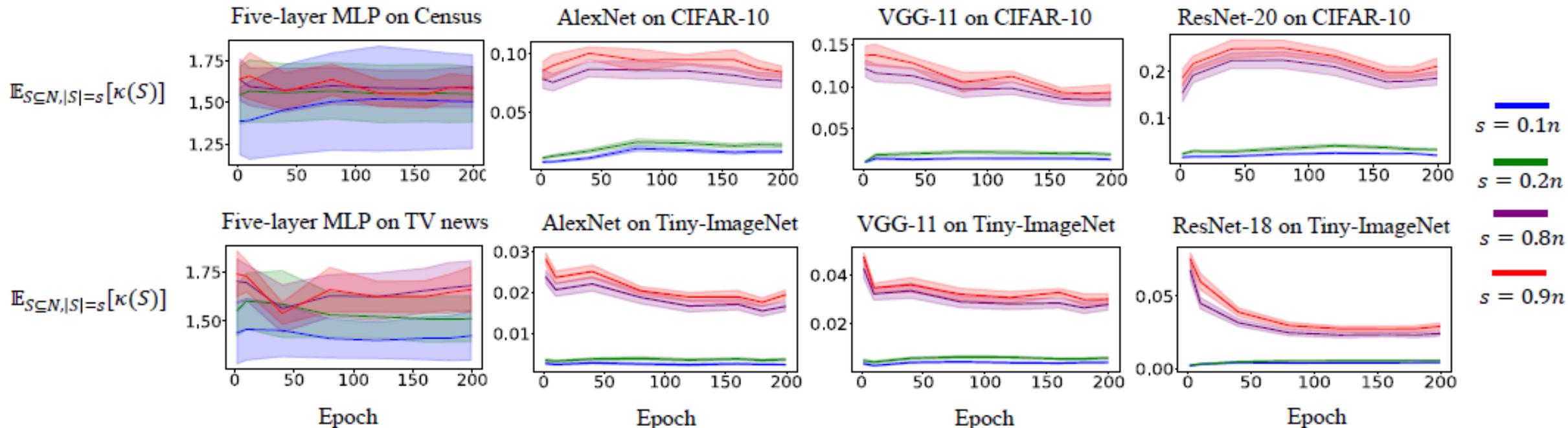


➤ DNNs mainly learn low-order interactive concepts

Experiment 1: verifying the claim that high-order interactive concepts are more sensitive to data variations than low-order interaction.

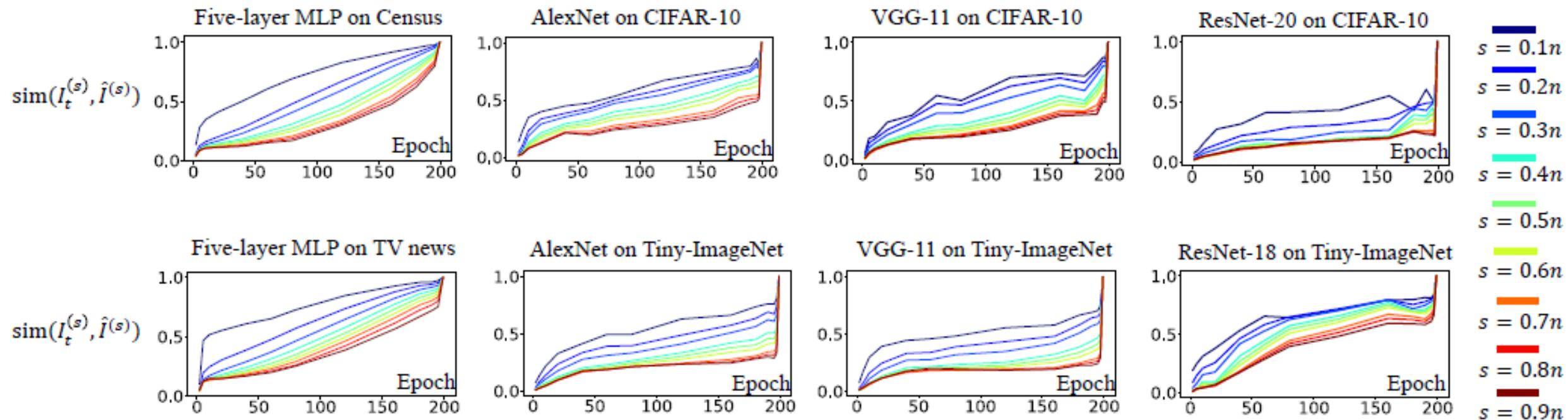
Instability of the interactive concept S

$$\kappa(S) = \mathbb{E}_{\mathbf{x} \in \mathbf{X}} [\mathbb{E}_{\epsilon} [|I(S|\mathbf{x} + \epsilon) - I(S|\mathbf{x})|]] / \mathbb{E}_{\mathbf{x} \in \mathbf{X}} [|I(S|\mathbf{x})|]$$



➤ DNNs mainly learn low-order interactive concepts

Experiment 2: verifying the phenomenon that low-order interactive concepts are usually learned faster than high-order concepts. Specifically, we examine the similarity between interactive concepts in the network v_t and the interactive concepts in the finally-learned DNN \hat{v} .

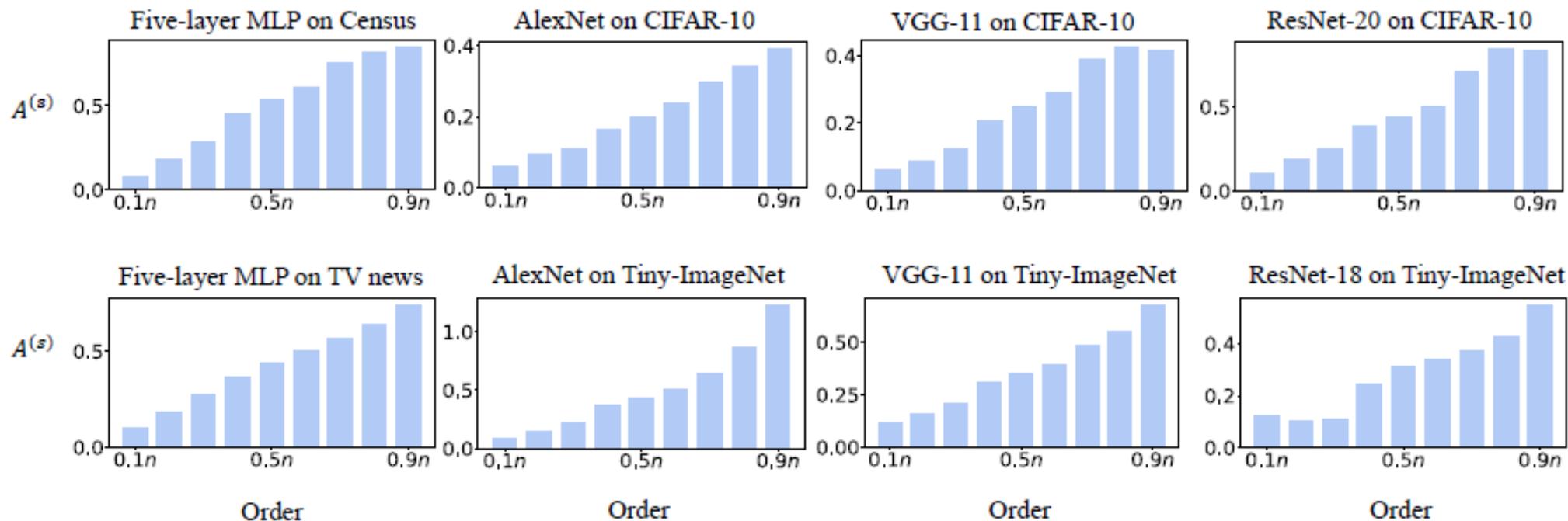


➤ Explaining generalization power and adversarial robustness

Low-order interactive concepts are less sensitive to adversarial perturbations

- Sensitivity of the interactive concept S to adversarial perturbations δ

$$\alpha(S) = \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [|I(S|\mathbf{x} + \delta) - I(S|\mathbf{x})|] / \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [|I(S|\mathbf{x})|] \quad A^{(s)} = \mathbb{E}_{S \subseteq N, |S|=s} [\alpha(S)]$$



Explaining existing findings about what are learned first

We discuss some related studies on which kind of knowledge is usually first learned by a DNN. We find that our theorems can partially explain mechanisms behind some previous findings[1-4].

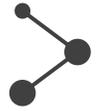
[1] Devansh Arpit, Stanisław Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In International conference on machine learning, pp. 233–242. PMLR, 2017..

[2] Karttikeya Mangalam and Vinay Uday Prabhu. Do deep neural networks learn shallow learnable examples first? 2019..

[3] Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle: Fourier analysis sheds light on deep neural networks. arXiv preprint arXiv:1901.06523, 2019..

[4] Chen Liu, Zhichao Huang, Mathieu Salzmann, Tong Zhang, and Sabine Süsstrunk. On the impact of hard adversarial instances on overfitting in adversarial training. arXiv preprint arXiv:2112.07324, 2021.

[5] Xu Cheng, Chuntung Chu, Yi Zheng, Jie Ren, and Quanshi Zhang. A game-theoretic taxonomy of visual concepts in dnns. arXiv preprint arXiv:2106.10938, 2021.



Conclusion

- We prove that low-order interactive concepts in the data are much **more stable than** high order concepts, which makes low-order interactive concepts more likely to be encoded.
- **We provide new insights into several empirical findings** *w.r.t.* the conceptual representation of DNNs.