

PAPER REVIEW PRESENTATION

# Where We Have Arrived in Proving the Emergence of Sparse Interaction Primitives in DNNs

Qihan Ren, Jiayang Gao, Wen Shen, and Quanshi Zhang  
ICLR, 2024

2026.4.30

SNU IDEA Lab. 정유립

# 01 Interaction-based works ( prof. Quanshi Zhang )

## [ Foundation (Mathematical Methodology)]

1. AAI 2021 Interpreting Multivariate Shapley Interactions in DNNs

## [ Interaction-based Understanding of DNNs ]

2. CVPR 2023 Defining and Quantifying the Emergence of Sparse Concepts in DNNs

3. ICML 2023 Does a Neural Network Really Encode Symbolic Concepts?

4. ICLR 2024 [Where We Have Arrived in Proving the Emergence of Sparse Interaction Primitives in DNNs](#)

5. NeurIPS 2024 Towards the Dynamics of a DNN Learning Symbolic Interactions

## [ Application ]

6. ICML 2023 Bayesian Neural Networks Avoid Encoding Complex and Perturbation-Sensitive Concepts

7. AAI 2024 Explaining Generalization Power of a DNN Using Interactive Concepts

## [ Architectural Innovation ]

8. ICML 2023 HarsanyiNet: Computing Accurate Shapley Values in a Single Forward Propagation

# 01 Interaction-based works ( prof. Quanshi Zhang )

1. AAI 2021 Interpreting Multivariate Shapley Interactions in DNNs
2. CVPR 2023 Defining and Quantifying the Emergence of Sparse Concepts in DNNs
3. ICML 2023 Does a Neural Network Really Encode Symbolic Concepts?
4. ICLR 2024 Where We Have Arrived in Proving the Emergence of Sparse Interaction Primitives in DNNs

1. AAI 2021

$$v(x) = \sum_{S \subseteq N} I(S) + v(x_\emptyset)$$

3. ICLR 2024

$$v(x) \approx \sum_{S \in \Omega_{\text{saliient}}} I(S) + v(x_\emptyset)$$

$$|\Omega_{\text{saliient}}| \ll 2^n$$

## Notation

Harsanyi Interaction  $I(S) \stackrel{\text{def}}{=} \sum_{T \subseteq S} (-1)^{|S|-|T|} \cdot u(T)$   
masked value function  $u(T) \stackrel{\text{def}}{=} v(x_T) - v(x_\emptyset)$

Trained DNN  $v$

Output of the DNN on the sample  $x$   $v(x) \in \mathbb{R}$

Input sample  $x = [x_1, \dots, x_n]^T$

Baseline value  $b = [b_1, \dots, b_n]^T = x_\emptyset$

Masked input sample  $x_T = [z_1, \dots, z_n]^T$ , Where  $z_i = \begin{cases} x_i, & i \in T \\ b_i, & i \in N \setminus T \end{cases}$

Set  $N = \{1, \dots, n\}$

Saliient Subset  $\Omega_{\text{saliient}} \subseteq N$

## 02 INTERACTION AS A SYMBOLIC CONCEPTS

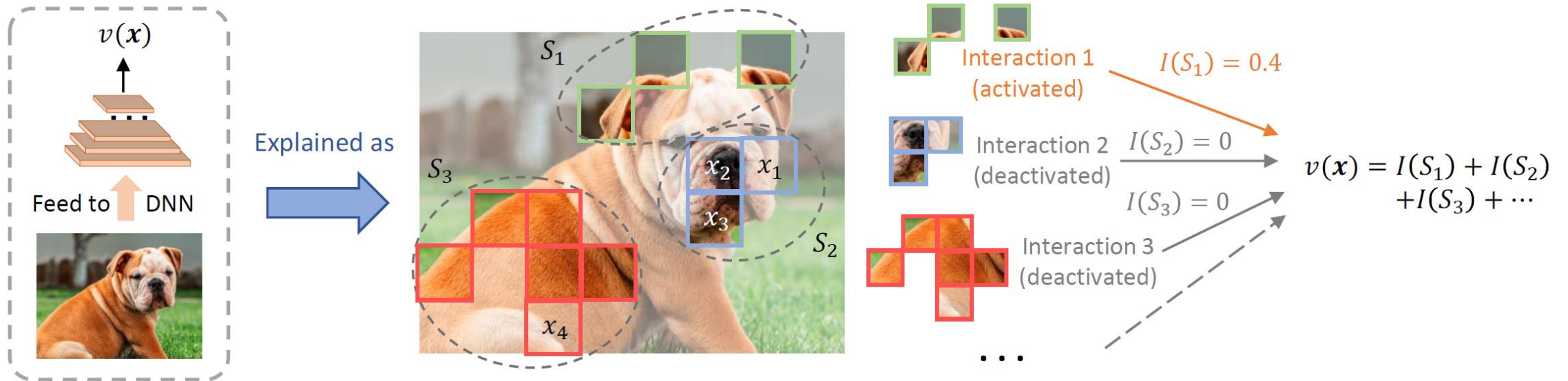


Figure 1: Illustration of interactions encoded by a DNN. Each interaction  $S$  corresponds to an AND relationship among a specific set  $S$  of input variables (image patches). The patches  $x_1$  and  $x_4$  are masked, so that interactions  $S_2$  and  $S_3$  are deactivated

### Axiomatic properties to define interactions as representations of faithful inference patterns (or concepts) encoded by a DNN

*Ren et al. (ICML 2023)*

#### (1) Sparsity property.

A DNN is supposed to encode few salient interactions on a specific sample.

#### (2) Universal matching property.

The network output on any arbitrarily masked sample is supposed to be well matched by the effects of specific interactions.

**Theorem 1** (Proven in [Ren et al. \(2023a\)](#) and [Appendix B.1](#)). *Let the input sample  $\mathbf{x}$  be arbitrarily masked to obtain a masked sample  $\mathbf{x}_S$ . The output of the DNN on masked sample  $\mathbf{x}_S$  can be disentangled into the sum of all interaction effects within  $S$ :  $\forall S \subseteq N, v(\mathbf{x}_S) = \sum_{T \subseteq S} I(T) + v(\mathbf{x}_\emptyset)$ .*

$$\ast v(\mathbf{x}) = \sum_{S \subseteq N} I(S) + v(\mathbf{x}_\emptyset)$$

#### (3) Sample-wise transferability property.

Salient interactions are supposed to be shared across different samples in the same category

# 03 PROVING THE SPARSITY OF INTERACTIONS

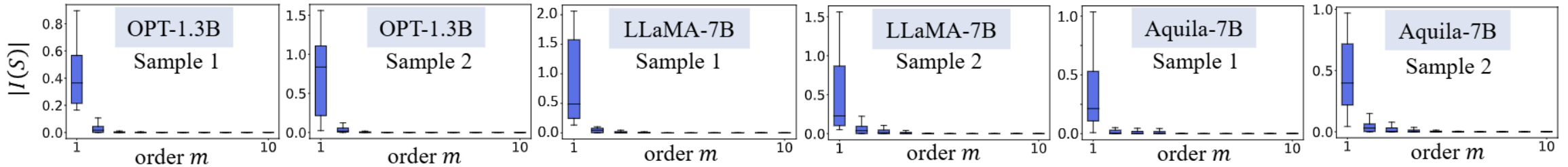
## Assumption 1

The DNN has at most  $M$ -th order non-zero derivatives, where  $M < n$ , and  $n$  is the number of input variables to the DNN

$$\forall S \subseteq N, |S| \geq M + 1, I(S) = 0$$

### Justification

(a) empirical evidence



(b) heuristic reasoning

(c) literature precedent

: 기존 game-theoretic interaction 연구들도 고차항 truncate 함.

# 03 PROVING THE SPARSITY OF INTERACTIONS

## Assumption 2

The DNN can be used on occluded samples (e.g., an image with some patches being masked), and yields a higher classification confidence when the sample is less occluded.

$$\forall m' \leq m \Rightarrow \bar{u}^{(m')} \leq \bar{u}^{(m)} \quad \bar{u}^{(m)} \stackrel{def}{=} \mathbb{E}_{|S|=m}[u(S)]$$
$$u(S) = v(x_S) - v(x_\emptyset)$$

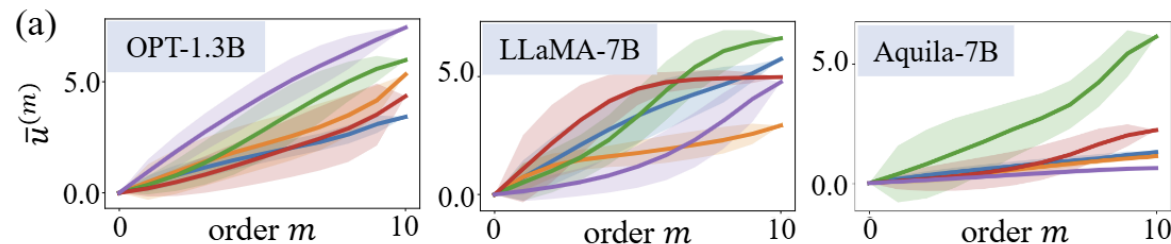
## Assumption 3

The classification confidence of the DNN does not significantly degrade on occluded samples

$$\bar{u}^{(m')} \geq \left(\frac{m'}{m}\right)^p \bar{u}^{(m)} \text{ for } m' \leq m, p > 0 \text{ (positive constant)}$$

## Justification

(a) empirical evidence



(b) heuristic reasoning

	OPT-1.3B	LLaMA-7B	Aquila-7B
Average $p$	$1.18 \pm 0.22$	$1.18 \pm 0.42$	$1.11 \pm 0.39$

## 03 PROVING THE SPARSITY OF INTERACTIONS

**Theorem 2** (Proven in Appendix B.3). *There exists  $m_0 \in \{n, n - 1, \dots, n - M\}$ , such that for all  $1 \leq k \leq M$ , the sum of the effects of all  $k$ -order interactions can be written as*

$$A^{(k)} = (\lambda^{(k)} n^{p+\delta} + a_{[p]-1}^{(k)} n^{[p]-1} + \dots + a_1^{(k)} n + a_0^{(k)}) \bar{u}^{(1)}, \quad (3)$$

where  $|\lambda^{(k)}| \leq 1$ ,  $|a_0^{(k)}| < n$ ,  $|a_i^{(k)}| \in \{0, 1, \dots, n - 1\}$  for  $i = 1, \dots, [p] - 1$ , and

$$\delta \leq \log_n \left( \frac{1}{\lambda} \left( 1 - \frac{a_{[p]-1}}{n^{p-[p]+1}} - \dots - \frac{a_0}{n^p} \right) \right), \quad \text{if } \lambda > 0, \quad (4)$$

$$\delta \leq \log_n \left( \frac{1}{-\lambda} \left( \frac{a_{[p]-1}}{n^{p-[p]+1}} + \dots + \frac{a_0}{n^p} \right) \right), \quad \text{if } \lambda < 0. \quad (5)$$

Here,  $\lambda \stackrel{\text{def}}{=} \sum_{k=1}^M \frac{\binom{m_0}{k}}{\binom{n}{k}} \lambda^{(k)} \neq 0$ ,  $a_i \stackrel{\text{def}}{=} \sum_{k=1}^M \frac{\binom{m_0}{k}}{\binom{n}{k}} a_i^{(k)}$  for  $i = 0, 1, \dots, [p] - 1$ , and  $[p]$  denotes the greatest integer that is less than or equal to  $p$ .

The above theorem indicates that the sum of effects of all  $k$ -order interactions is  $O(n^{p+\delta})$ .

## 03 PROVING THE SPARSITY OF INTERACTIONS

**Theorem 3** (Proven in Appendix B.4).  $R^{(k)}$  has the following upper bound:

$$R^{(k)} \leq \frac{\bar{u}^{(1)}}{\tau |\eta^{(k)}|} |\lambda^{(k)} n^{p+\delta} + a_{[p]-1}^{(k)} n^{[p]-1} + \dots + a_0^{(k)}|. \quad (6)$$

$$\eta^{(k)} \stackrel{\text{def}}{=} \frac{\sum_{|S|=k} I(S)}{\sum_{|S|=k} |I(S)|} \quad R^{(k)} \stackrel{\text{def}}{=} |\{S \subseteq N \mid |S| = k, |I(S)| \geq \tau\}|$$

The above theorem indicates that if positive interactions do not fully cancel with negative interactions (i.e.,  $|\eta^{(k)}|$  is not extremely small), then the number of valid interactions  $R^{(k)}$  of the  $k$ -th order has an upper bound of  $O(np+\delta / |\tau \eta^{(k)}|)$ , which is much less than the total number of  $n^k$  potential interactions of the  $k$ -th order in most cases

Table 2: Comparison between the number of valid interactions and the derived upper bound.

	OPT-1.3B	LLaMA-7B	Aquila-7B	MLP (tabular dataset)
Real # of valid interactions	28.73±52.37	50.53±40.37	30.13±26.20	54.42±36.81
Upper bound	197.84±188.87	293.20±287.28	184.23±124.71	229.11±139.52

감사합니다!