

# TOWARDS THE DYNAMICS OF REPRESENTATION CHANGES DURING PARAMETER PRUNING FOR NETWORK COMPRESSION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This study explores the dynamics of interactions encoded by deep neural networks (DNNs) when we prune network parameters with an increasing pruning ratio. We discover a three-phase dynamics of the generalizability of the interactions removed by the parameter pruning operation, which clarifies a central issue in symbolic generalization, i.e., how interactions serve as the underlying factors that determine the change of a DNN’s performance. Experimental results demonstrate that the pruning operation mainly removes high-order interactions at low pruning ratios. Because the removed high-order interactions are usually unlikely to generalize, the removal of high-order interactions has a negligible impact on testing performance. In contrast, under higher pruning ratios, both low-order and high-order interactions are gradually removed. The high generalizability of the removed low-order interactions leads to a noticeable decline in testing performance.

## 1 INTRODUCTION

In the field of post-hoc explanation of deep neural networks (DNNs), there has been a trend toward more precise and nuanced approaches (Li & Zhang, 2023b; Ren et al., 2024a). Mechanistic interpretability studies (surveyed in Appendix A) are developed to explain individual neurons. In comparison, another emerging direction, namely *symbolic generalizability*, has provided a new strategy to explain the generalizability of DNNs (Kang et al., 2024; Ren et al., 2023a; 2024a; Tsai et al., 2023). As Figure 1 shows, this theory aims to define and use the generalizability of inference patterns encoded in a DNN to explain the generalizability of the entire DNN.

**Background: symbolic generalizability.** Recent advances in symbolic generalizability (surveyed in Appendix A) have revealed a counterintuitive phenomenon: **the complex inference logic of a DNN can be precisely explained by a small set of AND-OR interactions in mathematics**. Given an input sample, each interaction represents an AND relationship (or an OR relationship) among input variables that is equivalently encoded by the DNN. For instance as illustrated in Figure 1, the large language model encodes an AND interaction among input words  $S=\{\text{“Electric”, “currents”, “around”}\}$ . If and only if all words in  $S$  appear in the input prompt, this interaction is triggered and contributes an effect of 0.45 to boost the confidence of generating the target word “conductors.” Masking any words in  $S$  will deactivate the interaction and remove its effect. It is proven by Chen et al. (2024) that **people can use numerical effects of such interactions to accurately predict the DNN’s classification scores on exponentially many diverse samples**, which guarantees the faithfulness of such explanation.

Therefore, the next fundamental issue in this field is to use the generalizability of the **compositional interactions** to explain the generalizability of the **entire DNN**<sup>1</sup> (Zhang et al., 2024; Zhou et al., 2024; Ren et al., 2024b; Cheng et al., 2025). For example, as Figure 2 shows, the interaction among image patches of “wing” consistently appears in both training samples and testing samples, and contributes similar interaction effects. Therefore, this interaction is considered generalizable. It is found that interactions in a highly generalizable DNN are more likely to generalize to testing samples than those in a less well-trained DNN (Zhou et al., 2024).

<sup>1</sup>It is because Eq. (3) shows the DNN’s output  $v(x)$  can be represented as the sum of interaction effects.



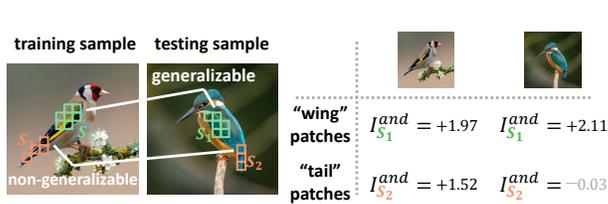


Figure 2: The interaction among “wing” patches consistently contributes similar interaction effects, thus being considered generalizable. The interaction among “tail” patches cannot generalize to testing samples, thus being considered non-generalizable.

## 2 DYNAMICS OF REPRESENTATION CHANGES DURING PARAMETER PRUNING

### 2.1 PRELIMINARIES: INTERACTIONS

Given an input sample with  $n$  input variables (indexed by the set  $N = \{1, 2, \dots, n\}$ ), denoted by  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ , we use  $v(\cdot)$  to represent the scalar output of the DNN. There exists different settings for the scalar output. Typically, for multi-category classification, we follow Deng et al. (2022) to set the widely used classification confidence as follows.

$$v(\mathbf{x}) \stackrel{\text{def}}{=} \log \frac{p(y = y^* | \mathbf{x})}{1 - p(y = y^* | \mathbf{x})} \quad (1)$$

where  $p(y = y^* | \mathbf{x})$  represent the predicted classification probability on the ground-truth category.

**Problem setting.** As an emerging explanation direction (surveyed in Appendix A), symbolic generalizability studies (Li & Zhang, 2023b; Ren et al., 2024a; Chen et al., 2024; Cheng et al., 2025) are proposed to explain the generalizability of primitive inference patterns encoded by a DNN. The basic idea is to use a logical model  $g(\cdot)$  to explain inference patterns in the DNN  $v(\cdot)$ , and the faithfulness of explanation is ensured by requiring a sufficiently concise logical model  $g$  to accurately explain a sufficiently large number of network outputs  $v$  w.r.t. the following two requirements. **(1) Fidelity requirement:** the logical model  $g(\cdot)$  is powerful enough to well match the network outputs  $v(\cdot)$  on all diverse samples in a sufficient large set  $\Psi$ . **(2) Conciseness requirement:** the logical model is supposed to only consist of a small number of inference patterns. These requirements can be formally expressed as

$$\forall \mathbf{x}' \in \Psi, g(\mathbf{x}') = v(\mathbf{x}') \quad \text{subject to} \quad \text{complexity}(g) \leq M, \quad (2)$$

where  $M$  is the upper complexity bound of the logical model  $g(\cdot)$ .

**First, the logical model  $g(\cdot)$  is implemented to encode AND-OR interactions among the input variables in  $N$  (Li & Zhang, 2023b; Chen et al., 2024).** Figure 1 depicts how each interactions corresponds to a non-linear relationship among a set of input variables. The formal definition of the model is given below. **See Appendix I for examples of logical models that explains LLMs. A video demo that introduces the logical model is also attached as a supplementary material.**

$$g(\mathbf{x}') \stackrel{\text{def}}{=} \sum_{S \in \Omega^{\text{and}}} I_S^{\text{and}} \cdot \delta_{\text{and}}(S | \mathbf{x}') + \sum_{S \in \Omega^{\text{or}}} I_S^{\text{or}} \cdot \delta_{\text{or}}(S | \mathbf{x}') + b \quad (3)$$

The AND trigger function  $\delta_{\text{and}}(S | \mathbf{x}')$  activates (returns 1) if and only if all variables in the subset  $S \subseteq N$  are present in  $\mathbf{x}'$ ; if any variable in  $S$  is masked<sup>2</sup>, the function returns 0.  $I_S^{\text{and}}$  is the scalar weight of the AND interaction  $S$ . Similarly, the OR trigger function  $\delta_{\text{or}}(S | \mathbf{x}')$  activates whenever at least one variable in the subset  $S$  is present.  $I_S^{\text{or}}$  is the scalar weight of the OR interaction  $S$ .  $\Omega^{\text{and}}$  and  $\Omega^{\text{or}}$  denote the set of AND interactions and that of OR interactions, respectively, extracted from the input  $\mathbf{x}'$ .  $b$  is a scalar bias term.

**Second, the fidelity requirement is guaranteed by the universal matching property in Theorem 2.1.** Under a specific settings of weights, the logical model  $g(\cdot)$  can accurately reproduce the outputs of the DNN  $v(\cdot)$ , no matter how we augment the input  $\mathbf{x}$  by randomly masking a subset of input variables. In other words, the sample set  $\Psi = \{\mathbf{x}_T | T \subseteq N\}$  contains  $2^n$  masked states, which is sufficiently large. Here,  $\mathbf{x}_T$  denotes the masked version of the sample  $\mathbf{x}$  that retains only the variables in  $T$ , with others in  $N \setminus T$  masked<sup>2</sup>.

<sup>2</sup>The masking of the  $i$ -th input variables is implemented by setting  $x_i$  to the baseline values  $b$ , which is commonly set as the average value of the variable on multiple samples (Dabkowski & Gal, 2017).

**Theorem 2.1 (Universal matching property, proven by (Chen et al., 2024) and Appendix G)** Given a sample  $\mathbf{x}$  and a DNN  $v(\cdot)$ , let us set the scalar weights  $I_S^{and}$  and  $I_S^{or}$  in  $g(\cdot)$  as  $\forall S \subseteq N$ ,  $I_S^{and} = \sum_{T \subseteq S} (-1)^{|S|-|T|} \cdot u_T^{and}$ ,  $I_S^{or} = -\sum_{T \subseteq S} (-1)^{|S|-|T|} \cdot u_{N \setminus T}^{or}$ , subject to  $u_T^{and} = 0.5 \cdot v(\mathbf{x}_T) + \gamma_T$  and  $u_T^{or} = 0.5 \cdot v(\mathbf{x}_T) - \gamma_T$ . bias  $b = v(\mathbf{x}_\emptyset)$ , and  $\{\gamma_T\}$  is a set of learnable parameters. Then we have  $\forall T \subseteq N$ ,  $v(\mathbf{x}_T) = g(\mathbf{x}_T)$ .

**Interaction extraction.** We adopt the method of Chen et al. (2024) to learn parameters  $\{\gamma_T\}$ , so as to extract a set of AND-OR interactions. Detailed pseudocode are provided in Appendix F.

**Third, the conciseness requirement is ensured by the sparsity property of interactions.** Ren et al. (2024a) have shown that a well-trained DNN encodes only a relatively small number of salient interactions, roughly  $O(n^p/\tau)$  with  $p \in [1.5, 2]$ , under three common conditions introduced in Appendix E. We can construct a much more compact model  $\hat{g}(\cdot)$  with the small number of salient interactions in  $\Omega^{and} = \{S \subseteq N : |I_S^{and}| > \tau\}$  and  $\Omega^{or} = \{S \subseteq N : |I_S^{or}| > \tau\}$ . Thus, we can accurately predict the network outputs over differently masked inputs.  $\tau$  is a scalar threshold, which is set to  $\tau = 0.05 \cdot \sum_{\mathbf{x}} |v(\mathbf{x}_T) - v(\mathbf{x}_\emptyset)|$ .

**Order of interactions.** The order of an AND/OR interaction  $S$ , written as  $\text{order}(S) = |S|$ , is defined by the quantity of variables in  $S$ . It serves as an indicator of the interaction’s complexity.

## 2.2 EXPLAINING THE GENERALIZABILITY OF DNNs IN TERMS OF INTERACTIONS

As introduced in Section 2.1, interactions provide a new perspective for explaining the root causes for the performance of DNNs. Prior studies (Zhou et al., 2024) have found that the generalizability of the **entire** DNN can be explained by the overall generalizability of its **compositional** interactions<sup>1</sup>. In particular, it has been found (Zhou et al., 2024; Liu et al., 2023) that the **complexity** and the **generalizability** of interactions are the key factors that determine network performance.

- Zhou et al. (2024) found that **interactions of higher orders were usually less generalizable**.
- Liu et al. (2023) found that **high-order interactions showed more vulnerability to feature noises**.
- The generalizability of interactions could explain the overfitting of DNNs** (Cheng et al., 2025).

Despite above achievements, existing studies (Ren et al., 2023a; Li & Zhang, 2023b; Ren et al., 2024a; Chen et al., 2024; Cheng et al., 2025) mainly passively and empirically explained the generalizability and complexity of interactions, and they still failed to *investigate the precise relationship between network performance and interaction patterns in terms of interventional studies*.

**Therefore, in this study, we focus on neural network parameter pruning, and aim to explore how the complexity and generalizability of interactions encoded by a DNN are affected when we apply a progressively increasing pruning ratio to it.** Furthermore, we aim to uncover the underlying mechanisms behind DNNs’ performance changes under parameter pruning. It is because parameter pruning is a classical strategy for removing trivial parameters from neural networks, which provides a valuable perspective for analyzing neural network behaviors.

**Quantifying the complexity of interactions.** Definition 2.1 uses the distribution of interaction strength over different orders to measure a DNN’s representational complexity (see Figure 3).

**Definition 2.1** Given a set of AND interactions  $\hat{\Omega}^{and}$  and a set of OR interactions  $\hat{\Omega}^{or}$ , the total strength of all  $m$ -order positive interactions  $\mathbf{I}^{(m),+}(\hat{\Omega}^{and}, \hat{\Omega}^{or})$  and the total strength of all  $m$ -order negative interactions  $\mathbf{I}^{(m),-}(\hat{\Omega}^{and}, \hat{\Omega}^{or})$  are defined as follows:

$$\begin{aligned} \mathbf{I}^{(m),+}(\hat{\Omega}^{and}, \hat{\Omega}^{or}) &= \sum_{S \in \hat{\Omega}^{and}; |S|=m} \max(I_S^{and}, 0) + \sum_{S \in \hat{\Omega}^{or}; |S|=m} \max(I_S^{or}, 0), \\ \mathbf{I}^{(m),-}(\hat{\Omega}^{and}, \hat{\Omega}^{or}) &= \sum_{S \in \hat{\Omega}^{and}; |S|=m} \min(I_S^{and}, 0) + \sum_{S \in \hat{\Omega}^{or}; |S|=m} \min(I_S^{or}, 0). \end{aligned} \quad (4)$$

**Quantifying the generalizability of interactions.** According to the original definition (Zhou et al., 2024), the interaction  $S$  is considered generalizable if it consistently contributes similar salient interaction effects across different testing samples. However, it requires to check interactions through all testing samples, which is computational infeasible. Thus, we adopt an approximate yet more efficient method (He et al., 2025) to identify generalizable interactions.

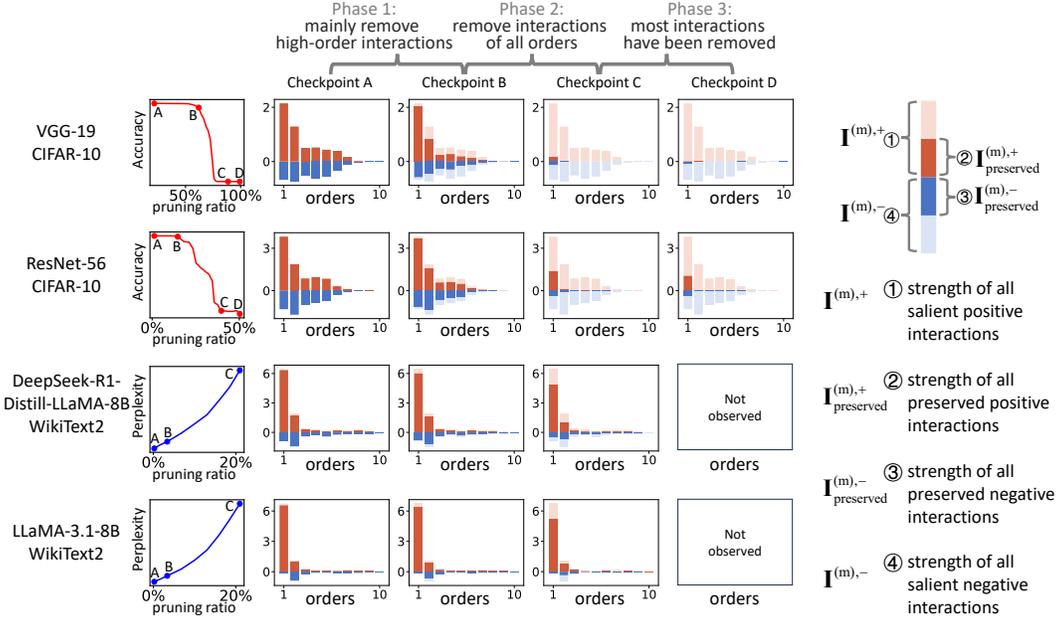


Figure 3: Changes of the distributions of  $I^{(m),+}$ ,  $I^{(m),-}$ ,  $I^{(m),+}_{\text{preserved}}$ , and  $I^{(m),-}_{\text{preserved}}$  when the pruning ratio increases. These changes can be divided into three phases. **See Appendix J for more results.**

Specifically, we train a separate DNN on testing samples, referred to as a *reference DNN*. An interaction  $S$  is defined generalizable if it is also utilized by this reference DNN for inference. This is because all interactions encoded by the reference DNN are learned from the testing samples.

**Definition 2.2** Let us be given a salient AND interaction  $S$  extracted by a DNN  $v(\cdot)$  from the input sample  $\mathbf{x}$ , subject to  $|I_S^{\text{and}}| > \tau$ . If this AND interaction is also extracted as a salient interaction by reference DNN  $v^*(\cdot)$ , and produces a consistent effect (i.e.,  $|I_S^{*,\text{and}}| > \tau$  and  $I_S^{\text{and}} \cdot I_S^{*,\text{and}} > 0$ ), then we consider this interaction generalizable. The generalizability of an OR interaction is defined similarly. Generalizable AND/OR interactions can be identified by the following two binary metrics:

$$G_S^{\text{and}} = \mathbb{1}(|I_S^{*,\text{and}}| > \tau \text{ and } I_S^{*,\text{and}} \cdot I_S^{\text{and}} > 0), \quad G_S^{\text{or}} = \mathbb{1}(|I_S^{*,\text{or}}| > \tau \text{ and } I_S^{*,\text{or}} \cdot I_S^{\text{or}} > 0) \quad (5)$$

where  $G_S^{\text{and}} \in \{0, 1\}$  and  $G_S^{\text{or}} \in \{0, 1\}$ .  $\mathbb{1}(\cdot)$  is an indicator function. It returns 1 if the condition is true, and returns 0 otherwise.  $I_S^{*,\text{and}}$  and  $I_S^{*,\text{or}}$  represent the numerical effects of the AND interaction  $S$  and the OR interaction  $S$ , respectively, extracted by the reference DNN.

### 2.3 THREE-PHASE DYNAMICS OF INTERACTIONS DURING PROGRESSIVE PRUNING

In this paper, we explore the dynamics of interactions. When we progressively increase the pruning ratio applied to a DNN, we discover a distinct three-phase dynamics of interaction complexities (orders), which well explains the performance change of the pruned DNN.

**Quantifying the removed and preserved interactions after parameter pruning.** To this end, we propose a set of metrics to quantify the change of interactions after parameter pruning. Given an original DNN  $v(\cdot)$  and a pruned DNN  $v'(\cdot)$ , we follow (Chen et al., 2024) to extract two sets of interaction weights  $\{I_S^{\text{and}}, I_S^{\text{or}} \mid S \subseteq N\}$  and  $\{I_S^{\text{and}'}, I_S^{\text{or}'}\mid S \subseteq N\}$  from each model, respectively. Then, we define the interactions removed after parameter pruning as follows.

**Definition 2.3** Let us be given a salient AND interaction  $S$  extracted by the original DNN  $v(\cdot)$  from the input  $\mathbf{x}$ , subject to  $|I_S^{\text{and}}| > \tau$ . If this interaction is no longer salient in the pruned DNN  $v'(\cdot)$  (subject to  $|I_S^{\text{and}'}| \leq \tau$ ) or exhibits a contrary effect ( $I_S^{\text{and}'} \cdot I_S^{\text{and}} < 0$ ), then this interaction is

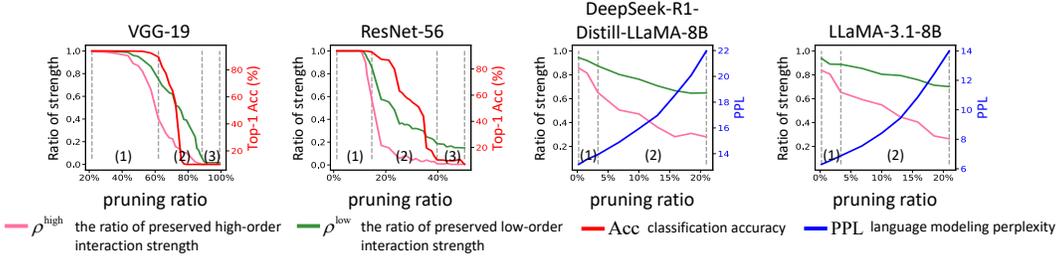


Figure 4: Changes of the ratio of preserved low-order interaction strength  $\rho^{\text{low}}$ , the ratio of preserved high-order interaction strength  $\rho^{\text{high}}$ , classification accuracy (or language modeling perplexity) when the pruning ratio increases. **Please see Appendix J for results on more DNNs.**

considered to have been removed by the parameter pruning operation. Thus, we use the following binary metrics to identify the removed AND/OR interactions.

$$R_S^{\text{and}} = \mathbb{1}(|I_S^{\text{and}}| \leq \tau \text{ or } I_S^{\text{and}} \cdot I_S^{\text{and}} < 0), \quad R_S^{\text{or}} = \mathbb{1}(|I_S^{\text{or}}| \leq \tau \text{ or } I_S^{\text{or}} \cdot I_S^{\text{or}} < 0) \quad (6)$$

Thus, let us use  $\Omega^{\text{and}} = \{S \subseteq N : |I_S^{\text{and}}| > \tau\}$  and  $\Omega^{\text{or}} = \{S \subseteq N : |I_S^{\text{or}}| > \tau\}$  to denote the sets of salient AND-OR interactions extracted from the original DNN  $v(\cdot)$ . Accordingly, we can use  $\Omega_{\text{preserved}}^{\text{and}} = \{S \in \Omega^{\text{and}} : R_S^{\text{and}} = 0\}$  and  $\Omega_{\text{preserved}}^{\text{or}} = \{S \in \Omega^{\text{or}} : R_S^{\text{or}} = 0\}$  to denote the preserved AND-OR interactions extracted from the pruned DNN  $v'(\cdot)$ .

**Evaluation metric.** We measure the ratio of low-order interactions (1-st to 3-rd order) that are preserved after parameter pruning, denoted by  $\rho^{\text{low}}$ , and the ratio of the preserved high-order interactions (4-th to  $n$ -th order), denoted by  $\rho^{\text{high}}$ , to evaluate the representation quality of the pruned DNN. This is because, according to Section 2.2 and Liu et al. (2023); Zhou et al. (2024), high-order interactions are less generalizable to testing samples and more sensitive to feature noises than low-order interactions. Thus, if the pruning operation mainly removes high-order interactions (causing a low  $\rho^{\text{high}}$  value), then such removal usually has little effect on the DNN’s performance, because most removed interactions are offsetting high-order interactions. Conversely, if the pruning operation mainly removes low-order interactions (causing a low  $\rho^{\text{low}}$  value), the DNN’s performance will be significantly degraded.

$$\rho^{\text{low}} = \frac{\sum_{\text{type} \in \{\text{and}, \text{or}\}} \sum_{S \in \Omega^{\text{type}}: 1 \leq |S| \leq 3} |I_S^{\text{type}}| \cdot (1 - R_S^{\text{type}})}{\sum_{\text{type} \in \{\text{and}, \text{or}\}} \sum_{S \in \Omega^{\text{type}}: 1 \leq |S| \leq 3} |I_S^{\text{type}}|} \quad (7)$$

$$\rho^{\text{high}} = \frac{\sum_{\text{type} \in \{\text{and}, \text{or}\}} \sum_{S \in \Omega^{\text{type}}: 4 \leq |S| \leq n} |I_S^{\text{type}}| \cdot (1 - R_S^{\text{type}})}{\sum_{\text{type} \in \{\text{and}, \text{or}\}} \sum_{S \in \Omega^{\text{type}}: 4 \leq |S| \leq n} |I_S^{\text{type}}|}$$

Figure 3 compares the distribution of interaction strength over different orders between those extracted from the original DNN  $v(\cdot)$  (measured by  $\mathbf{I}^{(m),+} = \mathbf{I}^{(m),+}(\Omega^{\text{and}}, \Omega^{\text{or}})$  and  $\mathbf{I}^{(m),-} = \mathbf{I}^{(m),-}(\Omega^{\text{and}}, \Omega^{\text{or}})$ ) and those extracted from the pruned DNN  $v'(\cdot)$  (measured by  $\mathbf{I}_{\text{preserved}}^{(m),+} = \mathbf{I}^{(m),+}(\Omega_{\text{preserved}}^{\text{and}}, \Omega_{\text{preserved}}^{\text{or}})$  and  $\mathbf{I}_{\text{preserved}}^{(m),-} = \mathbf{I}^{(m),-}(\Omega_{\text{preserved}}^{\text{and}}, \Omega_{\text{preserved}}^{\text{or}})$ ). Figure 4 shows the changes of the ratio of low-order interaction strength that are preserved after pruning  $\rho^{\text{low}}$  and the ratio of preserved high-order interaction strength  $\rho^{\text{high}}$  when the pruning ratio increases, along with the testing performance of DNNs (measured by classification accuracy or language modeling perplexity). Based on the above experimental results, we can divide the entire dynamics of interactions into three phases.

•**Phase 1:** As Figure 3 and Figure 4 shows, when the pruning ratio is low, only high-order interactions are removed, while low-order ones remain largely unaffected (*i.e.*,  $\rho^{\text{high}}$  noticeably decreases but  $\rho^{\text{low}} \approx 1$ ). Note that the removal of high-order interactions usually does not affect the testing accuracy, which has been supported by multiple lines of evidence. (1) Zhou et al. (2024) have found that high-order interactions are less generalizable to testing samples than low-order interactions. (2) The interactions removed in this phase exhibit a spindle-shaped distribution over different orders. A half of these removed interactions produce positive effects and boost the classification score, while the other half produce negative effects and reduce the classification score.

Table 1: Analysis of three-phase dynamics of interactions.

Perspectives	Phase 1	Phase 2	Phase 3
Testing accuracy	barely affected	significantly decreased	remains low
Low-order interactions	slightly removed	significantly removed	most have been removed
High-order interactions	significantly removed	significantly removed	most have been removed
Generalizability of the removed	slightly increase	significantly increase	most have been removed
Generalizability of the emerged	remains low	remains low	remains low

•*Phase 2*: When the pruning ratio further increases, low-order interactions are also removed besides the removal of high-order interactions. The testing accuracy begins to decrease. The decrease of the testing accuracy is caused by the removal of low-order interactions, because most low-order interactions have been found by Zhou et al. (2024) to represent well-trained inference patterns that can generalize to testing samples.

•*Phase 3*: When the pruning ratio increases to even higher levels, both low-order and high-order interactions show only minor changes because most of them have already been removed. The testing accuracy remains extremely low under excessive parameter pruning.

In sum, when the pruning ratio is low (during Phase 1), parameter pruning mainly removes high-order interactions, which barely affects the performance of DNNs. In contrast, when the pruning ratio is high (during Phases 2 and 3), low-order interactions start to be removed, which explains the performance degradation of DNNs. This discovery provides further support for the assumption in prior work on symbolic generalizability (Zhou et al., 2024; Cheng et al., 2025).

**Settings for network pruners.** We conducted experiments on eight DNNs with three network pruners. The eight DNNs included two LLMs for language generation, two ViT models, and four CNNs for image classification. Considering the adaptability of pruning algorithms to models, (1) we applied LLM-Pruner (Ma et al., 2023) to two different LLMs, including DeepSeek-R1-Distill-LLaMA-8B model (DeepSeek-AI et al., 2025) and LLaMA-3.1-8B model (Grattafiori et al., 2024); (2) we applied Isomorphic-Pruning (Fang et al., 2024) to two residual networks and two ViT models pretrained on the ImageNet-1k dataset (Russakovsky et al., 2015); (3) we applied DepGraph (Fang et al., 2023) to two small networks for image classification on the CIFAR-10 dataset (Krizhevsky, 2009). **Please refer to Appendix H for detailed settings and Appendix J for detailed results.**

## 2.4 UNDERSTANDING THE THREE-PHASE DYNAMICS IN TERMS OF GENERALIZABILITY

In this subsection, we aim to explore the three-phase dynamics of the generalizability of interactions. The precise and fine-grained explanation of detailed performance changes of a DNN presents the core challenge for symbolic generalization.

Specifically, the dynamics of generalizability of interactions can well explain the change of generalizability of the DNN. Table 1 shows the analytical results of the three-phase dynamics of interactions from multiple perspectives. Detailed results can be found in Appendix J.

**Distribution of generalizable interactions.** Based on Definition 2.1, the distribution of generalizable interactions across each  $m$ -th order can be quantified by the total strength of all positive generalizable interactions  $\mathbf{I}_{\text{generalizable}}^{(m),+} = \mathbf{I}^{(m),+}(\Omega_{\text{generalizable}}^{\text{and}}, \Omega_{\text{generalizable}}^{\text{or}})$ , and the total strength of all negative generalizable interactions  $\mathbf{I}_{\text{generalizable}}^{(m),-} = \mathbf{I}^{(m),-}(\Omega_{\text{generalizable}}^{\text{and}}, \Omega_{\text{generalizable}}^{\text{or}})$ , where  $\Omega_{\text{generalizable}}^{\text{and}} = \{S \in \Omega^{\text{and}} : G_S^{\text{and}} = 1\}$  and  $\Omega_{\text{generalizable}}^{\text{or}} = \{S \in \Omega^{\text{or}} : G_S^{\text{or}} = 1\}$  denote the sets of generalizable AND interactions and generalizable OR interactions, respectively.

Figure 26 in Appendix K shows the distributions of  $\mathbf{I}^{(m),+}$ ,  $\mathbf{I}^{(m),-}$ ,  $\mathbf{I}_{\text{generalizable}}^{(m),+}$  and  $\mathbf{I}_{\text{generalizable}}^{(m),-}$  over different orders, which are extracted from eight DNNs. Most generalizable interactions are low-order interactions, and most high-order interactions cannot generalize to the reference DNN. Thus, these high-order interactions are considered as non-generalizable patterns.

**Exploring the efficiency of parameter pruning in terms of generalizability of the removed interactions.** We investigate generalizability of interactions removed by parameter pruning, and

Table 2: Efficiency of parameter pruning in the three-phases. See Appendix J for more results.

Architecture	Phase 1	Phase 2	Phase 3
ResNet-56	$e^{1.21 \rightarrow 14.65} = 0.92$	$e^{14.65 \rightarrow 39.68} = 0.78$	$e^{39.68 \rightarrow 50.27} = 0.38$
VGG-19	$e^{21.64 \rightarrow 62.02} = 0.89$	$e^{62.02 \rightarrow 88.54} = 0.61$	$e^{88.54 \rightarrow 99.23} = 0.36$
DeepSeek-R1-Distill-LLaMA-8B	$e^{0.23 \rightarrow 3.37} = 0.78$	$e^{3.37 \rightarrow 20.96} = 0.75$	–
LLaMA-3.1-8B	$e^{0.23 \rightarrow 3.37} = 0.74$	$e^{3.37 \rightarrow 20.96} = 0.70$	–

thereby analyze the detailed utility of a neural network compression method. In the first phase, pruning mainly removes non-generalizable interactions. In the second and third phases, it mainly removes generalizable interactions, which affects the DNNs’ performance. To verify this, let us first quantify the total strength of all removed interactions  $\mathbf{A}_{\text{removed}}$ , as well as the strength of removed generalizable interactions  $\mathbf{G}_{\text{removed}}$ .

$$\mathbf{A}_{\text{removed}} = \sum_{\substack{\text{type} \in \\ \{\text{and}, \text{or}\}}} \sum_{S \in \Omega^{\text{type}}} |I_S^{\text{type}}| \cdot R_S^{\text{type}}, \quad \mathbf{G}_{\text{removed}} = \sum_{\substack{\text{type} \in \\ \{\text{and}, \text{or}\}}} \sum_{S \in \Omega^{\text{type}}} |I_S^{\text{type}}| \cdot R_S^{\text{type}} \cdot G_S^{\text{type}} \quad (8)$$

In this way, let  $\mathbf{A}_{\text{removed}}(r_1)$  and  $\mathbf{A}_{\text{removed}}(r_2)$  denote the total strength of the removed interactions at two compression ratios  $r_1$  and  $r_2$  ( $r_2 > r_1$ ). Correspondingly, we compute  $\mathbf{G}_{\text{removed}}(r_1)$  and  $\mathbf{G}_{\text{removed}}(r_2)$  based on the generalizable interactions. Then, the efficiency of parameter pruning  $e^{r_1 \rightarrow r_2}$  within the interval  $[r_1, r_2]$  is defined as follows.

$$e^{r_1 \rightarrow r_2} = 1 - \frac{\mathbf{G}_{\text{removed}}(r_2) - \mathbf{G}_{\text{removed}}(r_1)}{\mathbf{A}_{\text{removed}}(r_2) - \mathbf{A}_{\text{removed}}(r_1)} \quad (9)$$

This metric reflects the ratio of non-generalizable interactions to all the removed interactions. A higher value of  $e^{r_1 \rightarrow r_2}$  indicates that less generalizable interactions are newly removed at pruning ratio  $r_2$ , compared with the DNN obtained with the pruning ratio  $r_1$ . We conducted experiments to illustrate the generalizability of the removed interactions after parameter pruning. According to Definition 2.2, we need to train reference DNNs from scratch on the testing samples. However, for large-scale models, it is often difficult to identify their training and testing samples. More importantly, retraining the model from scratch is impractical. Therefore, we used other classic models that were independently trained as the reference DNNs. Please see Appendix H for details.

Table 2 shows that parameter pruning exhibits fully different efficiency values  $e^{r_1 \rightarrow r_2}$  in different phases. Specifically, (1) in the first phase when the pruning ratio increases is low, the pruning operation mainly removes high-order interactions, but few generalizable interaction are removed, thereby exhibiting a high pruning efficiency  $e^{r_1 \rightarrow r_2}$ . The performance of the DNN is not affected significantly, and generalizability of the entire DNN is largely preserved after parameter pruning. (2) In the second phase, both low-order interactions and high-order interactions are gradually removed. In this way, generalizable interactions (mainly of low orders) are removed, and the efficiency value  $e^{r_1 \rightarrow r_2}$  begins to decrease. As a result, the testing accuracy of the DNN decreases significantly. (3) In the third phase, since most interactions in the original DNN have been removed, the performance of the DNN does not decrease a lot, and the very few removed interactions lead to an unstable efficiency  $e^{r_1 \rightarrow r_2}$ .

**Exploring generalizability of the emerged interactions during the three-phase dynamics.** In addition to removing interactions, we notice that the parameter-pruning operation may also bring in some new interactions. *Generally, we consider these newly emerged interactions to be purely noise patterns, because parameter pruning does not involve any data-driven learning process. In order to verify this hypothesis, we define and quantify the emerged interactions as follows.*

**Definition 2.4** *Let us be given a salient AND interaction  $S$  extracted by the pruned DNN  $v'(\cdot)$  from the input  $\mathbf{x}$ , subject to  $|I_S^{\text{and}}| > \tau$ . If this interaction is negligible in the original DNN  $v(\cdot)$  (subject to  $|I_S^{\text{and}}| \leq \tau$ ), then this interaction is considered as an emerged interaction after parameter pruning. The emergence of an OR interaction is defined similarly. Thus, we use the following binary metrics  $E_S^{\text{and}}, E_S^{\text{or}} \in \{0, 1\}$  to represent the emergence of AND-OR interactions.*

$$E_S^{\text{and}} = \mathbb{1}(|I_S^{\text{and}}| \leq \tau \text{ and } |I_S^{\text{and}}| > \tau), \quad E_S^{\text{or}} = \mathbb{1}(|I_S^{\text{or}}| \leq \tau \text{ and } |I_S^{\text{or}}| > \tau) \quad (10)$$

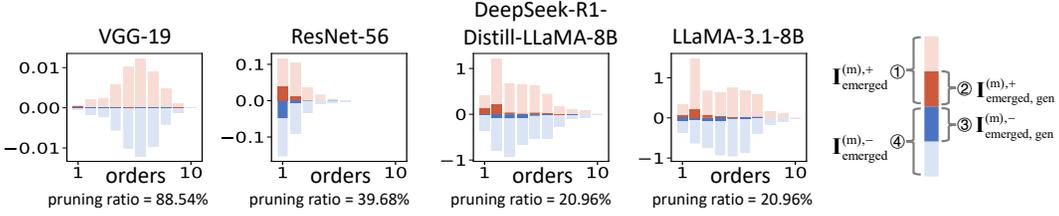


Figure 5: The distributions of  $\mathbf{I}_{\text{emerged}}^{(m),+}$ ,  $\mathbf{I}_{\text{emerged}}^{(m),-}$ ,  $\mathbf{I}_{\text{emerged,gen}}^{(m),+}$  and  $\mathbf{I}_{\text{emerged,gen}}^{(m),-}$ . Emerged interactions are barely generalizable. **Please see Appendix J for results on more DNNs.**

In this way, the distribution of emerged interactions over different orders after parameter pruning can be quantified by the total strength of all positive and negative emerged interactions, denoted by  $\mathbf{I}_{\text{emerged}}^{(m),+} = \mathbf{I}^{(m),+}(\Omega_{\text{emerged}}^{\text{and}}, \Omega_{\text{emerged}}^{\text{or}})$  and  $\mathbf{I}_{\text{emerged}}^{(m),-} = \mathbf{I}^{(m),-}(\Omega_{\text{emerged}}^{\text{and}}, \Omega_{\text{emerged}}^{\text{or}})$ , respectively, and the total strength of all positive and negative emerged generalizable interactions, denoted by  $\mathbf{I}_{\text{emerged,gen}}^{(m),+} = \mathbf{I}^{(m),+}(\Omega_{\text{emerged,gen}}^{\text{and}}, \Omega_{\text{emerged,gen}}^{\text{or}})$  and  $\mathbf{I}_{\text{emerged,gen}}^{(m),-} = \mathbf{I}^{(m),-}(\Omega_{\text{emerged,gen}}^{\text{and}}, \Omega_{\text{emerged,gen}}^{\text{or}})$ , respectively.

$$\begin{aligned} \Omega_{\text{emerged}}^{\text{and}} &= \{S \in \Omega^{\text{and}} : E_S^{\text{and}} = 1\}, & \Omega_{\text{emerged}}^{\text{or}} &= \{S \in \Omega^{\text{or}} : E_S^{\text{or}} = 1\} \\ \Omega_{\text{emerged,gen}}^{\text{and}} &= \{S \in \Omega^{\text{and}} : E_S^{\text{and}} = G_S^{\text{and}} = 1\}, & \Omega_{\text{emerged,gen}}^{\text{or}} &= \{S \in \Omega^{\text{or}} : E_S^{\text{or}} = G_S^{\text{or}} = 1\} \end{aligned} \quad (11)$$

where  $\Omega_{\text{emerged}}^{\text{and}}$  and  $\Omega_{\text{emerged}}^{\text{or}}$  denote the sets of emerged interactions, and  $\Omega_{\text{emerged,gen}}^{\text{and}} \subseteq \Omega_{\text{emerged}}^{\text{and}}$  and  $\Omega_{\text{emerged,gen}}^{\text{or}} \subseteq \Omega_{\text{emerged}}^{\text{or}}$  denote the subsets of interactions that are generalizable, respectively.

We conducted experiments to illustrate the generalizability of the emerged interactions after parameter pruning. We used the same reference DNN as in Section 2.3, when we analyzed generalizability of removed interactions. Figure 5 shows the distribution of the emerged interactions brought by the parameter pruning operation at a certain pruning ratio. Most emerged interactions exhibits mutually offsetting effects and have low generalizability. A half of emerged interactions boost the classification score, and the other half decrease the classification score. According to He et al. (2025), such newly emerged mutually offsetting interactions usually represent fully noise patterns.

In summary, parameter pruning affects the DNN in two ways. (1) It removes both generalizable and non-generalizable interactions, and (2) it also brings in some new non-generalizable interactions. The observed alignment between the change of a DNN’s generalizability and the dynamics of interactions further validates strong connection between the generalizability of the entire DNN and the generalizability of its compositional interactions.

### 3 CONCLUSION

In this paper, we have discovered a specific three-phase dynamics of interactions encoded by DNNs, when parameters are pruned under progressively increasing pruning ratios. Compared to traditional black-box evaluations, interaction-based analysis offers a more insightful perspective to understand the explicit influence of parameter pruning on the complexity and generalizability of a DNN’s interaction representations. Specifically, we find that low-order interactions generally exhibit substantially stronger generalizability than high-order ones, and they tend to persist until relatively high pruning ratios are reached. Since we have discovered that performance degradation in neural networks can be largely attributed to the removal of generalizable low-order interactions, this work provides a new perspective to analyze the optimal compression ratio at the level of interactions.

While the three-phase dynamics of interactions has been consistently observed across multiple neural networks and tasks, a rigorous theoretical account of its root cause remains absent. In particular, it is unclear why low-order interactions exhibit significantly greater robustness to pruning.

Our method is of considerable application value. In preliminary follow-up works, we penalize all non-generalizable interactions during the training of the DNN, which enables the DNN to be compressed at higher pruning ratios while maintaining its performance (see Appendix B).

## REFERENCES

- 486  
487  
488 Lu Chen, Siyu Lou, Benhao Huang, and Quanshi Zhang. Defining and extracting generalizable  
489 interaction primitives from dnns. In *International Conference on Learning Representations*, 2024.
- 490  
491 Lei Cheng, Junpeng Zhang, Qihan Ren, and Quanshi Zhang. Revisiting generalization power of a  
492 dnn in terms of symbolic interactions. *arXiv preprint arXiv:2502.10162*, 2025.
- 493  
494 Piotr Dabkowski and Yarin Gal. Real time image saliency for black box classifiers. In I. Guyon,  
495 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.),  
496 *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,  
497 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/  
file/0060ef47b12160b9198302ebdb144dcf-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/0060ef47b12160b9198302ebdb144dcf-Paper.pdf).
- 498  
499 DeepSeek-AI et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learn-  
500 ing, 2025. URL <https://arxiv.org/abs/2501.12948>.
- 501  
502 Huiqi Deng, Qihan Ren, Hao Zhang, and Quanshi Zhang. Discovering and explaining the represen-  
503 tation bottleneck of dnns. In *International Conference on Learning Representations*, 2022. URL  
<https://openreview.net/forum?id=iRCUlgmdfHJ>.
- 504  
505 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
506 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-  
507 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recogni-  
508 tion at scale. In *International Conference on Learning Representations*, 2021. URL [https://  
openreview.net/forum?id=YicbFdNTTy](https://openreview.net/forum?id=YicbFdNTTy).
- 509  
510 Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards  
511 any structural pruning. In *Proceedings of the IEEE/CVF conference on computer vision and  
512 pattern recognition*, pp. 16091–16101, 2023.
- 513  
514 Gongfan Fang, Xinyin Ma, Michael Bi Mi, and Xinchao Wang. Isomorphic pruning for vision  
515 models. In *Proceedings of the European Conference on Computer Vision (ECCV) 2024 Posters*,  
516 2024. URL <https://eccv.ecva.net/virtual/2024/poster/2295>. Poster.
- 517  
518 Aaron Grattafiori et al. The llama 3 herd of models, 2024. URL [https://arxiv.org/abs/  
2407.21783](https://arxiv.org/abs/2407.21783).
- 519  
520 Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural network  
521 with pruning, trained quantization and huffman coding. In *International Conference on Learning  
522 Representations*, 2016.
- 523  
524 John C Harsanyi. A simplified bargaining model for the n-person cooperative game. *International  
525 Economic Review*, 4(2):194–220, 1963.
- 526  
527 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
528 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
770–778, 2016.
- 529  
530 Yuxuan He, Junpeng Zhang, Hongyuan Zhang, and Quanshi Zhang. Technical report: Quantifying  
531 and analyzing the generalization power of a dnn. *arXiv preprint arXiv:2505.06993*, 2025.
- 532  
533 Justin Singh Kang, Yigit Efe Erginbas, Landon Butler, Ramtin Pedarsani, and Kannan Ramchan-  
534 dran. Learning to understand: Identifying interactions via the möbius transform. In A. Globerson,  
535 L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in  
536 Neural Information Processing Systems*, volume 37, pp. 46160–46202. Curran Associates, Inc.,  
537 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/  
file/520b379123d16e41f85472e766846486-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/520b379123d16e41f85472e766846486-Paper-Conference.pdf).
- 538  
539 Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical re-  
port, University of Toronto, 2009. URL [https://www.cs.toronto.edu/~kriz/  
learning-features-2009-TR.pdf](https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf).

- 540 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep  
541 convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger  
542 (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc.,  
543 2012. URL [https://proceedings.neurips.cc/paper\\_files/paper/2012/  
544 file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
- 545 Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. [https://tinyimagenet.  
546 stanford.edu/](https://tinyimagenet.stanford.edu/), 2015.
- 547
- 548 Mingjie Li and Quanshi Zhang. Defining and quantifying and-or interactions for faithful and concise  
549 explanation of dnns. *arXiv preprint arXiv:2304.13312*, 2023a.
- 550
- 551 Mingjie Li and Quanshi Zhang. Does a neural network really encode symbolic concepts? In  
552 *International conference on machine learning*, pp. 20452–20469. PMLR, 2023b.
- 553
- 554 Dongrui Liu, Huiqi Deng, Xu Cheng, Qihan Ren, Kangrui Wang, and Quanshi Zhang. Towards the  
555 difficulty for a deep neural network to learn concepts of different complexities. In *Advances in  
556 Neural Information Processing Systems*, volume 36, pp. 41283–41304, 2023.
- 557
- 558 Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large  
559 language models. In *Advances in Neural Information Processing Systems*, volume 36, pp. 21702–  
560 21720, 2023.
- 561 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and  
562 editing factual associations in gpt. In S. Koyejo, S. Mohamed, A. Agarwal,  
563 D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Process-  
564 ing Systems*, volume 35, pp. 17359–17372. Curran Associates, Inc., 2022. URL  
565 [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/  
566 6f1d43d5a82a37e89b0665b33bf3a182-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/6f1d43d5a82a37e89b0665b33bf3a182-Paper-Conference.pdf).
- 567
- 568 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
569 models. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Byj72udxe>.
- 570
- 571 Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan  
572 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,  
573 Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin  
574 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li,  
575 Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,  
576 Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2024.  
577 URL <https://arxiv.org/abs/2412.15115>.
- 578
- 579 Jie Ren, Mingjie Li, Qirui Chen, Huiqi Deng, and Quanshi Zhang. Defining and quantifying the  
580 emergence of sparse concepts in dnns. In *Proceedings of the IEEE/CVF conference on computer  
581 vision and pattern recognition*, pp. 20280–20289, 2023a.
- 582
- 583 Jie Ren, Zhanpeng Zhou, Qirui Chen, and Quanshi Zhang. Can we faithfully represent absence states  
584 to compute shapley values on a DNN? In *The Eleventh International Conference on Learning  
585 Representations*, 2023b. URL <https://openreview.net/forum?id=YV8tP7bW6Kt>.
- 586
- 587 Qihan Ren, Huiqi Deng, Yunuo Chen, Siyu Lou, and Quanshi Zhang. Bayesian neural networks  
588 avoid encoding complex and perturbation-sensitive concepts. In *International Conference on  
589 Machine Learning*, pp. 28889–28913. PMLR, 2023c.
- 590
- 591 Qihan Ren, Jiayang Gao, Wen Shen, and Quanshi Zhang. Where we have arrived in proving the  
592 emergence of sparse interaction primitives in dnns. In *International Conference on Learning  
593 Representations*, 2024a.
- 594
- 595 Qihan Ren, Junpeng Zhang, Yang Xu, Yue Xin, Dongrui Liu, and Quanshi Zhang. Towards the  
596 dynamics of a dnn learning symbolic interactions. In *Advances in Neural Information Processing  
597 Systems*, volume 37, pp. 50653–50688, 2024b.

- 594 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
595 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C Berg, and Li Fei-Fei.  
596 Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*  
597 (*IJCV*), 115(3):211–252, 2015.
- 598 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image  
599 recognition. In *International Conference on Learning Representations*, 2015. URL <https://arxiv.org/abs/1409.1556>.
- 600 Andreas Peter Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit,  
601 and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision  
602 transformers. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL  
603 <https://openreview.net/forum?id=4nPswr1KcP>.
- 604 Mukund Sundararajan, Kedar Dhamdhere, and Ashish Agarwal. The shapley taylor interaction  
605 index. In *International conference on machine learning*, pp. 9259–9268. PMLR, 2020.
- 606 Che-Ping Tsai, Chih-Kuan Yeh, and Pradeep Ravikumar. Faith-shap: The faithful shapley interac-  
607 tion index. *Journal of Machine Learning Research*, 24(94):1–42, 2023.
- 608 Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.  
609 Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In  
610 *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=NpsVSN6o4ul>.
- 611 Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training  
612 procedure in timm. In *NeurIPS 2021 Workshop on ImageNet: Past, Present, and Future*, 2021.  
613 URL <https://arxiv.org/abs/2110.00476>. Workshop paper.
- 614 Junpeng Zhang, Qing Li, Liang Lin, and Quanshi Zhang. Two-phase dynamics of interactions  
615 explains the starting point of a dnn learning over-fitted features. *arXiv preprint arXiv:2405.10262*,  
616 2024.
- 617 Bolei Zhou, David Bau, Aude Oliva, and Antonio Torralba. Interpreting deep visual representations  
618 via network dissection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(9):  
619 2131–2145, 2019. doi: 10.1109/TPAMI.2018.2858759.
- 620 Huilin Zhou, Hao Zhang, Huiqi Deng, Dongrui Liu, Wen Shen, Shih-Han Chan, and Quanshi Zhang.  
621 Explaining generalization power of a dnn using interactive concepts. In *Proceedings of the AAAI*  
622 *Conference on Artificial Intelligence*, volume 38, pp. 17105–17113, 2024.

## 623 A RELATED WORK ON SYMBOLIC GENERALIZABILITY

624 **Using interactions to explain the detailed inference logic encoded by a DNN.** Prior works on  
625 symbolic generalizability (Kang et al., 2024; Li & Zhang, 2023b; Tsai et al., 2023; Sundarara-  
626 jan et al., 2020) proposed to explain the inference logic of a DNN by quantifying the interactions  
627 among input variables encoded by the DNN. Building on this, Ren et al. (2024a) further proved that  
628 networks producing relatively smooth outputs across different input perturbations typically encode a  
629 small number of interactions. Chen et al. (2024) developed an algorithm to extract interactions that  
630 were commonly encoded by multiple different neural networks. Ren et al. (2023b) further improved  
631 a method to optimize the baseline value, which was used for interaction extraction.

632 **Using interactions to explain the generalizability of a DNN.** Furthermore, symbolic generalizabil-  
633 ity research explain a DNN’s generalizability as the overall generalizability of its compositional in-  
634 teractions. Zhou et al. (2024) found that high-order interactions usually generalize worse to testing  
635 samples than lower-order interactions. citetren2021towards further found that these higher-order  
636 interactions also showed poorer adversarial robustness. Ren et al. (2023c) found that mean-field  
637 Bayesian neural networks typically struggled more than regular neural networks in modeling higher-  
638 order interactions.

639 Although symbolic generalizability offers a powerful strategy to explain the inference logic of DNNs  
640 and their generalizability across different data domains, existing studies still leave a blank in the  
641

field of post-training interventional techniques. To this end, by investigating parameter pruning techniques, this paper bridges a crucial gap between symbolic generalizability theory and parameter pruning techniques, and offers a new perspective for understanding how interventional techniques, such as parameter pruning operations, reshape network representation quality and generalizability.

**Comparison with mechanistic interpretability.** Mechanistic interpretability treats a neural network as a collection of distinct components rather than a single black box (Zhou et al., 2019; Meng et al., 2022; Wang et al., 2023). It studies the roles of individual neurons, attention heads, or sub-networks, and how these parts work together to produce the model’s overall behavior. Instead of physically dividing a neural network into sub-components, symbolic generalizability treats a neural network as a whole and explains the inference logic equivalently encoded by the network using a set of AND-OR interactions.

## B PRACTICAL VALUE OF OUR FINDINGS

According to our findings, low-order generalizable interactions are removed at higher pruning ratios than high-order non-generalizable interactions. Based on this observation, in follow-up studies we penalize the strength of all non-generalizable interactions during DNN training. This method encourages the DNN to encode more generalizable interactions.

We trained VGG-16 (Simonyan & Zisserman, 2015) and AlexNet (Krizhevsky et al., 2012) on the TinyImageNet dataset (Le & Yang, 2015), both with and without the proposed penalty.

As shown in Figure 6, DNNs trained with this penalty maintained their performance under higher pruning ratios compared to DNNs trained without it.

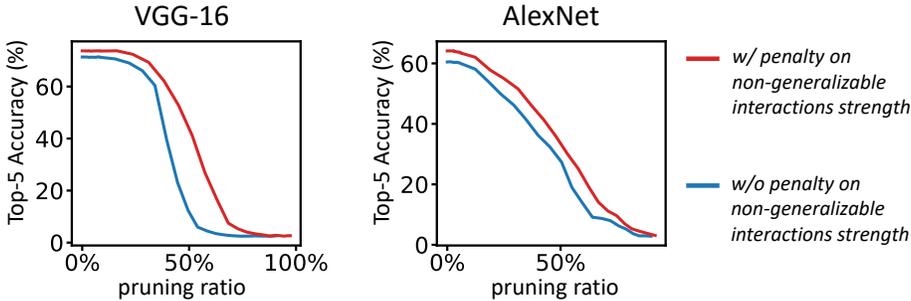


Figure 6: Top-5 classification accuracy of pruned VGG-16 and pruned AlexNet on TinyImageNet under different pruning ratio.

## C PROPERTIES OF THE AND INTERACTION

The Harsanyi interaction Harsanyi (1963) (referred to as the AND interaction in this work) has been a conventional metric for measuring the effect of the AND relationship that a DNN encodes among input variables. In this section, we introduce several desirable axioms that the AND interaction  $I_T^{\text{and}}$  adheres to. These properties further underscore the reliability of using AND interactions to explain the inference score of a DNN.

(1) *Efficiency axiom* (proven by Harsanyi (1963)). The output score of a model can be decomposed into interaction effects of different patterns, i.e.  $v(\mathbf{x}) = \sum_{T \subseteq N} I_T^{\text{and}}$ .

(2) *Linearity axiom*. If we merge output scores of two models  $v_1$  and  $v_2$  as the output of model  $v$ , i.e.  $\forall S \subseteq N, v(\mathbf{x}_S) = v_1(\mathbf{x}_S) + v_2(\mathbf{x}_S)$ , then their interaction effects  $I_{T,v_1}^{\text{and}}$  and  $I_{T,v_2}^{\text{and}}$  can also be merged as  $\forall T \subseteq N, I_{T,v}^{\text{and}} = I_{T,v_1}^{\text{and}} + I_{T,v_2}^{\text{and}}$ .

(3) *Dummy axiom*. If a variable  $i \in N$  is a dummy variable, i.e.  $\forall S \subseteq N \setminus \{i\}, v(\mathbf{x}_{S \cup \{i\}}) = v(\mathbf{x}_S) + v(\mathbf{x}_{\{i\}})$ , then it has no interaction with other variables,  $\forall \emptyset \neq T \subseteq N \setminus \{i\}, I_{T \cup \{i\}}^{\text{and}} = 0$ .

(4) *Symmetry axiom.* If input variables  $i, j \in N$  cooperate with other variables in the same way,  $\forall S \subseteq N \setminus \{i, j\}, v(\mathbf{x}_{S \cup \{i\}}) = v(\mathbf{x}_{S \cup \{j\}})$ , then they have same interaction effects with other variables,  $\forall T \subseteq N \setminus \{i, j\}, I_{T \cup \{i\}}^{\text{and}} = I_{T \cup \{j\}}^{\text{and}}$ .

(5) *Anonymity axiom.* For any permutations  $\pi$  on  $N$ , we have  $\forall T \subseteq N, I_{T, v}^{\text{and}} = I_{\pi T, \pi v}^{\text{and}}$ , where  $\pi T \stackrel{\text{def}}{=} \{\pi(i) | i \in T\}$ , and the new model  $\pi v$  is defined by  $(\pi v)(\mathbf{x}_{\pi S}) = v(\mathbf{x}_S)$ . This indicates that interaction effects are not changed by permutation.

(6) *Recursive axiom.* The interaction effects can be computed recursively. For  $i \in N$  and  $T \subseteq N \setminus \{i\}$ , the interaction effect of the pattern  $T \cup \{i\}$  is equal to the interaction effect of  $T$  with the presence of  $i$  minus the interaction effect of  $T$  with the absence of  $i$ , i.e.  $\forall T \subseteq N \setminus \{i\}, I_{T \cup \{i\}}^{\text{and}} = I_{T, i \text{ present}}^{\text{and}} - I_T^{\text{and}}$ .  $I_{T, i \text{ present}}^{\text{and}}$  denotes the interaction effect when the variable  $i$  is always present as a constant context, i.e.  $I_{T, i \text{ present}}^{\text{and}} = \sum_{L \subseteq T} (-1)^{|T| - |L|} \cdot v(\mathbf{x}_{L \cup \{i\}})$ .

(7) *Interaction distribution axiom.* This axiom characterizes how interactions are distributed for “interaction functions” Sundararajan et al. (2020). An interaction function  $v_T$  parameterized by a subset of variables  $T$  is defined as follows.  $\forall S \subseteq N$ , if  $T \subseteq S, v_T(\mathbf{x}_S) = c$ ; otherwise,  $v_T(\mathbf{x}_S) = 0$ . The function  $v_T$  models pure interaction among the variables in  $T$ , because only if all variables in  $T$  are present, the output value will be increased by  $c$ . The interactions encoded in the function  $v_T$  satisfies  $I_T^{\text{and}} = c$ , and  $\forall S \neq T, I_S^{\text{and}} = 0$ .

## D DIFFERENT PROPERTIES OF LOW-ORDER AND HIGH-ORDER INTERACTIONS

As Figure 3 and Figure 26 show, interactions of different orders exhibit distinct properties.

(1) Among all low-order interactions (including interactions of the 1st and 3rd orders), there are more low-order interactions with positive effects than those with negative effects, suggesting that low-order interactions serve as the primary inference patterns to boost the classification confidence.  
 (2) Most high-order interactions, which range from the 4th to the  $n$ -th order, usually exhibit mutually offsetting effects, i.e., a half of the high-order interactions increase the classification confidence, while the other half decrease the classification confidence.

## E COMMON CONDITIONS FOR SPARSE INTERACTIONS

Ren et al. (2024a) have proved three sufficient conditions for the sparsity of AND interactions.

**Condition 1.** *The DNN does not encode extremely high-order interactions:*  $\forall T \in \{T \subseteq N \mid |T| \geq M + 1\}, I_T^{\text{and}} = 0$ .

Condition 1 is common because extremely high-order interactions usually represent very complex and over-fitted patterns, which are unlikely to be learned by a well-trained DNN in real scenarios.

**Condition 2.** *Let  $\bar{u}^{(k)} \stackrel{\text{def}}{=} \mathbb{E}_{|S|=k} [v(\mathbf{x}_S) - v(\mathbf{x}_\emptyset)]$  denote the average classification confidence of the DNN over all masked samples  $\mathbf{x}_S$  with  $k$  unmasked input variables. This average classification confidence monotonically increases when  $k$  increases:  $\forall k' \leq k, \bar{u}^{(k')} \leq \bar{u}^{(k)}$ .*

Condition 2 implies that a well-trained DNN is likely to have higher average classification confidence for less masked input samples.

**Condition 3.** *Given the average classification confidence  $\bar{u}^{(k)}$  of samples with  $k$  unmasked input variables, there is a polynomial lower bound for the average classification confidence with  $k'$  ( $k' \leq k$ ) unmasked input variables:  $\forall k' \leq k, \bar{u}^{(k')} \geq (\frac{k'}{k})^p \bar{u}^{(k)}$ , where  $p > 0$  is a constant.*

Condition 3 suggests that the classification confidence of the DNN remains relatively stable even when presented with masked input samples. In real-world applications, the classification or detection of masked or occluded samples frequently occurs. As a result, a well-trained DNN typically develops the ability to classify such masked inputs by leveraging local information, which can be derived from the visible portions of the input. Consequently, the model should not produce a substantially reduced confidence score for masked samples.

## F DETAILS OF EXTRACTING THE SPARSEST AND-OR INTERACTIONS

A method is proposed Chen et al. (2024); Li & Zhang (2023a) to simultaneously extract AND interactions  $I_T^{\text{and}}$  and OR interactions  $I_T^{\text{or}}$  from the network output. Given a masked sample  $\mathbf{x}_L$ , Li & Zhang (2023a) proposed to learn a decomposition  $v(\mathbf{x}_L) = u_L^{\text{and}} + u_L^{\text{or}}$  towards the sparsest interactions. The component  $u_L^{\text{and}}$  was explained by AND interactions, and the component  $u_L^{\text{or}}$  was explained by OR interactions. Specifically, they decomposed  $v(\mathbf{x}_L)$  into  $u_L^{\text{and}} = 0.5 \cdot v(\mathbf{x}_L) + \gamma_L$  and  $u_L^{\text{or}} = 0.5 \cdot v(\mathbf{x}_L) - \gamma_L$ , where  $\{\gamma_L : L \subseteq N\}$  is a set of learnable variables that determine the decomposition. In this way, the AND interactions and OR interactions can be computed according to Theorem 2.1, i.e.,  $I_T^{\text{and}} = \sum_{L \subseteq T} (-1)^{|T|-|L|} u_L^{\text{and}}$ , and  $I_T^{\text{or}} = -\sum_{L \subseteq T} (-1)^{|T|-|L|} u_L^{\text{or}}$ .

The parameters  $\{\gamma_L\}$  were learned by minimizing the following LASSO-like loss to obtain sparse interactions:

$$\min_{\{\gamma_L\}} \sum_{T \subseteq N} |I_T^{\text{and}}| + |I_T^{\text{or}}| \quad (12)$$

The following pseudocode 1 outlines the core procedure of extracting AND-OR interactions.

---

### Algorithm 1 Computing AND-OR interactions

---

- 1: **Input:** Input sample  $\mathbf{x}$ , the DNN  $v(\cdot)$
  - 2: **Output:** A set of interactions  $I_S^{\text{AND}}$  and  $I_S^{\text{OR}}$
  - 3: **for**  $S \subseteq N$  **do**
  - 4:   For each masked sample  $\mathbf{x}_S$ , compute the confidence score  $v(\mathbf{x}_S)$  based on Eq. (1);
  - 5: **end for**
  - 6: **for**  $S \subseteq N$  **do**
  - 7:   Given  $v(\mathbf{x}_S)$  for all combinations  $S \subseteq N$ , compute each AND interaction effect  $I_S^{\text{AND}}$  and each OR interaction effect  $I_S^{\text{OR}}$  via  $\min_{\{\gamma_T\}} \sum_{S \subseteq N, S \neq \emptyset} [|I_S^{\text{AND}}| + |I_S^{\text{OR}}|]$ ;
  - 8: **end for**
  - 9: **return**  $I_S^{\text{AND}}, I_S^{\text{OR}}$
- 

Furthermore, to extract interactions shared between the original DNN  $v(\mathbf{x})$  and the pruned DNN  $v'(\mathbf{x})$ , we modified the objective as:

$$\begin{aligned} \min_{\{\gamma_L\}} \sum_{T \subseteq N} & \left( |I_T^{\text{and}}| + |I_T^{\text{or}}| + |I_T^{\text{and}'}| + |I_T^{\text{or}'}| \right) \\ & + \sum_{T \subseteq N} \left( |\min(0, I_T^{\text{and}}, I_T^{\text{and}'})| + |\max(0, I_T^{\text{and}}, I_T^{\text{and}'})| \right) \\ & + |\min(0, I_T^{\text{or}}, I_T^{\text{or}'})| + |\max(0, I_T^{\text{or}}, I_T^{\text{or}'})| \end{aligned} \quad (13)$$

**Removing small noises.** A small noise  $\delta$  in the network output may significantly affect the extracted interactions, especially for high-order interactions. Thus, Li et al. Li & Zhang (2023a) proposed to learn to remove a small noise term  $\delta_T$  from the computation of AND-OR interactions. Specifically, the decomposition was rewritten as  $u_L^{\text{and}} = 0.5(v(\mathbf{x}_L) - \delta_L) + \gamma_L$  and  $u_L^{\text{or}} = 0.5(v(\mathbf{x}_L) - \delta_L) + \gamma_L$ . Thus, the parameters  $\{\delta_L\}$  and  $\{\gamma_L\}$  are simultaneously learned by minimizing the loss function in Eq. (12). The values of  $\{\delta_L\}$  were constrained in  $[-\zeta, \zeta]$  where  $\zeta = 0.02 \cdot |v(\mathbf{x}) - v(\mathbf{x}_\emptyset)|$ .

## G PROOF OF THEOREM 2.1

### Proof G.1 (1) Universal matching property of AND interactions.

We will prove that output component  $u_S^{\text{AND}}$  on all  $2^n$  masked samples  $\{\mathbf{x}_S : S \subseteq N\}$  could be universally explained by the all interactions in  $S \subseteq N$ , i.e.,  $\forall \emptyset \neq S \subseteq N, u_S^{\text{AND}} = \sum_{\emptyset \neq T \subseteq S} I_T^{\text{AND}} + v(\mathbf{x}_\emptyset)$ . In particular, we define  $u_\emptyset^{\text{AND}} = v(\mathbf{x}_\emptyset)$  (i.e., we attribute output on an empty sample to AND interactions).

Specifically, the AND interaction is defined as  $I_T^{AND} = \sum_{L \subseteq T} (-1)^{|T|-|L|} u_L^{AND}$ . To compute the sum of AND interactions  $\sum_{\emptyset \neq T \subseteq S} I_T^{AND} = \sum_{\emptyset \neq T \subseteq S} \sum_{L \subseteq T} (-1)^{|T|-|L|} u_L^{AND}$ , we first exchange the order of summation of the set  $L \subseteq T \subseteq S$  and the set  $T \supseteq L$ . That is, we compute all linear combinations of all sets  $T$  containing  $L$  with respect to the model outputs  $u_L^{AND}$  given a set of input variables  $L$ , i.e.,  $\sum_{T: L \subseteq T \subseteq S} (-1)^{|T|-|L|} u_L^{AND}$ . Then, we compute all summations over the set  $L \subseteq S$ .

In this way, we can compute them separately for different cases of  $L \subseteq T \subseteq S$ . In the following, we consider the cases (1)  $L = S = T$ , and (2)  $L \subseteq T \subseteq S, L \neq S$ , respectively.

(1) When  $L = S = T$ , the linear combination of all subsets  $T$  containing  $L$  with respect to the model output  $u_L^{AND}$  is  $(-1)^{|S|-|S|} u_L^{AND} = u_L^{AND}$ .

(2) When  $L \subseteq T \subseteq S, L \neq S$ , the linear combination of all subsets  $T$  containing  $L$  with respect to the model output  $u_L^{AND}$  is  $\sum_{T: L \subseteq T \subseteq S} (-1)^{|T|-|L|} u_L^{AND}$ . For all sets  $T: S \supseteq T \supseteq L$ , let us consider the linear combinations of all sets  $T$  with number  $|T|$  for the model output  $u_L^{AND}$ , respectively. Let  $m := |T| - |L|$ , ( $0 \leq m \leq |S| - |L|$ ), then there are a total of  $C_{|S|-|L|}^m$  combinations of all sets  $T$  of order  $|T|$ . Thus, given  $L$ , accumulating the model outputs  $u_L^{AND}$  corresponding to all  $T \supseteq L$ , then  $\sum_{T: L \subseteq T \subseteq S} (-1)^{|T|-|L|} u_L^{AND} = u_L^{AND} \cdot \underbrace{\sum_{m=0}^{|S|-|L|} C_{|S|-|L|}^m (-1)^m}_{=0} = 0$ . Please see the complete

derivation of the following formula.

$$\begin{aligned}
\sum_{\emptyset \neq T \subseteq S} I_T^{AND} &= \sum_{\emptyset \neq T \subseteq S} \sum_{L \subseteq T} (-1)^{|T|-|L|} u_L^{AND} \\
&= \sum_{L \subseteq S} \sum_{T: L \subseteq T \subseteq S} (-1)^{|T|-|L|} u_L^{AND} - u_{\emptyset}^{AND} \\
&= \underbrace{u_S^{AND}}_{L=S} + \sum_{L \subseteq S, L \neq S} u_L^{AND} \cdot \underbrace{\sum_{m=0}^{|S|-|L|} C_{|S|-|L|}^m (-1)^m}_{=0} - u_{\emptyset}^{AND} \\
&= u_S^{AND} - u_{\emptyset}^{AND} = u_S^{AND} - v(\mathbf{x}_{\emptyset})
\end{aligned} \tag{14}$$

Thus, we have  $\forall \emptyset \neq S \subseteq N, u_S^{and} = \sum_{\emptyset \neq T \subseteq S} I_T^{and} + v(\mathbf{x}_{\emptyset})$ .

## (2) Universal matching property of OR interactions.

According to the definition of OR interactions, we will derive that  $\forall S \subseteq N, u_S^{OR} = \sum_{T: T \cap S \neq \emptyset} I_T^{OR}$ , where we define  $u_{\emptyset}^{OR} = 0$  (recall that in Step (1), we attribute the output on empty input to AND interactions).

Specifically, the OR interaction is defined as  $I_T^{OR} = -\sum_{L \subseteq T} (-1)^{|T|-|L|} u_{N \setminus L}^{OR}$ . Similar to the above derivation of the universal matching theorem of AND interactions, to compute the sum of OR interactions  $\sum_{T: T \cap S \neq \emptyset} I_T^{OR} = \sum_{T: T \cap S \neq \emptyset} \left[ -\sum_{L \subseteq T} (-1)^{|T|-|L|} u_{N \setminus L}^{OR} \right]$ , we first exchange the order of summation of the set  $L \subseteq T \subseteq N$  and the set  $T: T \cap S \neq \emptyset$ . That is, we compute all linear combinations of all sets  $T$  containing  $L$  with respect to the model outputs  $u_{N \setminus L}^{OR}$  given a set of input variables  $L$ , i.e.,  $\sum_{T: T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR}$ . Then, we compute all summations over the set  $L \subseteq N$ .

In this way, we can compute them separately for different cases of  $L \subseteq T \subseteq N, T \cap S \neq \emptyset$ . In the following, we consider the cases (1)  $L = N \setminus S$ , (2)  $L = N$ , (3)  $L \cap S \neq \emptyset, L \neq N$ , and (4)  $L \cap S = \emptyset, L \neq N \setminus S$ , respectively.

(1) When  $L = N \setminus S$ , the linear combination of all subsets  $T$  containing  $L$  with respect to the model output  $u_{N \setminus L}^{OR}$  is  $\sum_{T: T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR} = \sum_{T: T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_S^{OR}$ . For all sets  $T: T \supseteq L, T \cap S \neq \emptyset$  (then  $T \neq N \setminus S, T \neq L$ ), let us consider the linear combinations of all sets  $T$  with number  $|T|$  for the model output  $u_S^{OR}$ , respectively. Let  $|T'| := |T| - |L|$ , ( $1 \leq |T'| \leq |S|$ ), then there are a total of  $C_{|S|}^{|T'|}$  combinations of all sets  $T'$  of order  $|T'|$ . Thus, given  $L$ , accumulating

864 the model outputs  $u_S^{OR}$  corresponding to all  $T \supseteq L$ , then  $\sum_{T:T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR} = u_S^{OR}$ .  
 865  
 866 
$$\underbrace{\sum_{|T'|=1}^{|S|} C_{|S|}^{|T'|} (-1)^{|T'|}}_{=-1} = -u_S^{OR}.$$
  
 867  
 868

869 (2) When  $L = N$  (then  $T = N$ ), the linear combination of all subsets  $T$  containing  $L$  with respect  
 870 to the model output  $u_{N \setminus L}^{OR}$  is  $\sum_{T:T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR} = (-1)^{|N|-|N|} u_{\emptyset}^{OR} = u_{\emptyset}^{OR}$ .

871  
 872 (3) When  $L \cap S \neq \emptyset, L \neq N$ , the linear combination of all subsets  $T$  containing  $L$  with respect to  
 873 the model output  $u_{N \setminus L}^{OR}$  is  $\sum_{T:T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR}$ . For all sets  $T : T \supseteq L, T \cap S \neq \emptyset$ ,  
 874 let us consider the linear combinations of all sets  $T$  with number  $|T|$  for the model output  $u_S^{OR}$ ,  
 875 respectively. Let us split  $|T| - |L|$  into  $|T'|$  and  $|T''|$ , i.e.,  $|T| - |L| = |T'| + |T''|$ , where  $T' =$   
 876  $\{i | i \in T, i \notin L, i \in N \setminus S\}$ ,  $T'' = \{i | i \in T, i \notin L, i \in S\}$  (then  $0 \leq |T''| \leq |S| - |S \cap L|$ ) and  
 877  $|T'| + |T''| + |L| = |T|$ . In this way, there are a total of  $C_{|S|-|S \cap L|}^{|T''|}$  combinations of all sets  $T''$  of  
 878 order  $|T''|$ . Thus, given  $L$ , accumulating the model outputs  $u_{N \setminus L}^{OR}$  corresponding to all  $T \supseteq L$ , then  
 879 
$$\sum_{T:T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR} = u_{N \setminus L}^{OR} \cdot \sum_{T' \subseteq N \setminus S \setminus L} \underbrace{\sum_{|T''|=0}^{|S|-|S \cap L|} C_{|S|-|S \cap L|}^{|T''|} (-1)^{|T'+|T''|}}_{=0} =$$
  
 880  
 881  
 882 0.

883  
 884 (4) When  $L \cap S = \emptyset, L \neq N \setminus S$ , the linear combination of all subsets  $T$  containing  $L$  with respect  
 885 to the model output  $u_{N \setminus L}^{OR}$  is  $\sum_{T:T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR}$ . Similarly, let us split  $|T| - |L|$   
 886 into  $|T'|$  and  $|T''|$ , i.e.,  $|T| - |L| = |T'| + |T''|$ , where  $T' = \{i | i \in T, i \notin L, i \in N \setminus S\}$ ,  
 887  $T'' = \{i | i \in T, i \in S\}$  (then  $0 \leq |T''| \leq |S|$ ) and  $|T'| + |T''| + |L| = |T|$ . In this way,  
 888 there are a total of  $C_{|S|}^{|T''|}$  combinations of all sets  $T''$  of order  $|T''|$ . Thus, given  $L$ , accumulating  
 889 the model outputs  $u_{N \setminus L}^{OR}$  corresponding to all  $T \supseteq L$ , then  $\sum_{T:T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR} =$   
 890 
$$u_{N \setminus L}^{OR} \cdot \sum_{T' \subseteq N \setminus S \setminus L} \underbrace{\sum_{|T''|=0}^{|S|} C_{|S|}^{|T''|} (-1)^{|T'+|T''|}}_{=0} = 0.$$
  
 891  
 892  
 893

894 Please see the complete derivation of the following formula.

895  
 896  
 897  
 898  
 899  
 900  
 901  
 902  
 903  
 904  
 905  
 906  
 907  
 908  
 909  
 910  
 911  
 912  
 913  
 914  
 915  
 916  
 917

$$\begin{aligned}
 \sum_{T:T \cap S \neq \emptyset} I_T^{OR} &= \sum_{T:T \cap S \neq \emptyset} \left[ - \sum_{L \subseteq T} (-1)^{|T|-|L|} u_{N \setminus L}^{OR} \right] \\
 &= - \sum_{L \subseteq N} \sum_{T:T \cap S \neq \emptyset, T \supseteq L} (-1)^{|T|-|L|} u_{N \setminus L}^{OR} \\
 &= - \left[ \sum_{|T'|=1}^{|S|} C_{|S|}^{|T'|} (-1)^{|T'|} \right] \cdot \underbrace{u_S^{OR}}_{L=N \setminus S} - \underbrace{u_{\emptyset}^{OR}}_{L=N} \\
 &\quad - \sum_{L \cap S \neq \emptyset, L \neq N} \left[ \sum_{T' \subseteq N \setminus S \setminus L} \left( \sum_{|T''|=0}^{|S|-|S \cap L|} C_{|S|-|S \cap L|}^{|T''|} (-1)^{|T'+|T''|} \right) \right] \cdot u_{N \setminus L}^{OR} \\
 &\quad - \sum_{L \cap S = \emptyset, L \neq N \setminus S} \left[ \sum_{T' \subseteq N \setminus S \setminus L} \left( \sum_{|T''|=0}^{|S|} C_{|S|}^{|T''|} (-1)^{|T'+|T''|} \right) \right] \cdot u_{N \setminus L}^{OR} \tag{15} \\
 &= -(-1) \cdot u_S^{OR} - u_{\emptyset}^{OR} - \sum_{L \cap S \neq \emptyset, L \neq N} \left[ \sum_{T' \subseteq N \setminus S \setminus L} 0 \right] \cdot u_{N \setminus L}^{OR} \\
 &\quad - \sum_{L \cap S = \emptyset, L \neq N \setminus S} \left[ \sum_{T' \subseteq N \setminus S \setminus L} 0 \right] \cdot u_{N \setminus L}^{OR} \\
 &= u_S^{OR} - u_{\emptyset}^{OR} \\
 &= u_S^{OR}
 \end{aligned}$$

(3) Universal matching property of AND-OR interactions.

918 *With the universal matching theorem of AND interactions and the universal matching theorem of OR interac-*  
 919 *tions, we can easily get  $v(\mathbf{x}_S) = u_S^{and} + u_S^{or} = v(\mathbf{x}_\emptyset) + \sum_{\emptyset \neq T \subseteq S} I_T^{and} + \sum_{T: T \cap S \neq \emptyset} I_T^{or}$ , thus, we obtain the*  
 920 *universal matching theorem of AND-OR interactions.*

## 922 H EXPERIMENTAL DETAIL

### 924 H.1 TRAINING SETTINGS

925  
 926 In this paper, we followed Fang et al. (2023) and trained VGG-19 Simonyan & Zisserman (2015)  
 927 and ResNet-56 He et al. (2016) on the CIFAR-10 dataset Krizhevsky (2009) for 200 epochs using  
 928 stochastic gradient descent (SGD) with a learning rate of 0.1, momentum of 0.9, and a weight decay  
 929 of  $5 \times 10^{-4}$ . The learning rate was decayed by a factor of 0.1 at epochs 120, 150, and 180. We used  
 930 a batch size of 128 in all experiments.

931 All experiments were conducted on a compute node equipped with dual Intel Xeon Silver 4310  
 932 CPUs (48 logical cores) and a combination of four NVIDIA A800 and two NVIDIA A100 GPUs  
 933 (each with 80GB memory).  
 934

### 935 H.2 PRUNING SETTINGS

936  
 937 Considering the adaptability of pruning algorithms to models, (1) we applied LLM-Pruner (Ma et al.,  
 938 2023) to prune DeepSeek-R1-Distill-LLaMA-8B model (DeepSeek-AI et al., 2025) and LLaMA-  
 939 3.1-8B model (Grattafiori et al., 2024); (2) we applied Isomorphic-Pruning (Fang et al., 2024) to  
 940 prune ResNet-50, ResNet-101 (He et al., 2016; Wightman et al., 2021), ViT-Small and ViT-Base  
 941 (Dosovitskiy et al., 2021; Steiner et al., 2022) pretrained on the ImageNet-1k dataset (Russakovsky  
 942 et al., 2015); (3) we used DepGraph (Fang et al., 2023) to prune VGG-19, ResNet-56 on the CIFAR-  
 943 10 dataset (Krizhevsky, 2009).

944 All other pruning settings followed the original settings in Ma et al. (2023); Fang et al. (2024; 2023),  
 945 respectively. No post-pruning finetuning was performed.  
 946

### 947 H.3 DETAILS ABOUT HOW TO CALCULATE INTERACTIONS FOR DIFFERENT DNNs

948  
 949 • **For experiments on image classification models**, since the computational cost of interactions  
 950 was intolerable, we applied a sampling-based approximation method to calculate AND-OR interac-  
 951 tions. Specifically, we considered the feature map after the low-layer as intermediate-layer features  
 952 of DNNs. We uniformly split each intermediate-layer feature map into  $s \times s$  patches, and sampled  
 953 10 patches based on the highest significance values computed via gradient integration. For VGG-  
 954 19 and ResNet-56 on CIFAR-10,  $s = 8$ . For ResNet-50, ResNet-50, ViT-Small and ViT-Base on  
 955 ImageNet-1k,  $s = 7$ . These selected patches were considered as input variables for the correspond-  
 956 ing intermediate-layer feature. We used a specific feature vector as the baseline value to mask the  
 957 variables in  $N \setminus T$ . The duration of the experiments ranges from 4 to 8 hours.

958 • **For experiments on large language models**, we manually collected 100 sentences for interaction  
 959 extraction, which cover topics such as corporate restructuring, science and technology, economic  
 960 policy, global trade, geopolitics, public health, political developments, cryptocurrency collapse,  
 961 monetary policy, and climate and energy equity. The zero-shot perplexity (PPL) analysis is per-  
 962 formed on WikiText2 dataset (Merity et al., 2017). We considered the outputs of the low-layer  
 963 corresponding to input words as input features. We considered the embeddings corresponding to  
 964 input features as input variables for each input sentence, and we randomly sampled 10 words, which  
 965 must have a specific meaning and not be stop words, to calculate interactions. We used the average  
 966 embedding over different input variables to mask the tokens in  $N \setminus T$ . We used a specific feature  
 967 vector as the baseline value to mask the variables in  $N \setminus T$ . The duration of the experiments ranges  
 968 from 6 to 8 hours.

969 Specifically, we empirically considered the first 4 convolutional layer of VGG-19 as the low layers,  
 970 and considered all the other 13 layers<sup>3</sup> as the high layers. For ResNet-56, we consider the first 3  
 971

<sup>3</sup>VGG-19 for CIFAR-10 Krizhevsky (2009) has been simplified compared to the original VGG-19 archi-  
 tecture designed for ImageNet, by replacing the multi-layer perceptron (MLP) classifier with a single fully

convolution layers as low layers. For ResNet-50 and ResNet-101, we consider the first 11 convolution layers as low layers. For AlexNet, ViT-Small and ViT-Base, we consider the first convolution layer as low layers. Typically, we compute the mean distribution of interactions over 100 samples for ResNet-56 on the CIFAR-10 dataset. We compute the mean distribution of interactions over 100 samples for VGG-19 on the CIFAR-10 dataset. We compute the mean distribution of interactions over 100 samples for DeepSeek-R1-Distill-LLaMA-8B on the above manually collected data.

#### H.4 SETTINGS OF REFERENCE DNNs

We used VGG-19 and ResNet-56 as each other’s reference model. For ResNet-50 and ResNet-101, we used the ViT-Base model as reference DNN. For ViT-Small and ViT-Base, we used the ResNet-101 model as reference DNN. We used the Qwen-2.5-7B (Qwen et al., 2024) as the reference DNN for the DeepSeek-R1-Distill-LLaMA-8B and LLaMA-3.1-8B model.

## I ILLUSTRATION OF AND-OR LOGICAL MODEL

The following four figures show the logical models that mathematically represent the inference logic of DeepSeek-R1-Distill-LLaMA-8B model and Qwen-2.5-7B on two input prompts, respectively.

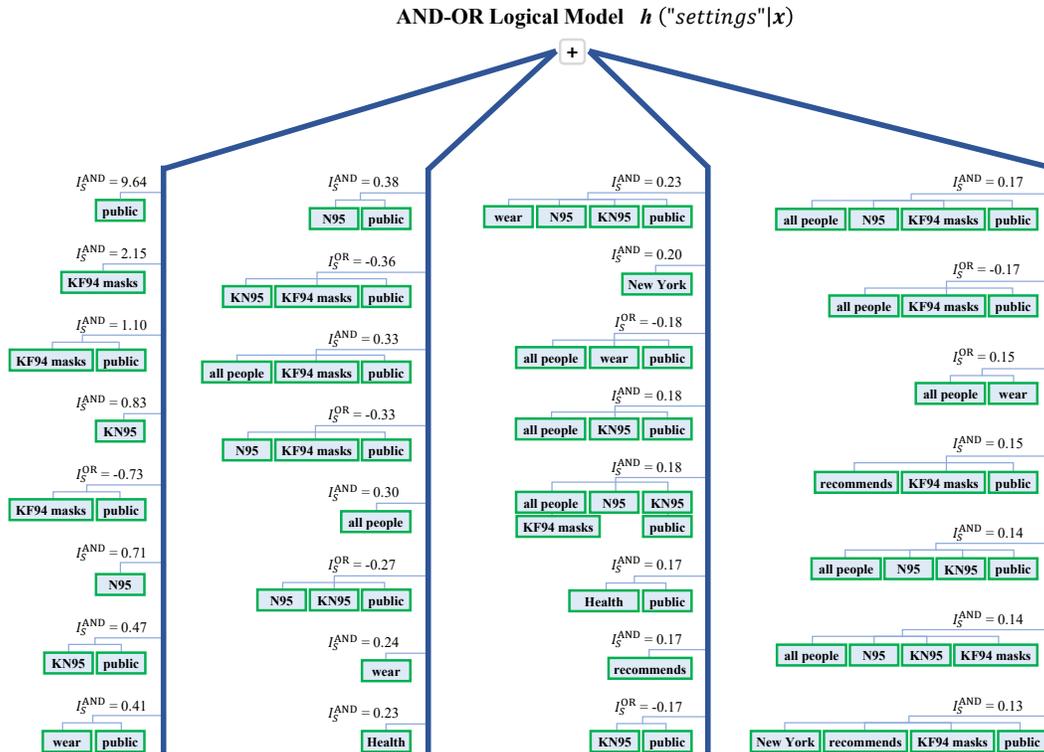
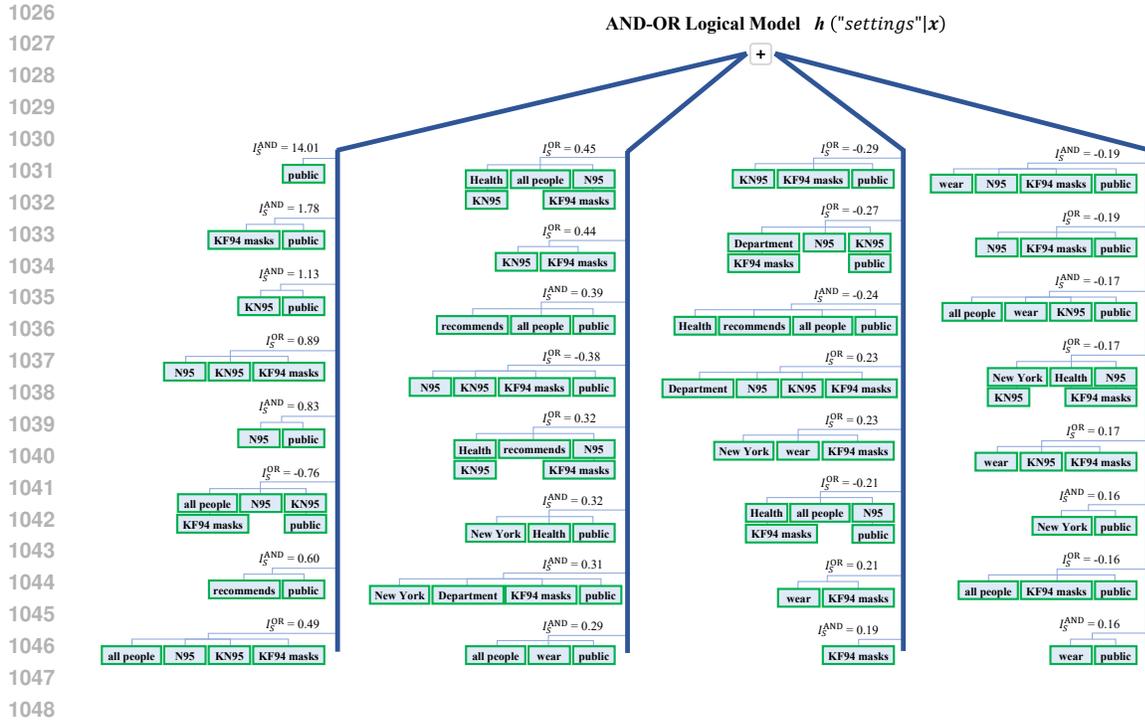
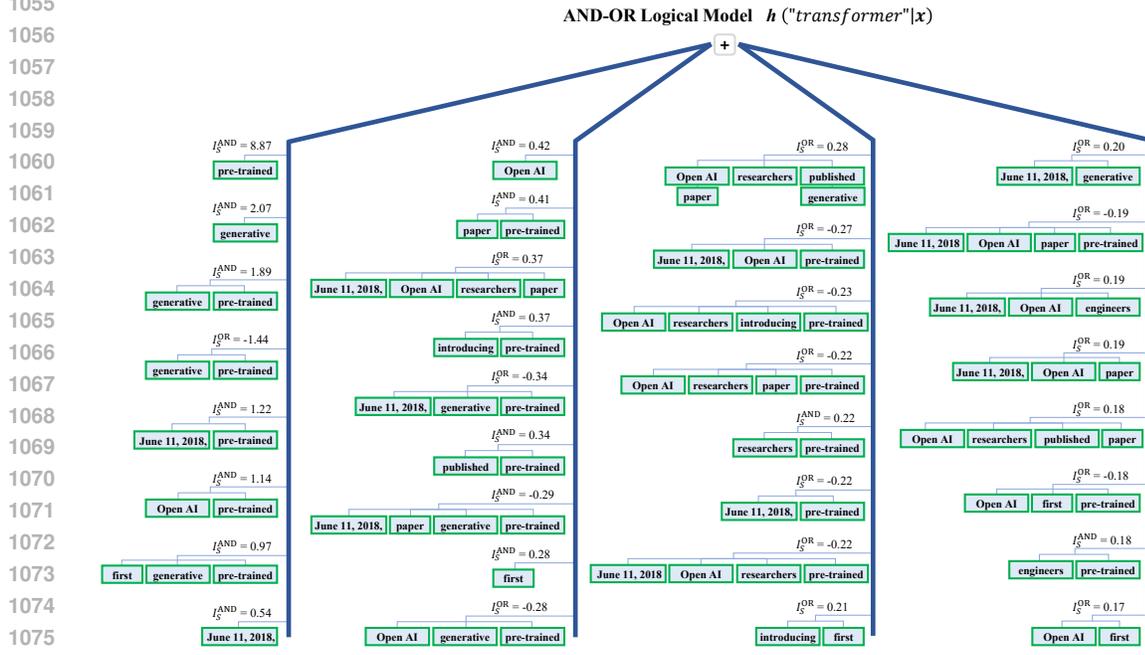


Figure 7: The logical model representing the inference logic of the DeepSeek-R1-Distill-LLaMA-8B model on the input prompt “New York Department of Health recommends that all people should wear N95, KN95, or KF94 masks in all public.” The predicted next word is “settings.”

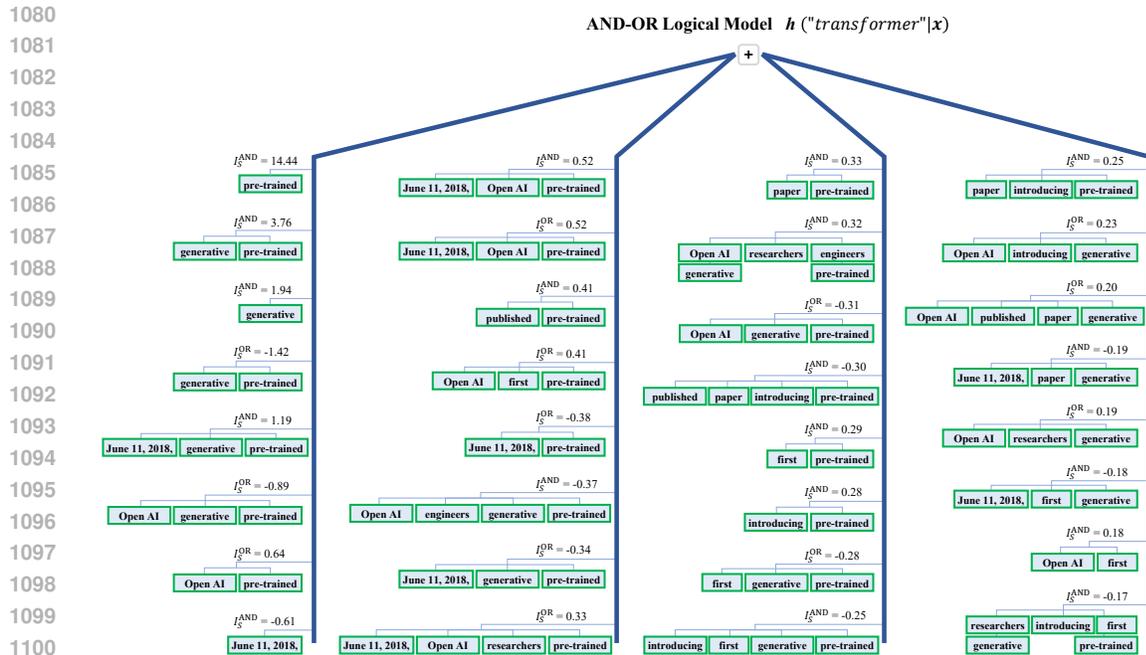
connected layer and adapting the input resolution. The total number of layers, including both convolutional and fully connected layers, remains 17 in the simplified VGG-19 model for CIFAR-10.



1049 Figure 8: The logical model representing the inference logic of the Qwen-2.5-7B model on the input  
1050 prompt “New York Department of Health recommends that all people should wear N95, KN95, or  
1051 KF94 masks in all public.” The predicted next word is “settings.”



1078 Figure 9: The logical model representing the inference logic of the DeepSeek-R1-Distill-LLaMA-  
1079 8B model on the input prompt “On June 11, 2018, OpenAI researchers and engineers published a  
paper introducing the first generative pre-trained.” The predicted next word is “transformer.”



1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131

Figure 10: The logical model representing the inference logic of the Qwen-2.5-7B model on the input prompt “On June 11, 2018, OpenAI researchers and engineers published a paper introducing the first generative pre-trained.” The predicted next word is “transformer.”

1132  
1133

The following four figures show the logical models that mathematically represent the inference logic of DeepSeek-R1-Distill-LLaMA-8B model (before and after pruning operation) on two input prompts, respectively.



1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

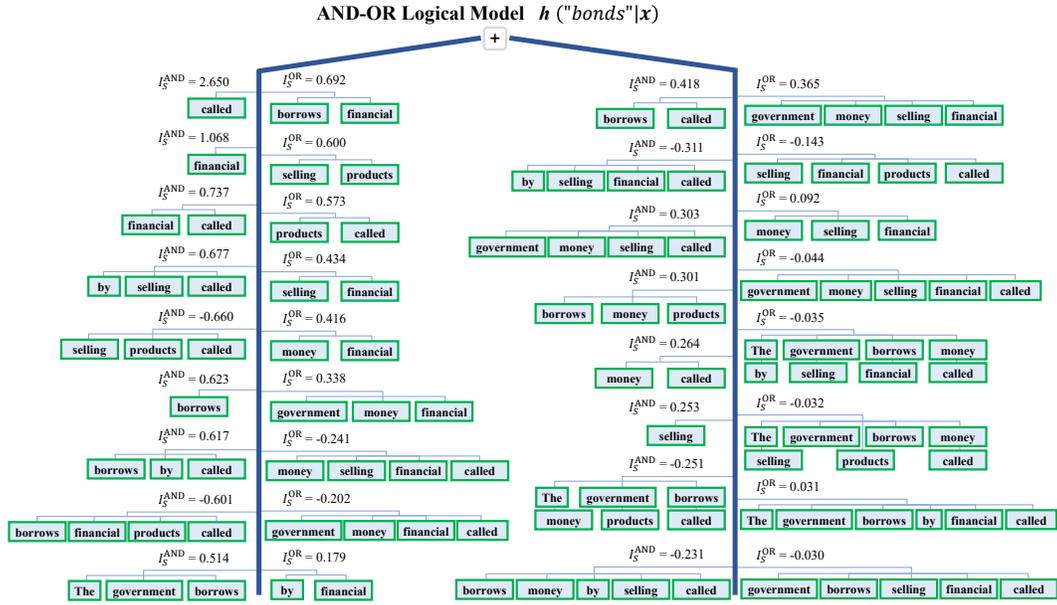


Figure 13: The logical model representing the inference logic of the original DeepSeek-R1-Distill-LLaMA-8B model on the input prompt “The government borrows money by selling financial products called.” The predicted next word is “bonds.”

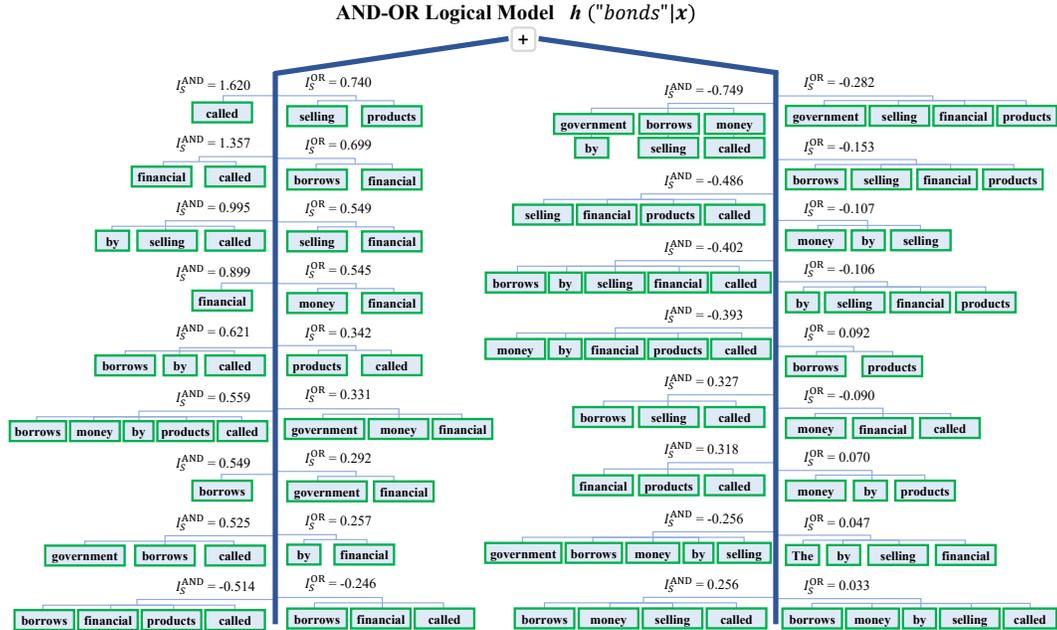


Figure 14: The logical model representing the inference logic of the pruned DeepSeek-R1-Distill-LLaMA-8B model (under a pruning ratio of 41.01) on the input prompt “The government borrows money by selling financial products called.” The predicted next word is “bonds.”

## J DETAILED RESULTS OF THE THREE-PHASE DYNAMICS OF INTERACTIONS

### J.1 DETAILED RESULTS ON VGG-19

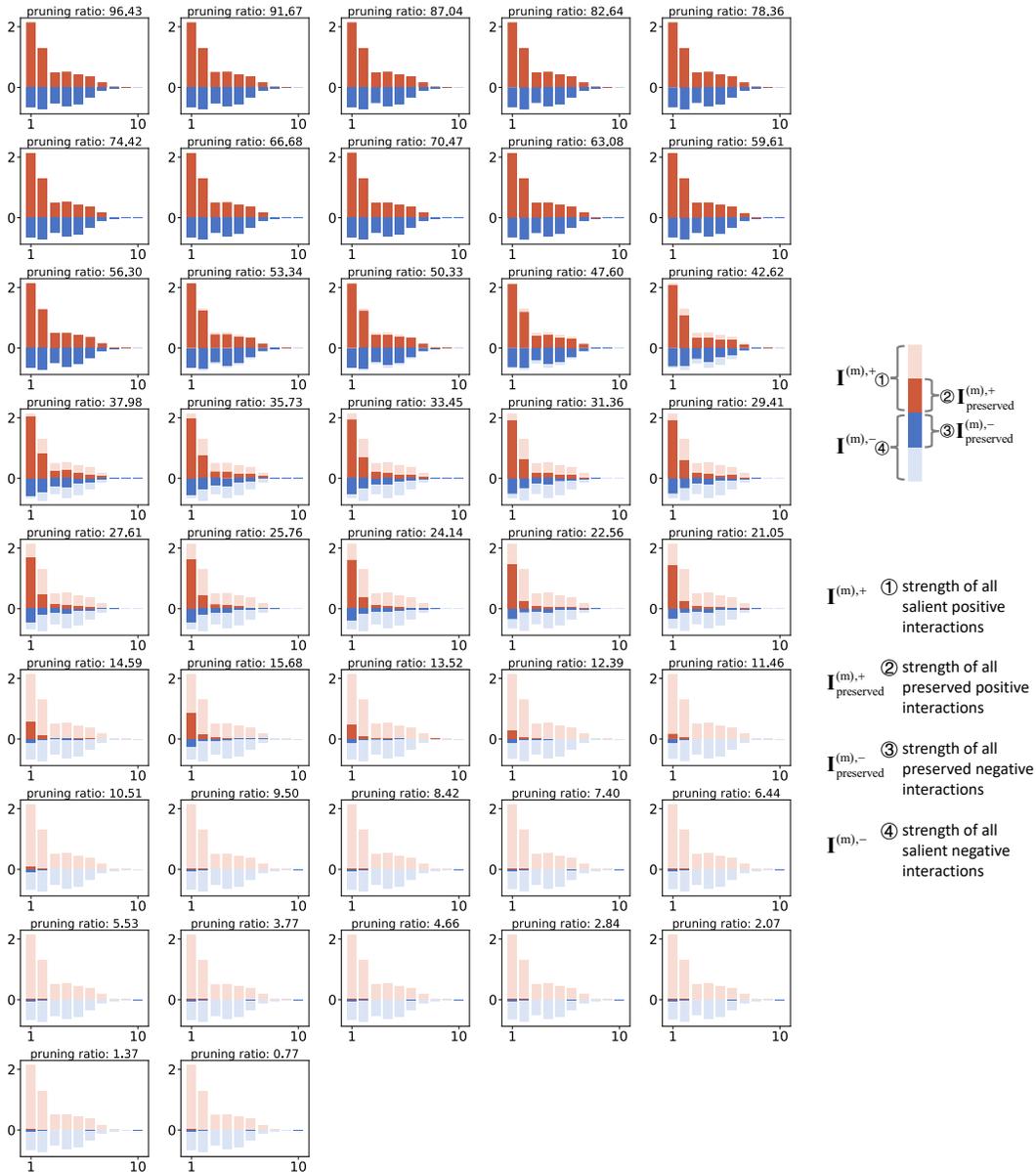


Figure 15: The distributions of interactions and preserved interactions across different orders encoded by pruned VGG-19 at different pruning ratios.

## J.2 DETAILED RESULTS ON RESNET-56

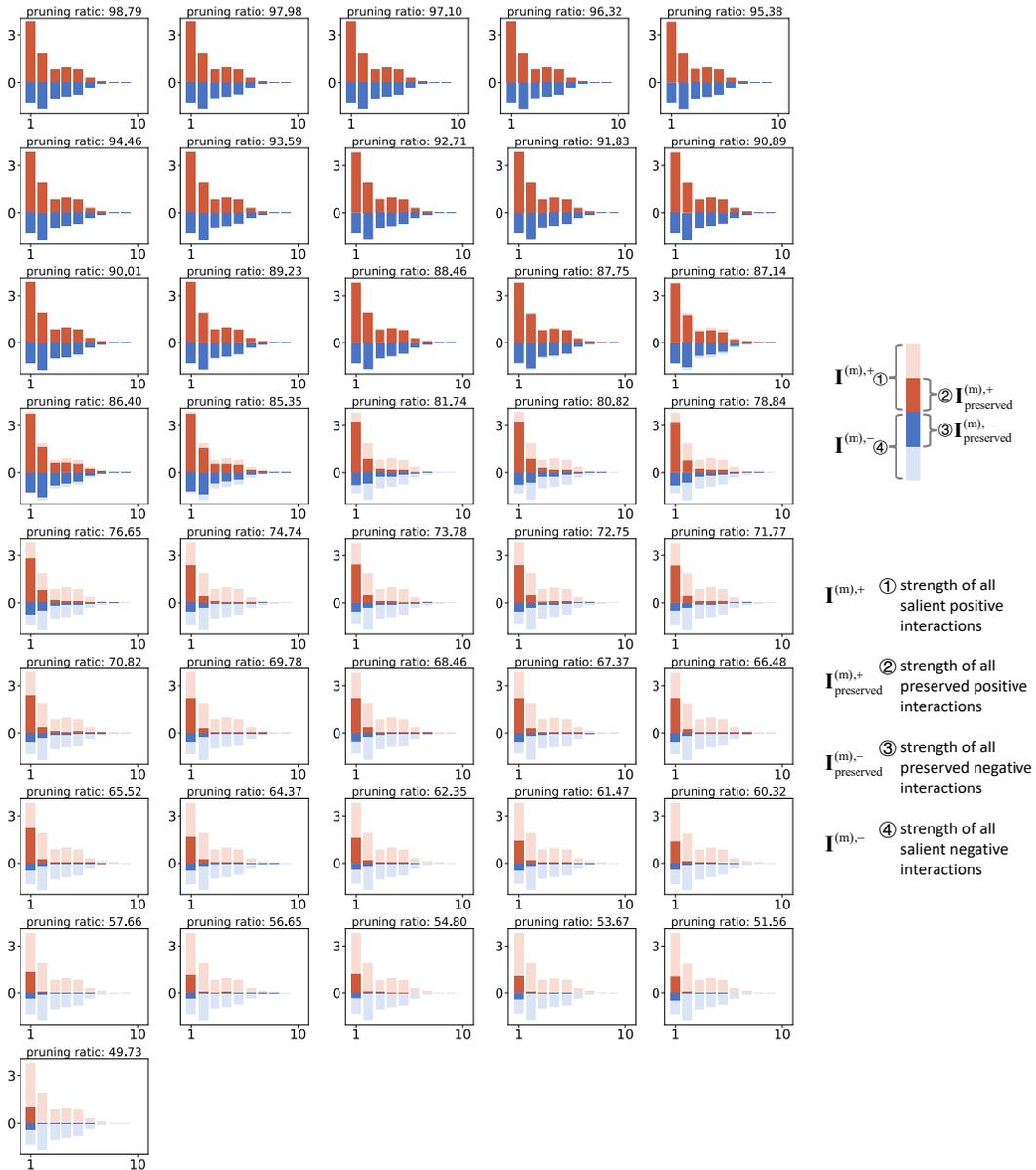


Figure 16: The distributions of interactions and preserved interactions across different orders encoded by pruned ResNet-56 at different pruning ratios.

J.3 DETAILED RESULTS ON RESNET-50

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

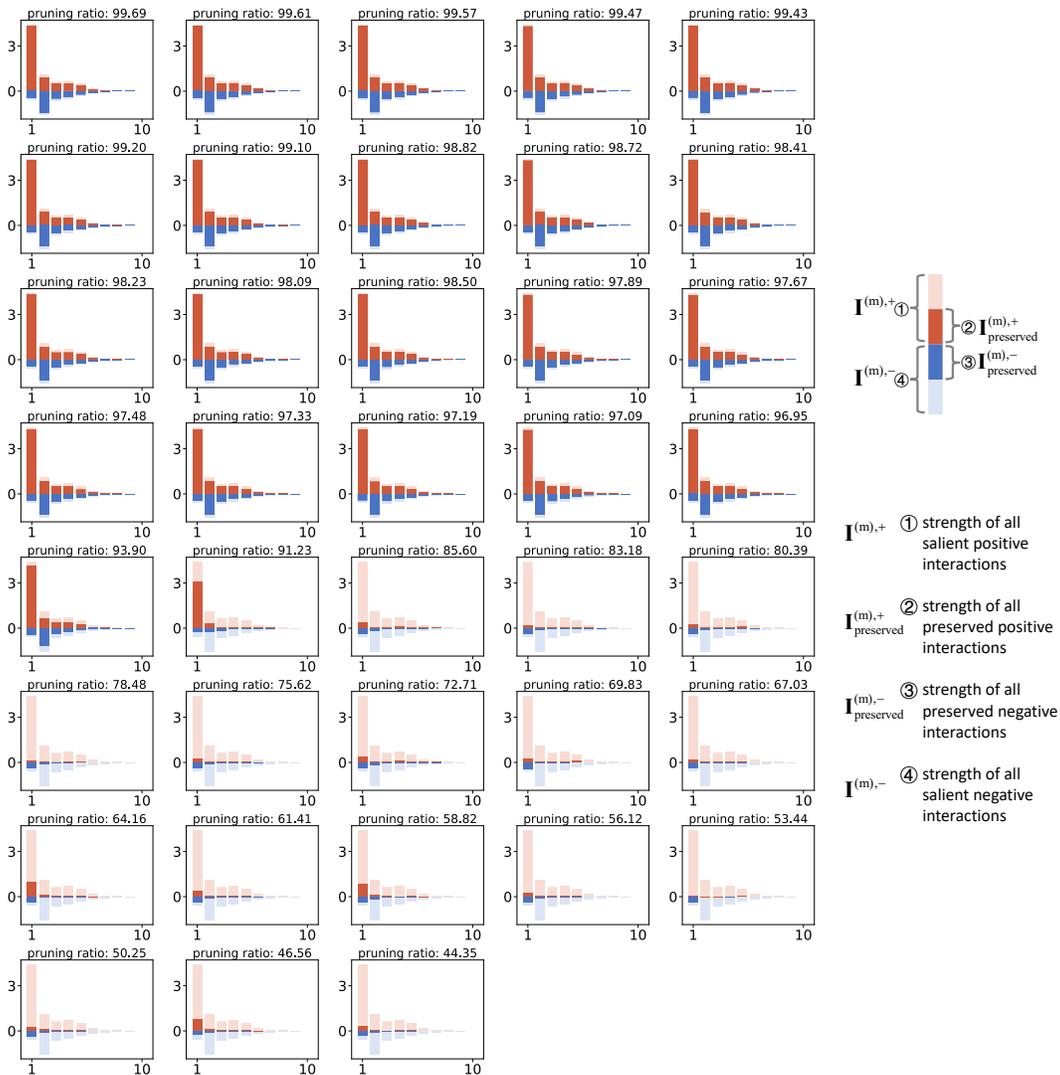


Figure 17: The distributions of interactions and preserved interactions across different orders encoded by pruned ResNet-50 at different pruning ratios.

## J.4 DETAILED RESULTS ON RESNET-101

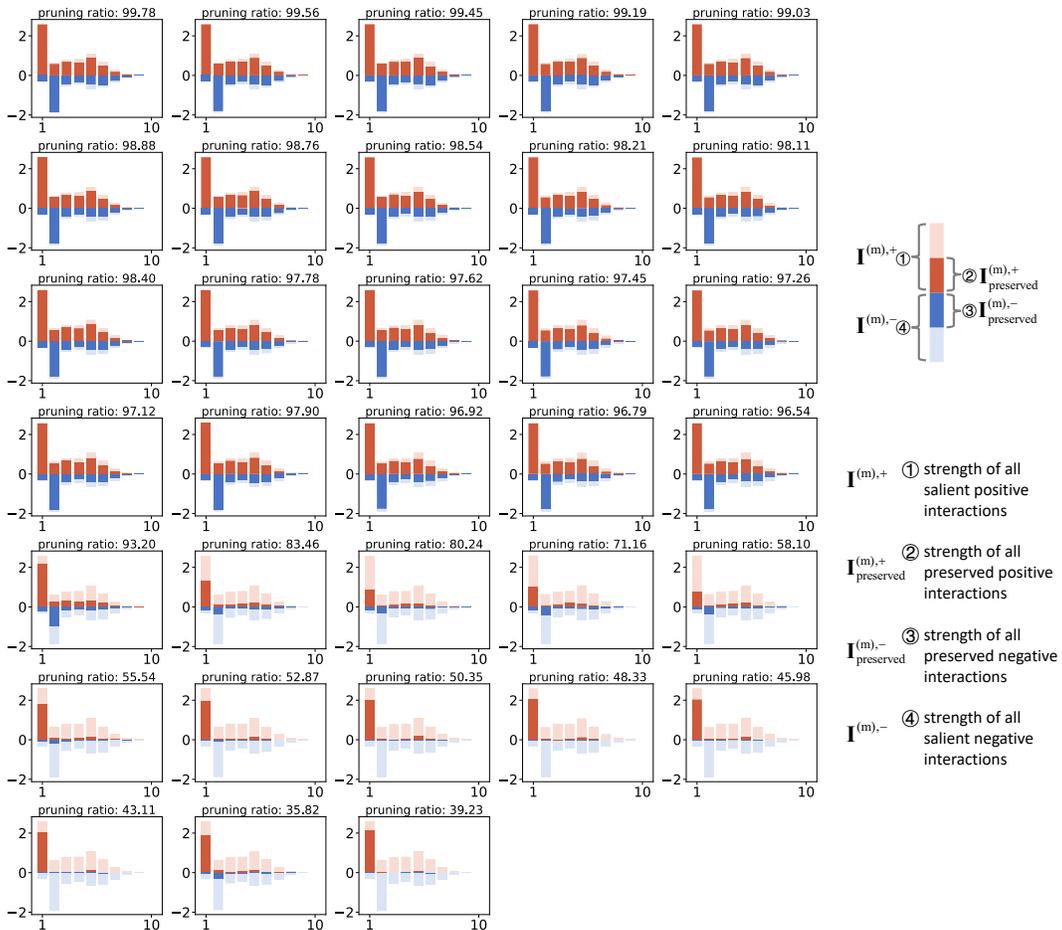


Figure 18: The distributions of interactions and preserved interactions across different orders encoded by pruned ResNet-101 at different pruning ratios.

1458 J.5 DETAILED RESULTS ON ViT-SMALL  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478

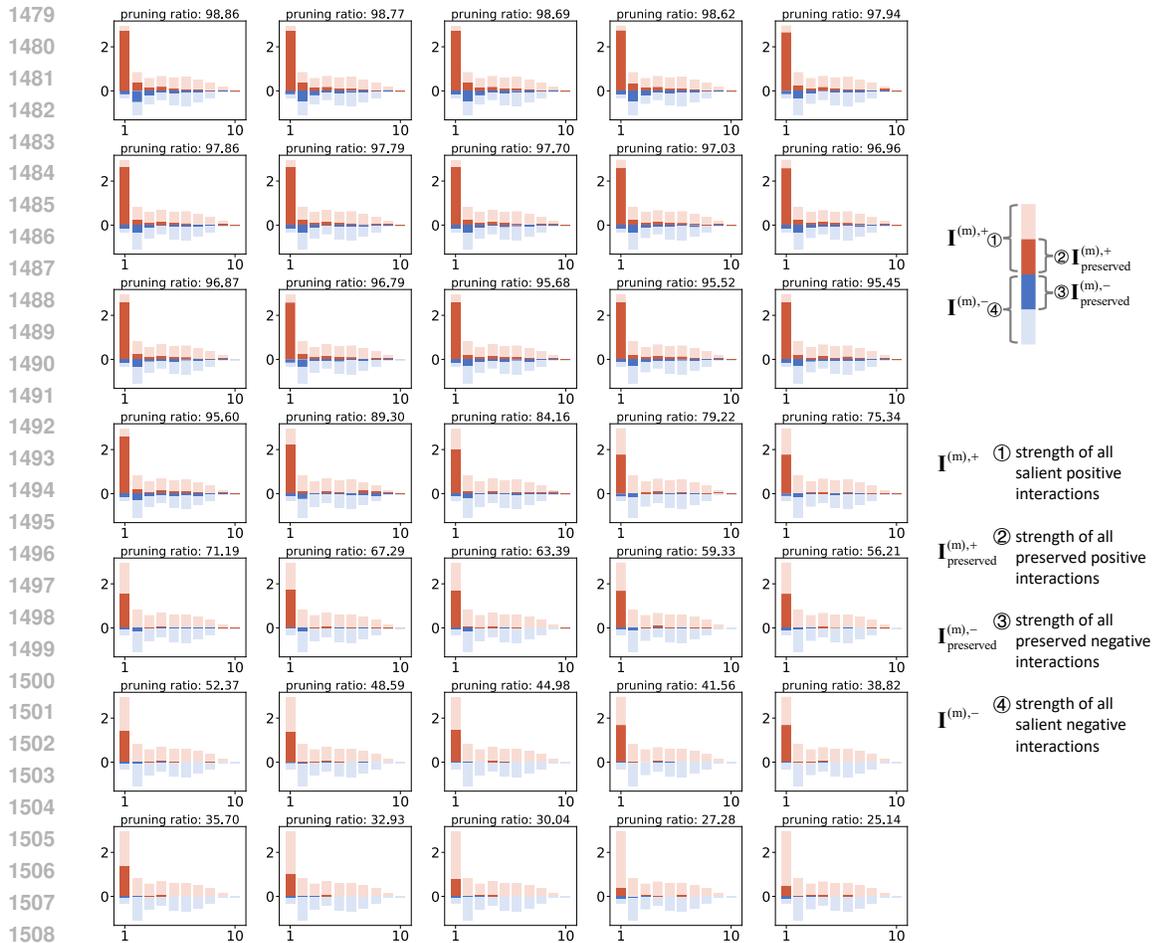


Figure 19: The distributions of interactions and preserved interactions across different orders encoded by pruned ViT-Small at different pruning ratios.

J.6 DETAILED RESULTS ON ViT-BASE

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

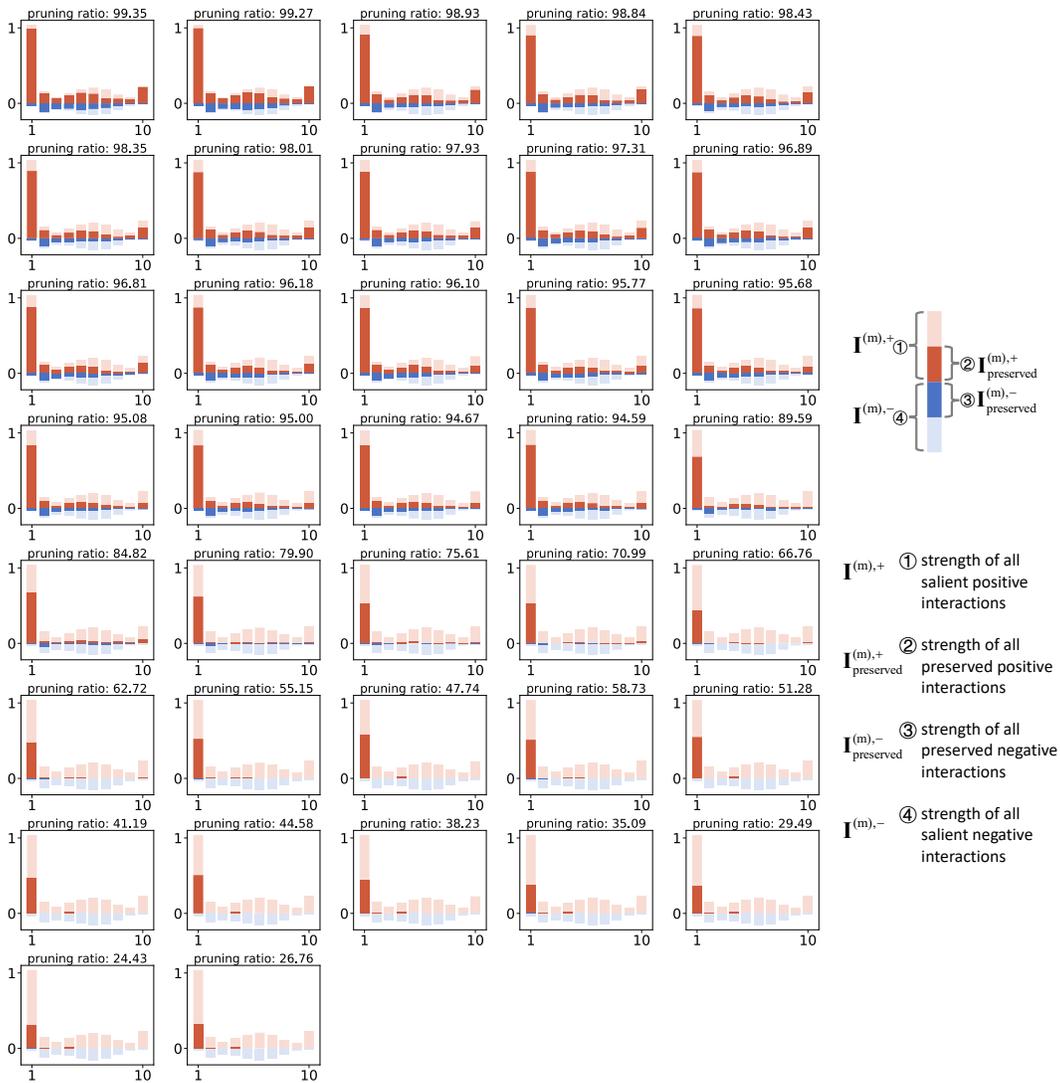
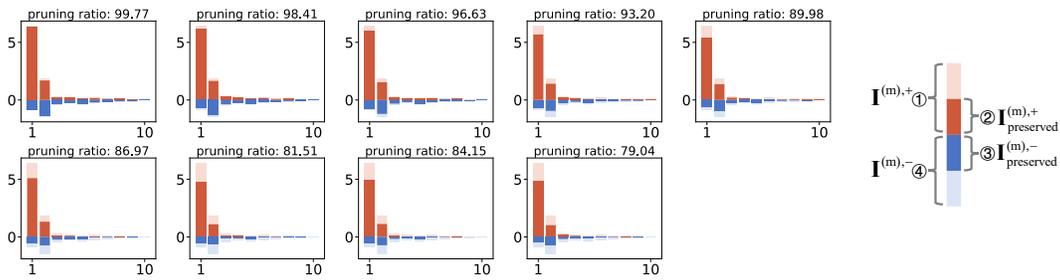


Figure 20: The distributions of interactions and preserved interactions across different orders encoded by pruned ViT-BASE at different pruning ratios.

J.7 DETAILED RESULTS ON DEEPSEEK-R1-DISTILL-LLAMA-8B

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576



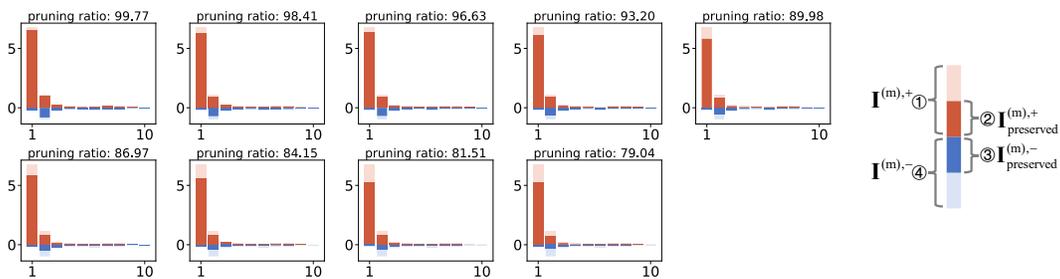
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586

Figure 21: The distributions of interactions and preserved interactions across different orders encoded by pruned DeepSeek-R1-Distill-LLaMA-8B at different pruning ratios.

1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596

J.8 DETAILED RESULTS ON LLAMA-3.1-8B

1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607



1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617

Figure 22: The distributions of interactions and preserved interactions across different orders encoded by pruned LLaMA-3.1-8B at different pruning ratios.

1618  
1619

## J.9 RESULTS ON ALL EIGHT DNNs

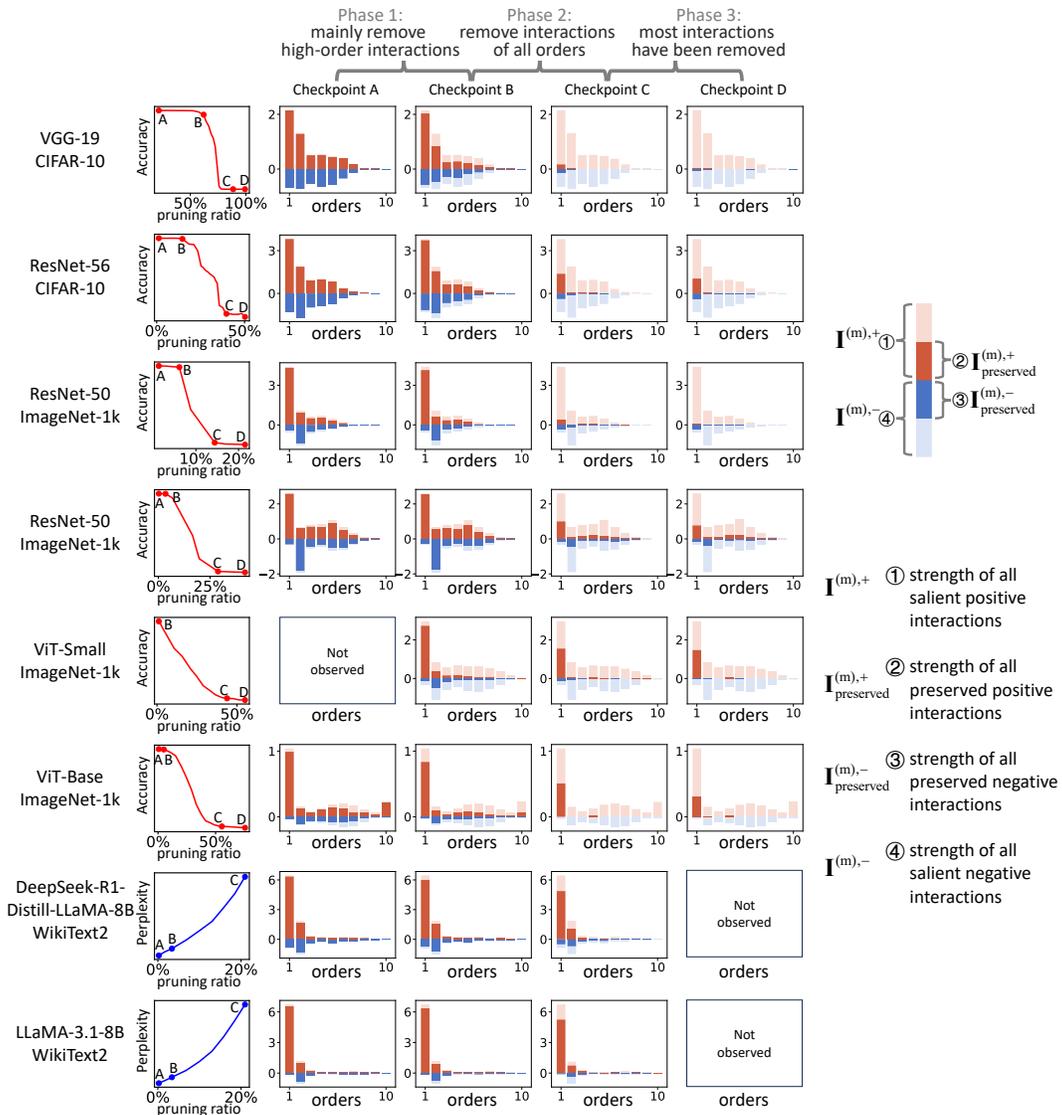


Figure 23: Changes of the distributions of  $I^{(m),+}$ ,  $I^{(m),-}$ ,  $I_{\text{preserved}}^{(m),+}$ , and  $I_{\text{preserved}}^{(m),-}$  on eight DNNs when the pruning ratio increases. These changes can be divided into three phases.

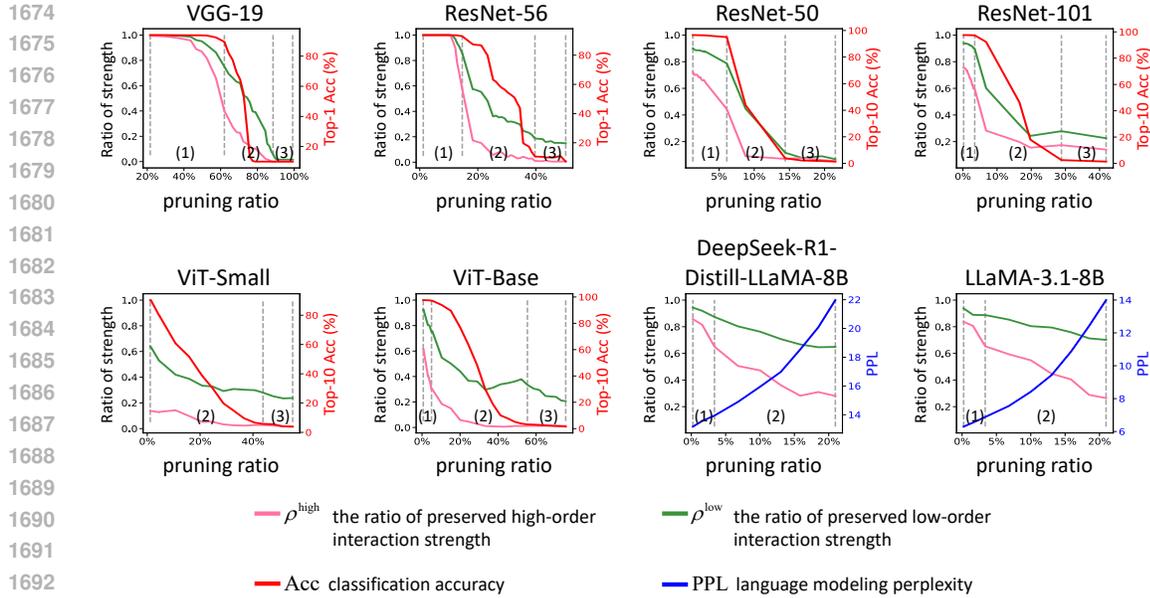


Figure 24: Changes of the ratio of preserved low-order interaction strength  $\rho^{\text{low}}$ , the ratio of preserved high-order interaction strength  $\rho^{\text{high}}$  and classification accuracy (or language modeling perplexity) on eight DNNs when the pruning ratio increases.

Table 3: Efficiency of parameter pruning in the three-phases.

Architecture	Phase 1	Phase 2	Phase 3
ResNet-56	$e^{1.21 \rightarrow 14.65} = 0.92$	$e^{14.65 \rightarrow 39.68} = 0.78$	$e^{39.68 \rightarrow 50.27} = 0.38$
VGG-19	$e^{21.64 \rightarrow 62.02} = 0.89$	$e^{62.02 \rightarrow 88.54} = 0.61$	$e^{88.54 \rightarrow 99.23} = 0.36$
ResNet-50	$e^{1.28 \rightarrow 6.10} = 0.91$	$e^{6.10 \rightarrow 14.40} = 0.61$	$e^{14.40 \rightarrow 21.52} = 0.72$
ResNet-101	$e^{0.22 \rightarrow 3.46} = 0.93$	$e^{3.46 \rightarrow 28.84} = 0.75$	$e^{28.84 \rightarrow 41.90} = 0.70$
ViT-Small	—	$e^{1.14 \rightarrow 43.79} = 0.57$	$e^{43.79 \rightarrow 55.02} = 0.27$
ViT-Base	$e^{0.65 \rightarrow 5.00} = 0.87$	$e^{5.00 \rightarrow 55.42} = 0.69$	$e^{55.42 \rightarrow 75.57} = 0.55$
DeepSeek-R1-Distill-LLaMA-8B	$e^{0.23 \rightarrow 3.37} = 0.78$	$e^{3.37 \rightarrow 20.96} = 0.75$	—
LLaMA-3.1-8B	$e^{0.23 \rightarrow 3.37} = 0.74$	$e^{3.37 \rightarrow 20.96} = 0.70$	—

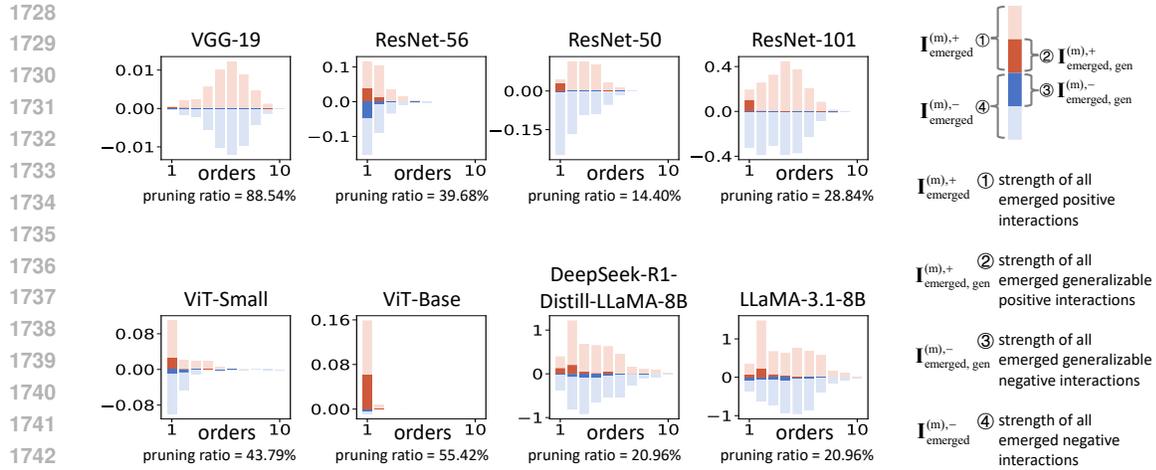


Figure 25: The distributions of the emerged interactions introduced by the parameter pruning operation  $\mathbf{I}_{\text{emerged}}^{(m),+}$ ,  $\mathbf{I}_{\text{emerged}}^{(m),-}$ ,  $\mathbf{I}_{\text{emerged,gen}}^{(m),+}$  and  $\mathbf{I}_{\text{emerged,gen}}^{(m),-}$  on eight pruned DNNs. Emerged interactions are barely generalizable.

## K DISTRIBUTIONS OF GENERALIZABLE INTERACTION

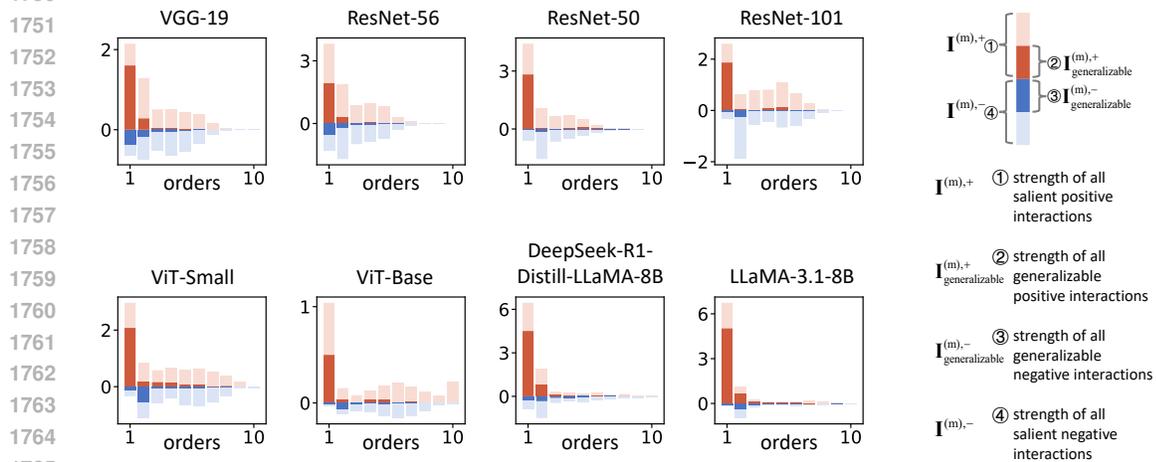


Figure 26: The distributions of interactions and generalizable interactions  $\mathbf{I}^{(m),+}$ ,  $\mathbf{I}^{(m),-}$ ,  $\mathbf{I}_{\text{generalizable}}^{(m),+}$  and  $\mathbf{I}_{\text{generalizable}}^{(m),-}$  encoded by eight original DNNs. Interactions of different orders exhibit different generalizability.

## L THE USE OF LARGE LANGUAGE MODELS (LLMs)

In this work, large language models (LLMs) were used solely as a general-purpose writing assistant to polish the grammar and improve the clarity of the text. No part of the research ideation, experiment design, data analysis, or substantive content generation relied on LLMs. The authors take full responsibility for the content of the paper.