

---

# Measuring Mechanistic Interpretability at Scale Without Humans

---

Roland S. Zimmermann<sup>1</sup> David Klindt<sup>2</sup> Wieland Brendel<sup>1</sup>

## Abstract

In today’s era, whatever we can measure at scale, we can optimize. So far, measuring the interpretability of units in deep neural networks (DNNs) for computer vision still requires direct human evaluation and is not scalable. As a result, the inner workings of DNNs remain a mystery despite the remarkable progress we have seen in their applications. In this work, we introduce the first scalable method to measure the per-unit interpretability in vision DNNs. This method does not require any human evaluations, yet its prediction correlates well with existing human interpretability measurements. We validate its predictive power through an interventional human psychophysics study. We demonstrate the usefulness of this measure by performing previously infeasible experiments: (1) A large-scale interpretability analysis across more than 70 million units from 835 computer vision models, and (2) an extensive analysis of how units transform during training. We find an anticorrelation between a model’s downstream classification performance and per-unit interpretability, which is also observable during model training. Furthermore, we see that a layer’s location and width influence its interpretability.

## 1. Introduction

With the arrival of the first non-trivial neural networks, researchers got interested in understanding their inner workings (Krizhevsky et al., 2012; Mahendran & Vedaldi, 2015). For one, this can be motivated by scientific curiosity; for another, a better understanding might lead to building more reliable, efficient, or fairer models. While the performance of machine learning models has seen a remarkable improvement over the last few years, our understanding of informa-

tion processing has progressed more slowly. Nevertheless, understanding how complex models — e.g., language models (Bricken et al., 2023) or vision models (Olah et al., 2017; Zimmermann et al., 2023) — work is still an active and growing field of research, coined *mechanistic interpretability* (Olah, 2022). A common approach in this field is to divide a network into atomic units, hoping they are easier to comprehend. Here, atomic units might refer to individual neurons or channels of (convolutional) layers (Olah et al., 2017), or general vectors in feature space (Elhage et al., 2022; Klindt et al., 2023). Besides this approach, mechanistic interpretability also includes the detection of neural circuits (Cammarata et al., 2020; Elhage et al., 2022) or analysis of global network properties (Nanda et al., 2023).

The goal of understanding the inner workings of a neural network is inherently human-centric: Irrespective of what tools have been used, in the end, humans should have a better comprehension of the network. However, human evaluations are time-consuming and costly due to their reliance on human labor (Zimmermann et al., 2023). This results in slower research progress, as validating novel hypotheses takes longer.

Removing the need for human labor by automating the interpretability evaluation can open up multiple high-impact research directions: One benefit is that it enables the creation of more interpretable networks by explicitly optimizing for interpretability — after all, what we can measure at scale, we can optimize. Moreover, it allows more efficient research on explanation methods and might lead to an increased overall understanding of neural networks. While efforts to build such measures for language models exist (Bills et al., 2023), there is no common approach yet for vision models.

The present work is the first to introduce a fully automated interpretability measure (see Fig. 1A & B): the Machine Interpretability Score (MIS). By leveraging the latest advances in image similarity functions aligned with human perception, we obtain a measure that is strongly predictive of human-perceived interpretability (see Fig. 1C). We verify our measure through both correlational and interventional experiments. By removing the need for human labor, we can scale existing evaluations up by multiple orders of magnitude. Finally, this work demonstrates potential workflows and use cases of our MIS.

---

<sup>1</sup>Max Planck Institute for Intelligent Systems, Tübingen AI Center, Tübingen, Germany <sup>2</sup>Stanford University, Stanford, USA. Correspondence to: Roland S. Zimmermann <[research@zimmermann.com](mailto:research@zimmermann.com)>.

Online version, code and interactive visualizations available at: [brendel-group.github.io/mis](https://brendel-group.github.io/mis)

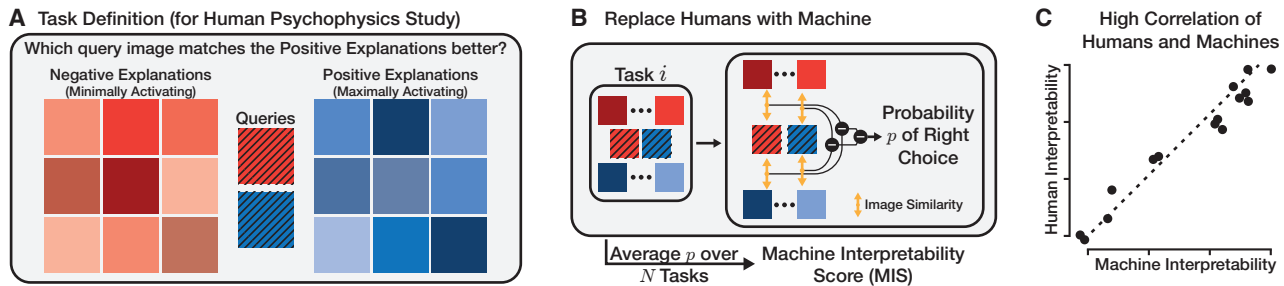


Fig. 1: **Definition of the Machine Interpretability Score.** **A.** We build on top of the established task definition proposed by Borowski et al. (2021) to quantify the per-unit interpretability via human psychophysics experiments. The task quantifies how well participants understand the sensitivity of a unit by asking them to match strongly activating query images to strongly activating *visual* explanations of the unit. See Fig. 13 for examples. **B.** Crucially, we remove the need for humans and fully automate the evaluation: We pass the explanations and query images through a feature encoder to compute pair-wise image similarities (DreamSim) before using a (hard-coded) binary classifier to solve the underlying task. Finally, the Machine Interpretability Score (MIS) is the average of the predicted probability of the correct choice over  $N$  tasks. **C.** The MIS proves to be highly correlated with human interpretability ratings and allows fast evaluations of new hypotheses.

## 2. Related Work

**Mechanistic Interpretability** While the overall field of explainable AI (XAI) tries to increase our understanding of neural networks, multiple subbranches with different foci exist (Gilpin et al., 2018). One of these branches, *mechanistic interpretability*, hopes to improve our understanding of neural networks by understanding their building blocks (Olah, 2022). An even more fine-grained branch aims to interpret individual units of vision models (Bau et al., 2017; Zhou et al., 2018; Bau et al., 2020; Morcos et al., 2018; Olah et al., 2017). We focus exclusively on this branch of research in the present work. This line of research for artificial neural networks was, arguably, inspired by similar efforts in neuroscience for biological neural networks (Hubel & Wiesel, 1962; Barlow, 1972; Quiroga et al., 2005).

Different studies set out to understand the behavior and sensitivity of individual units of vision networks – here, a unit can, e.g., be a channel in a convolutional neural network (CNN) or a neuron in a multilayer perceptron (MLP). With the recent progress in vision-language modeling, a few approaches started using textual descriptions of a unit’s behavior (Hernandez et al., 2022; Kalibhat et al., 2023). However, the majority still uses visual explanations which are either synthesized by performing activation maximization through, e.g., gradient ascent (Olah et al., 2017; Erhan et al., 2009; Mahendran & Vedaldi, 2015; Nguyen et al., 2014; Mordvintsev et al., 2015; Yosinski et al., 2015; Olah et al., 2017; Nguyen et al., 2017), or strongly activating dataset examples (Olah et al., 2017; Borowski et al., 2021).

With the onset of large language models (LLMs) and the increasing interest in them, there is also now an increasing interest in mechanistic interpretability of them (e.g., Elhage et al., 2021; Olsson et al., 2022; Bricken et al., 2023).

**Quantifying Interpretability** Rigorous evaluations, including falsifiable hypothesis testing, are critical for research on interpretability methods (Leavitt & Morcos, 2020). This also encompasses the need for human-centric evaluations (Borowski et al., 2021; Kim et al., 2022).

Nevertheless, such human-centric evaluations of interpretability methods are only available in some sub-fields. Specifically for the type of interpretability we are concerned about in this work, i.e., the per-unit interpretability of vision models, two methods for quantifying the helpfulness of explanations to humans were introduced before: Borowski et al. (2021) presented a two-alternative-forced-choice (2-AFC) psychophysics task that requires participants to determine which of two images elicits higher activation of the unit in question, given visual explanations (i.e., images that strongly activate or deactivate the unit, see Fig. 1A) of the unit’s behavior. Zimmermann et al. (2021) extended this paradigm to quantify how well participants can predict the influence of interventions in the form of occlusions in images. While these studies used their paradigms to evaluate the usefulness of different interpretability methods, Zimmermann et al. (2023) leveraged them to compare the interpretability of multiple models. Due to the reliance on human experiments, they could only probe the interpretability of 767 units from nine models. We now automatize this evaluation to scale it up by multiple orders of magnitude to more than 70 million units across 835 models.

**Automating Interpretability Research** To increase the efficiency of interpretability research and scale it to large modern-day networks, the concept of automated interpretability was proposed, first in the domain of natural language processing (Bills et al., 2023). This approach uses an LLM to generate textual descriptions of the behavior

of units in another LLM. Follow-up work by Huang et al. (2023), however, pointed out potential problems regarding the correctness of the explanations. To benchmark future fully automated interpretability tools, acting as independent agents, Schwettmann et al. (2023) introduced a synthetic benchmark suite inspired by the behavior of neural networks. In computer vision, there are also efforts to automate interpretability research (Hernandez et al., 2022; Zimmermann et al., 2023). Hernandez et al. (2022) and Oikarinen & Weng (2022) map visual to textual explanations of a unit’s behavior using automated tools, hoping to increase the efficiency of evaluations. Zimmermann et al. (2023) introduced the ImageNet Mechanistic Interpretability (IMI) dataset, containing per-unit interpretability annotations from humans for 767 units, meant to foster research on automating interpretability evaluations.

### 3. Method

We now introduce our fully automated interpretability measure, Machine Interpretability Score (MIS), visualized in Fig. 1. Borowski et al. (2021) proposed a setup that allows quantifying how well humans can infer the sensitivity of a unit in a vision model, e.g., a channel in a CNN, commonly averaged over space, or neuron in an MLP, from explanations: They leverage a 2-AFC task design in a psychophysics experiment (see left side of Fig. 1) to measure how well humans understand a unit by probing how well they can predict which of two extremely activating (query) images yields a higher activation, after seeing visual explanations. Specifically, two sets of explanations are displayed: highly and weakly activating images, called positive and negative explanations, respectively. See Appx. A.1 for a more detailed summary of the task. We build on top of this paradigm but replace human participants with machines, resulting in a fully automated interpretability metric that requires no humans.

**Definition of the Machine Interpretability Score** Let  $\mathcal{I}$  denote the space of valid input images for a model. For a specific explanation method and a unit in question, we denote the unit’s positive and negative visual explanations as sets of images  $\mathcal{E}^+ \subseteq \mathcal{I}$  and  $\mathcal{E}^- \subseteq \mathcal{I}$ , respectively. Further, let  $\mathcal{Q}^+ \subseteq \mathcal{I}$  and  $\mathcal{Q}^- \subseteq \mathcal{I}$  be the sets of query images with the most extreme (positive and negative) activations.

The task by Borowski et al. (2021) can now be expressed as: Given explanations  $\mathcal{E}^+$  and  $\mathcal{E}^-$  and two queries  $\mathbf{q}^+ \in \mathcal{Q}^+$  and  $\mathbf{q}^- \in \mathcal{Q}^-$ , which of the two queries matches  $\mathcal{E}^+$  and which  $\mathcal{E}^-$  more closely? An intuitive way to solve this binary decision task is to compare each query with every explanation and to match the query images to the sets of explanations based on the similarities of the images.

To formalize this, we introduce a perceptual (image) similar-

ity function  $f : \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$  computing the scalar similarity of two images (Zhang et al., 2018), and an aggregation function  $a : \mathbb{R}^K \rightarrow \mathbb{R}$  reducing a set of  $K$  similarities to a single one. This allows us to define the function  $s : \mathcal{I} \times \mathcal{I}^K \rightarrow \mathbb{R}$  that quantifies the similarity of a single query image to a set of explanations:

$$s(\mathbf{q}, \mathcal{E}) := a(\{f(\mathbf{q}, \mathbf{e}) \mid \mathbf{e} \in \mathcal{E}\}). \quad (1)$$

To decide whether a single query image is more likely to be the positive one, we can compute whether it is more similar to the positive than the negative explanations. We can compute this now for both the positive and the negative query images and get:

$$\Delta_+(\mathbf{q}^+, \mathcal{E}^+, \mathcal{E}^-) = s(\mathbf{q}^+, \mathcal{E}^+) - s(\mathbf{q}^+, \mathcal{E}^-), \quad (2)$$

$$\Delta_-(\mathbf{q}^-, \mathcal{E}^+, \mathcal{E}^-) = s(\mathbf{q}^-, \mathcal{E}^+) - s(\mathbf{q}^-, \mathcal{E}^-). \quad (3)$$

The classification problem will be solved correctly if the similarity of  $\mathbf{q}^+$  to  $\mathcal{E}^+$  relative to  $\mathcal{E}^-$  is stronger than those of  $\mathbf{q}^-$ . This means we can define the probability of solving the binary classification problem correctly as

$$p(\mathbf{q}^+, \mathbf{q}^-, \mathcal{E}^+, \mathcal{E}^-) := \sigma\left(\alpha \cdot (\Delta_+(\mathbf{q}^+, \mathcal{E}^+, \mathcal{E}^-) - \Delta_-(\mathbf{q}^-, \mathcal{E}^+, \mathcal{E}^-))\right) \quad (4)$$

where  $\sigma$  denotes the sigmoid function and  $\alpha$  is a free parameter to calibrate the classifier’s confidence.

We define the *Machine Interpretability Score* (MIS) as the predicted probability of making the right choice, averaged over  $N$  tasks for the same unit. Across these different tasks, the query images  $\mathbf{q}^+$ ,  $\mathbf{q}^-$  vary to cover a wider range of the unit’s behavior. If the explanation method used is stochastic, it is advisable to also average over different explanations:

$$\text{MIS} = \frac{1}{N} \sum_i^N p(\mathbf{q}_i^+, \mathbf{q}_i^-, \mathcal{E}_i^+, \mathcal{E}_i^-). \quad (5)$$

Note that the MIS is not a general property of a unit but depends on the method used for generating explanations. One might define a general score by computing the maximum MIS over multiple explanation methods.

**Choice of Hyperparameters.** We use the current state-of-the-art perceptual similarity, DreamSim (Fu et al., 2023), as  $f$ . See Appx. B for a sensitivity study on this choice. We use the mean to aggregate the distances between a query image and multiple explanations to a single scalar, i.e.,  $a(x_1, \dots, x_K) := 1/K \sum_i^K x_i$ . To choose  $\alpha$ , we use the interpretability annotations of IMI (Zimmermann et al., 2023): We optimize  $\alpha$  over a randomly chosen subset of just 5% of the annotated units to approximately match the value range of human interpretability scores, resulting in  $\alpha = 0.16$ .

Note that  $\alpha$  is, in fact, the only free parameter of our metric, resulting in very low chances of overfitting the metric to the IMI dataset. We use the same strategy as Borowski et al. (2021); Zimmermann et al. (2021) and Zimmermann et al. (2023) for generating new tasks (see Appx. A.2). As they used up to 20 tasks per unit, we average over  $N = 20$ . See Appx. C for a sensitivity study.

## 4. Results

This section is structured into two parts: First, we validate our Machine Interpretability Score (MIS) by showing that it is well correlated with existing interpretability annotations. Then, we demonstrate what type of experiments become feasible by having access to such an automated interpretability measure. Our experiments use the best-working — according to human judgements (Borowski et al., 2021) — visual explanation method, dataset examples, for computing the MIS. We demonstrate the applicability of our method to other interpretability methods (e.g., feature visualizations) in Appx. D. Note that different explanation methods might require different hyperparameters for computing the MIS. Both query images and explanations are chosen from the training set of ImageNet-2012 (Russakovsky et al., 2015). When investigating layers whose feature maps have spatial dimensions, we consider the spatial mean over a channel as one unit (e.g., Borowski et al., 2021). We ignore units with constant activations from our analysis as there is no behavior to understand (see Appx. E for details). The code for all experiments can be found at [URL included in camera-ready version].

### 4.1. Validating the Machine Interpretability Score

We validate our MIS measure by using the interpretability annotations in the IMI dataset (Zimmermann et al., 2023), which will be referred to as Human Interpretability Scores (HIS). The per-unit annotations are responses to the 2-AFC task described in Sec. 3, averaged over  $\approx 30$  participants. IMI contains scores for a subset of units for nine models.<sup>1</sup>

#### 4.1.1. MIS EXPLAINS EXISTING DATA

First, we reproduce the main result of Zimmermann et al. (2023): A comparison of nine models in terms of their the per-unit interpretability. We plot the HIS and MIS values (averaged over all units in a model) in Fig. 2 and find very strong correlations (Pearson’s  $r = 0.98$  and Spearman’s  $r = 0.94$ ). Reproducing the model ranking is strong evidence for the validity of the metric, as no information about these rankings was explicitly used to create our new measure.

<sup>1</sup>Zimmermann et al. (2023) investigate nine different models but test two of them in multiple settings, resulting in 14 distinct experimental conditions to compare.

Next, we can zoom in and look at individual units instead of per-model averages. Fig. 3 shows MIS and HIS for all units of IMI. The left figure clearly shows a strong correlation (Pearson’s and Spearman’s  $r = 0.80$ ). The interpretability scores in IMI are a (potentially noisy) estimate over a finite number of annotators. We estimate the ceiling performance due to noise (sampling 30 trials from a Bernoulli distribution) to equal a Pearson’s  $r = 0.82$  (see Appx. B for details). The right figure shows an alternative visualization, which bins the units according to their MIS and averages the HIS to reduce this noise — highlighting that the two scores correlate strongly. We can conclude that the MIS explains existing interpretability annotations well - both on a per-unit and on a per-model level.

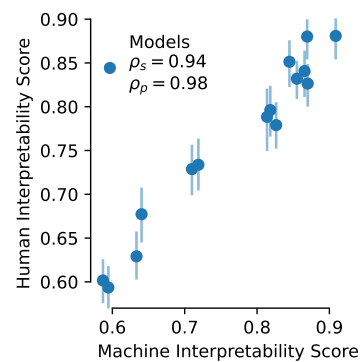


Fig. 2: **MIS Explains Interpretability Model Rankings.**

Our proposed Machine Interpretability Score (MIS) explains existing interpretability annotations (Human Interpretability Score, HIS) from IMI (Zimmermann et al., 2023) well: It reproduces the ranking of models presented in IMI while being fully automated and not requiring any human labor, as evident by the strong correlation between MIS and HIS.

#### 4.1.2. MIS MAKES NOVEL PREDICTIONS

While the previous results show a strong relation between MIS and human-perceived interpretability, they are of a descriptive (correlational) nature. To further test the match between MIS and HIS, we now turn to a causal (interventional) experiment: Instead of predicting the interpretability of units *after* a psychophysics evaluation produced their human scores, we now compute the MIS *before* conducting the psychophysics evaluation. We perform our experiment for two models: GoogLeNet and a ResNet-50. For each model, IMI contains interpretability scores for 96 randomly chosen units. We look at all the units not tested so far and find the 42 units yielding the highest (Easiest, average of 0.99 for both models) and lowest (Hardest, average of 0.63 and 0.59, respectively) MIS, respectively. Then, we use the same setup as Zimmermann et al. (2023) and perform a psychophysical evaluation on Amazon Mechanical Turk with 236 participants. We compare the HIS for the random units

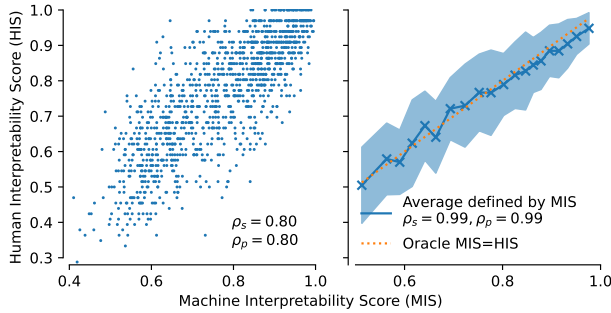


Fig. 3: **MIS Explains Per-unit Interpretability Annotations.** The proposed MIS does not only explain summary statistics for an entire model (see Fig. 2) but also individual per-unit interpretability annotations. The left side shows the calculated MIS and the recorded HIS for every unit in IMI. For the right side, the data points are grouped by their MIS into 20 bins of equal count, and the mean and standard deviation of the HIS are shown (blue). As a guideline, we display the curve (orange) an ideal metric would produce.

from the IMI dataset and the two newly recorded groups (easy, hard) of units in Fig. 4. The results are very clear again: As predicted by the MIS, the HIS is highest for the easiest and lowest for the hardest units. Further, the HIS is close to the *a priori* determined MIS given above. This demonstrates the strong predictive power of the MIS and its ability to be used for formulating novel hypotheses.

## 4.2. Analyzing & Comparing Hundreds of Models

After confirming the validity of the MIS, we now change gears and show use cases for it, i.e., experiments and analyses that were truly infeasible before due to the high cost, both time and money, of human evaluations.

### 4.2.1. COMPARISON OF MODELS

Zimmermann et al. (2023) investigated whether model or training design choices influence the interpretability of vision models. Although they invested a considerable amount of money in this investigation ( $\geq 12\,000$  USD), they could only compare nine models via a subset of units. We now scale up this line of work by two orders of magnitude and investigate all units of 835 models, almost all of which come from the well-established computer vision library timm (Wightman, 2019). See Appx. F for a list of models. Putting this scale into perspective, achieving the same scale by scaling up previous human psychophysics experiments would amount to the absurd costs of more than one billion USD. Following previous work we ignore the first and last layers of each model (Zimmermann et al., 2023).

When sorting the models according to their average MIS (Fig. 5) they span a value range of  $\approx 0.80 - 0.91$ . The strongest differences across models are present at the tails

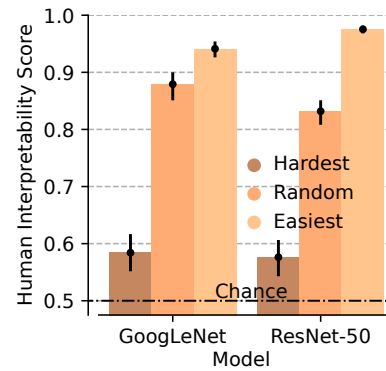


Fig. 4: **MIS Allows Detection of (Non-) Interpretable Units.** We use MIS to perform a causal intervention and determine the least (*hardest*) and most (*easiest*) interpretable units in a GoogLeNet and ResNet-50. We then use the psychophysics setup of Zimmermann et al. (2023) to measure their interpretability and compare them to randomly sampled units. Strikingly, the psychophysics results match the predicted properties: Units with the lowest MIS have significantly lower interpretability than random units, which have significantly lower interpretability than those with the highest MIS. Errorbars denote the 95 % confidence interval.

of the ranking. Note that GoogLeNet is ranked as the most interpretable model, resonating with the community’s interest in GoogLeNet as it is widely claimed to be more interpretable. The shaded area denotes the 5th to 95th percentile of the distribution across units. This reveals a strong difference in the variability of units for different models; further, as the upper end of the MIS is similar across models ( $\approx 95\%$ ), most of the change in the average score seems to stem from a change in the lower end, with decreasing width of the per-unit distribution for higher model rank.

To investigate the difference in how the MIS of units is distributed between different models, we select 15 exemplary models and visualize their per-unit MIS distribution in Fig. 6. Those models were chosen according to the distance between 5th and 95th percentile (five with highest, average, and lowest distance). While models with low and medium variability have unimodal left-skewed distributions, the ones with high variability have a rather bimodal distribution. Note that the distribution’s second, stronger mode has a similar mean and shape to the overall distribution for models with low variability. The first mode is placed at a value range slightly above 0.5, corresponding to the chance level in the task, indicating mostly uninterpretable units. This suggests that a subset of uninterpretable units (see Fig. 26 for examples) can explain most of the models’ differences in average MIS. We analyze this further in Fig. 23, where we compare the models in terms of their worst units. We see a similar shape as in Fig. 5, but with a larger value range

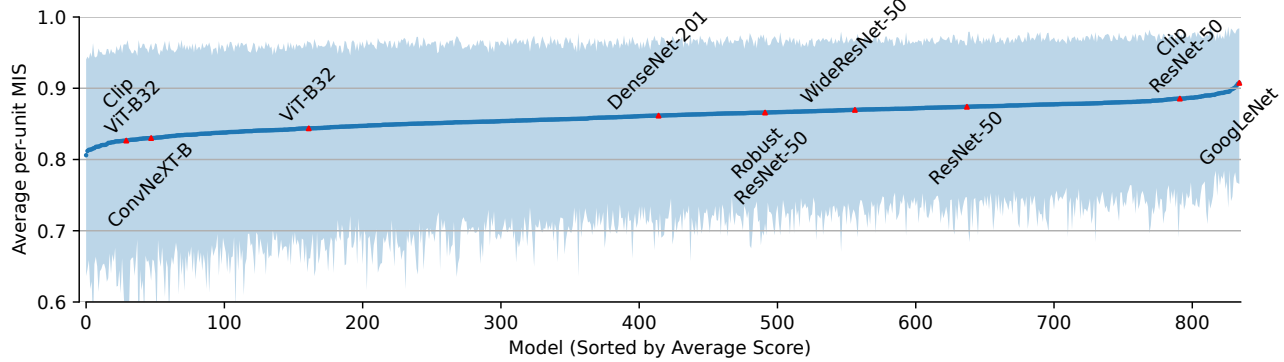


Fig. 5: **Comparison of the Average Per-unit MIS for Models.** We substantially extend the analysis of Zimmermann et al. (2023) from a noisy average over a few units for a few models to all units of 835 models. The models are compared regarding their average per-unit interpretability (as judged by MIS); the shaded area depicts the 5th to 95th percentile over units. We see that all models fall into an intermediate performance regime, with stronger changes in interpretability at the tails of the model ranking. Models probed by Zimmermann et al. (2023) are highlighted in red.

used, resulting in stronger model differences.

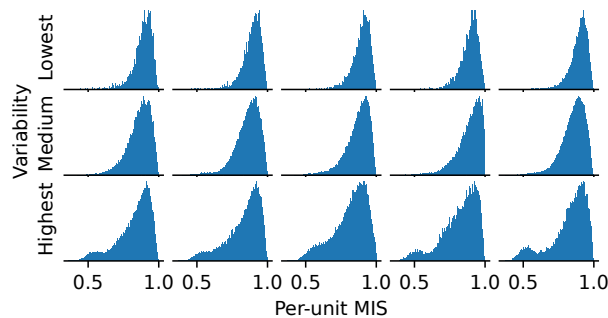


Fig. 6: **Distribution of per-unit MIS for Models.** Distribution of the per-unit MIS for 15 models, which were chosen based on the size of the error bar in Fig. 5: lowest (top row), medium (middle row), and highest variability (bottom row). While most models have a unimodal distribution, those with high variability have a second mode with lower MIS.

Previous work analyzed a potential correlation between interpretability and downstream classification performance. However, in a limited evaluation, it was found that better classifiers are not necessarily more interpretable (Zimmermann et al., 2023). A re-evaluation of this question is performed in Fig. 7 and paints an even darker picture: Here, better performing ImageNet classifiers are less interpretable (Pearson’s  $r = -0.5$  and Spearman’s  $r = -0.55$ ).

Among training procedures and architecture, the analyzed models also differ in the required resolution of their input. While previous work focused only on models with a single resolution (Zimmermann et al., 2023), we can now see whether the resolution influences interpretability. However, Fig. 21 suggests that there is no influence.

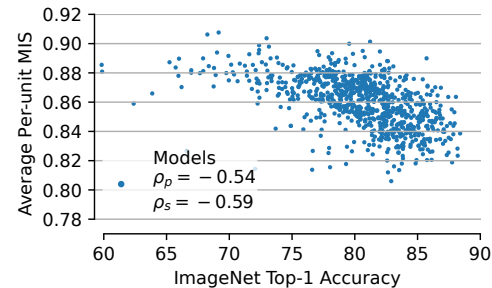


Fig. 7: **Relation Between ImageNet Accuracy and MIS.** The average per-unit MIS of a model is anticorrelated with the model’s top-1 ImageNet classification accuracy.

#### 4.2.2. COMPARISON OF LAYERS

Next, we can zoom into the results of Fig. 5 and investigate whether there are differences between different layers.

First, we are interested in testing whether the layer type is important, e.g., are convolutional more interpretable than normalization or linear layers? In Fig. 8, we sort the models by their average MIS over all layer types but show individual points for each of the five most common types (Conv, Linear, BatchNorm, LayerNorm, and GroupNorm). The number of points per model may vary, as not all models contain layers of all types. The figure shows a benefit of Conv over BatchNorm layers, which themselves are better than LayerNorm layers. Linear layers, if present, outperform both Batch- and LayerNorm as well as Conv layers. While the differences are small, they are statistically significant due to the large number of scores collected (see Fig. 20).

Second, we analyze whether the location of a layer inside a model plays a role, e.g., are earlier layers more interpretable than later ones? The average per-unit MIS (for each layer type) is shown in Fig. 9 as a function of the relative depth

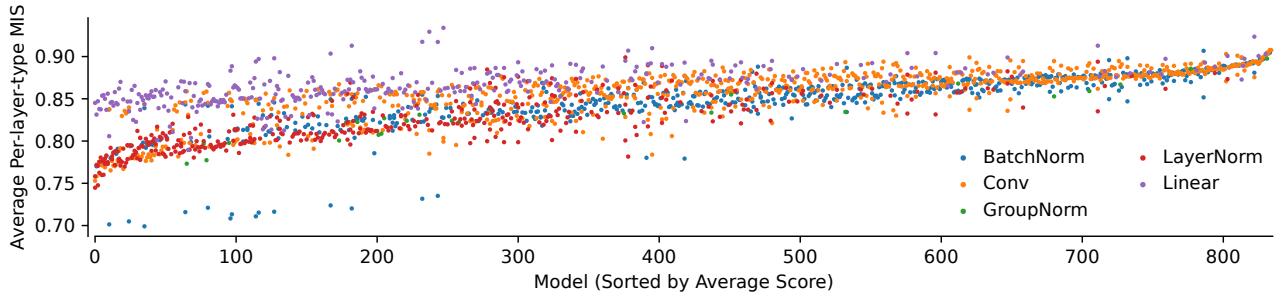


Fig. 8: **Comparison of the Average Per-unit MIS for Different Layer Types and Models.** We show the average interpretability of units from the most common layer types in vision models (BatchNorm, Conv, GroupNorm, LayerNorm, Linear). We follow Zimmermann et al. (2023) and restrict our analysis of Vision Transformers to the linear layers in each attention head. While not every layer type is used by every model, we still see some separation between types (see Fig. 20 for significance results): Linear and convolutional layers mostly outperform normalization layers. Models are sorted by average per-unit interpretability, as in Fig. 5.

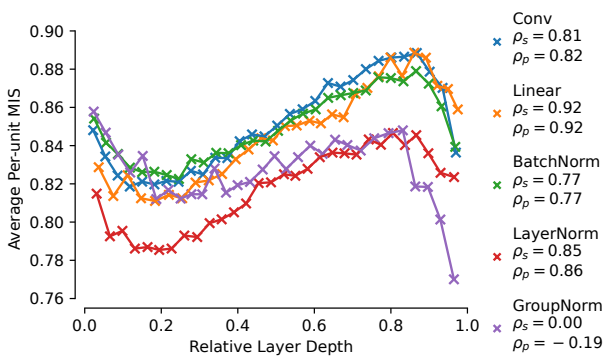


Fig. 9: **Deeper Layers are More Interpretable.** Average MIS per layer as a function of the relative depth of the layer within the network, grouped by layer types. For type, the values are grouped into 30 bins of equal count based on the relative depth; values shown correspond to the bin average.

of the layer. A value of zero corresponds to the first and a value of one to the last layer analyzed. The scores are averaged in bins of equal count defined by the relative layer depth to enhance readability. The resulting curves all follow a similar pattern: They start high, decrease in the first fifth, then increase steadily until they drop in the last tenth again, resulting in an almost sinusoidal shape.

Third, it is interesting to probe the influence of the width of layers on their average interpretability. Based on the superposition hypothesis (Elhage et al., 2022; Olah et al., 2020; Arora et al., 2018; Goh, 2016), one might expect wider layers to be more interpretable as features do not have to form in superposition (i.e., as *polysemantic* units) but can arise in a disentangled form (i.e., as *monosemantic* units). Fig. 10 shows the relation between MIS and relative layer width. We use the relative rather than the absolute width to reduce the influence of the overall model and show the results of

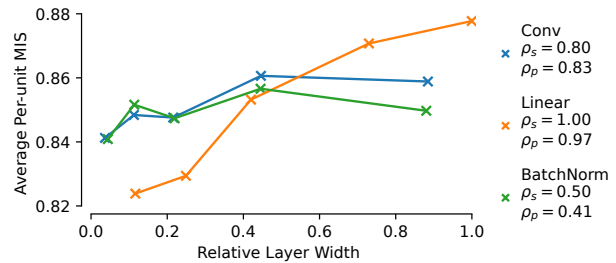


Fig. 10: **Wider Layers are More Interpretable.** Average MIS per layer as a function of the relative width of the layer compared to all layers of the same type in the network, grouped by layer types. For each type, the values are grouped into 5 bins of equal count based on the relative layer width; the values shown correspond to the bin average.

models with different architectures on the same axis. Note that, nevertheless, there might be other confounding factors correlated with the width e.g., the layer depth. While we see moderate correlations for Conv and BatchNorm layers, the one for Linear layers is much stronger. It is unclear what causes this difference in behavior. However, we see this as a hint that one way to increase a model’s interpretability is to increase the width (and not the number) of layers.

### 4.3. How Does the MIS Change During Training?

In the last set of experiments, we demonstrate how the MIS can be used to analyze models in a fine-grained way and obtain insights into their training dynamics. For this, we train a ResNet-50 on ImageNet-2012, following the training recipe A3 of Wightman et al. (2021), for 100 epochs.

Fig. 12 shows how the average per-unit MIS (left) changes during the training. Notably, the initial MIS (of the untrained network) is already substantially above chance level.

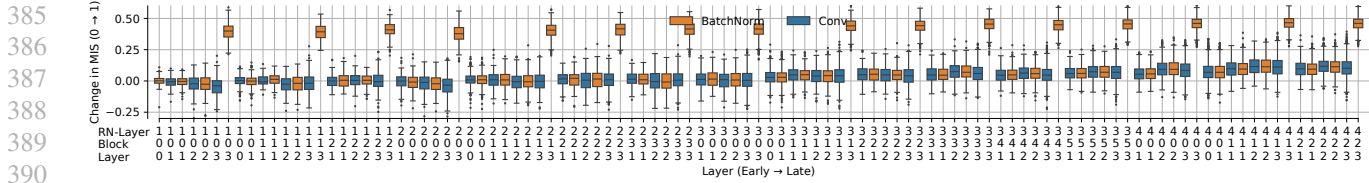


Fig. 11: **Change of Interpretability per Layer During Training.** To better understand the peak in interpretability after the first training epoch found in Fig. 12, we display the change in MIS during the first epoch, averaged over each layer. Note that layers are sorted by depth from left to right, and different colors encode different layer types. While the change in interpretability is moderately correlated with a layer’s depth, we consistently see big improvements for the last BatchNorm layer of each block (i.e., *BatchNorm*-\*-\*<sub>-3</sub>). For an extended visualization covering the full training, see Fig. 22.

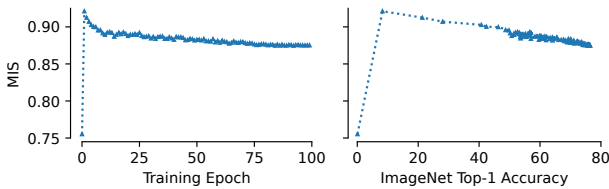


Fig. 12: **Change of Interpretability During Training.** For a ResNet-50 trained for 100 epochs on ImageNet, we track the MIS and top-1 accuracy after every epoch (epoch 0 refers to random initialization). While the MIS improves drastically in the first epoch, it monotonically decreases during the rest of the training (left). This results in an antiproportional relation between MIS and accuracy (right).

However, during the first epoch, the MIS still increases drastically to values around 0.93. Then, during the rest of the training, the score slowly decays. This indicates non-trivial dynamics of feature learning, which we analyze in Fig. 11. When showing the MIS as a function of ImageNet top-1 accuracy during training (right), a strong anticorrelation (ignoring the first points) becomes evident. This is in line with Fig. 7, also showing an anticorrelation.

To better understand the dynamics through the training — most importantly during the first epoch — we zoom in to find out which units cause this strong change in MIS. Fig. 11 shows the change in MIS during the first epoch for each layer separately (ordered by their depth within the network). Surprisingly, we see that the change in MIS is dominated by a set of BatchNorm layers, namely the last ones of each ResNetBlock, whose MIS increases drastically. Moreover, we detect a small trend of later layers improving more strongly than earlier ones but generally do not see a difference between Conv and BatchNorm layers.

## 5. Conclusion

This paper presented the first fully automated interpretability metric for vision models: the machine interpretability score

(MIS). We verified its alignment to human interpretability score (HIS) through both correlational and interventional experiments. We expect our MIS to enable experiments previously considered infeasible due to the costly reliance on human evaluations. To stress this, we demonstrated the metric’s usefulness for formulating and testing new hypotheses about a network’s behavior through a series of experiments: Based on the largest comparison of vision models in terms of their per-unit interpretability so far, we investigated potential influences on their interpretability, such as layer depth and width. Most importantly, we find an anticorrelation between a model’s downstream performance and its per-unit interpretability. Further, we performed the first detailed analysis of how the interpretability changes during training.

While this paper considerably advances the state of interpretability evaluations, there are some open questions and potential future research directions. Most importantly, the performance of our MIS on a per-unit level is close to the noise ceiling determined by the limited number of human interpretability annotations available. This means that future changes in the MIS measure (e.g., based on other image perceptual similarities) might require additional human labels to determine the significance of performance improvements. Additional human labels could also be leveraged to improve the MIS by following Fu et al. (2023) to fine-tune the image similarity directly on human judgments. In another direction, using vision language models for computing the MIS could be interesting as this might, in addition to a numerical score, also provide a textual description of a unit’s sensitivity (Hernandez et al., 2022). Finding a differentiable approximation of the MIS will be valuable for explicitly training models to be interpretable (Zimmermann et al., 2023). Note that while this paper looked at the interpretability of channels and neurons, it can also be used for analyzing arbitrary directions in activation space. Thus, we expect the MIS to also be valuable for researchers generally looking for more interpretable representations of (artificial) neural activations (e.g., Graziani et al., 2023). Finally, exploring whether this concept of interpretability quantification can be expanded to LLMs is an exciting direction.



## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning, specifically the field of Interpretable Machine Learning. The main contribution of our work is the presentation of a more time- and cost-efficient approach for quantifying how well humans can understand neural activations. A potential risk in automating interpretability research is that we will start optimizing for metrics that are never fully aligned with human judgments. It is conceivable that this will encourage the design of models that ace our metric but whose inner workings and decision making processes are still obscure to human observers. This would set false goal posts and potentially come with safety risks if a high score in MIS were mistaken for a white box model that comes with higher trustworthiness. Beyond that, we see many potential use cases for this result (see Sec. 5), that can all advance the state of machine learning. There are potential societal consequences of our work, however, none of which we feel must be specifically highlighted here.

## References

- Arora, S., Li, Y., Liang, Y., Ma, T., and Risteski, A. Linear Algebraic Structure of Word Senses, with Applications to Polysemy, December 2018. Cited on page 7.
- Barlow, H. Single units and sensation: A neuron doctrine for perceptual psychology? *Perception*, 1:371–94, 02 1972. doi: 10.1068/p010371. Cited on page 2.
- Bau, D., Zhou, B., Khosla, A., Oliva, A., and Torralba, A. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. Cited on page 2.
- Bau, D., Zhu, J.-Y., Strobel, H., Lapedriza, A., Zhou, B., and Torralba, A. Understanding the role of individual units in a deep neural network. *Proceedings of the National Academy of Sciences*, 117(48): 30071–30078, September 2020. doi: 10.1073/pnas.1907375117. URL <https://doi.org/10.1073/pnas.1907375117>. Cited on page 2.
- Bills, S., Cammarata, N., Mossing, D., Tillman, H., Gao, L., Goh, G., Sutskever, I., Leike, J., Wu, J., and Saunders, W. Language models can explain neurons in language models. <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>, 2023. Cited on pages 1 and 2.
- Borowski, J., Zimmermann, R. S., Schepers, J., Geirhos, R., Wallis, T. S. A., Bethge, M., and Brendel, W. Exemplary natural images explain cnn activations better than state-of-the-art feature visualization. In *Ninth International Conference on Learning Representations (ICLR 2021)*, 2021. Cited on pages 2, 3, 4, 12, and 13.
- Bricken, T., Templeton, A., Batson, J., Chen, B., Jermyn, A., Conerly, T., Turner, N., Anil, C., Denison, C., Askell, A., Lasenby, R., Wu, Y., Kravec, S., Schiefer, N., Maxwell, T., Joseph, N., Hatfield-Dodds, Z., Tamkin, A., Nguyen, K., McLean, B., Burke, J. E., Hume, T., Carter, S., Henighan, T., and Olah, C. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>. Cited on pages 1 and 2.
- Cammarata, N., Carter, S., Goh, G., Olah, C., Petrov, M., Schubert, L., Voss, C., Egan, B., and Lim, S. K. Thread: Circuits. *Distill*, 2020. doi: 10.23915/distill.00024. <https://distill.pub/2020/circuits>. Cited on page 1.
- Ding, K., Ma, K., Wang, S., and Simoncelli, E. P. Image Quality Assessment: Unifying Structure and Texture Similarity. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(5):2567–2581, May 2022. ISSN 1939-3539. doi: 10.1109/TPAMI.2020.3045810. Cited on page 13.
- Elhage, N., Nanda, N., Olsson, C., Henighan, T., Joseph, N., Mann, B., Askell, A., Bai, Y., Chen, A., Conerly, T., DasSarma, N., Drain, D., Ganguli, D., Hatfield-Dodds, Z., Hernandez, D., Jones, A., Kernion, J., Lovitt, L., Ndousse, K., Amodei, D., Brown, T., Clark, J., Kaplan, J., McCandlish, S., and Olah, C. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>. Cited on page 2.
- Elhage, N., Hume, T., Olsson, C., Schiefer, N., Henighan, T., Kravec, S., Hatfield-Dodds, Z., Lasenby, R., Drain, D., Chen, C., Grosse, R., McCandlish, S., Kaplan, J., Amodei, D., Wattenberg, M., and Olah, C. Toy models of superposition. *Transformer Circuits Thread*, 2022. Cited on pages 1 and 7.
- Erhan, D., Bengio, Y., Courville, A., and Vincent, P. Visualizing higher-layer features of a deep network. *Technical Report, Univeristé de Montréal*, 01 2009. Cited on page 2.
- Fu, S., Tamir, N., Sundaram, S., Chai, L., Zhang, R., Dekel, T., and Isola, P. DreamSim: Learning New Dimensions of Human Visual Similarity using Synthetic Data, December 2023. Cited on pages 3, 8, and 13.
- Gilpin, L. H., Bau, D., Yuan, B. Z., Bajwa, A., Specter, M. A., and Kagal, L. Explaining explanations: An overview of interpretability of machine learning. *2018 IEEE 5th International Conference on Data Science and*

- Advanced Analytics (DSAA), pp. 80–89, 2018. Cited on page 2.
- Goh, G. Decoding the Thought Vector. <https://gabgoh.github.io/ThoughtVectors/>, 2016. Cited on page 7.
- Graziani, M., O’Mahony, L., Nguyen, A.-p., Müller, H., and Andrearczyk, V. Uncovering unique concept vectors through latent space decomposition. *Transactions on Machine Learning Research*, 2023. Cited on page 8.
- Hernandez, E., Schwettmann, S., Bau, D., Bagashvili, T., Torralba, A., and Andreas, J. Natural Language Descriptions of Deep Visual Features, April 2022. Cited on pages 2, 3, and 8.
- Huang, J., Geiger, A., D’Oosterlinck, K., Wu, Z., and Potts, C. Rigorously assessing natural language explanations of neurons. *arXiv preprint arXiv:2309.10312*, 2023. Cited on page 3.
- Hubel, D. H. and Wiesel, T. N. Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex. *J. Physiol.*, 160(1):106–154, January 1962. Cited on page 2.
- Kalibhat, N., Bhardwaj, S., Bruss, C. B., Firooz, H., Sanjabi, M., and Feizi, S. Identifying interpretable subspaces in image representations. In *International Conference on Machine Learning*, pp. 15623–15638. PMLR, 2023. Cited on page 2.
- Kim, S. S. Y., Meister, N., Ramaswamy, V. V., Fong, R., and Russakovsky, O. HIVE: Evaluating the human interpretability of visual explanations. In *European Conference on Computer Vision (ECCV)*, 2022. Cited on page 2.
- Klindt, D., Sanborn, S., Acosta, F., Poitevin, F., and Miolane, N. Identifying interpretable visual features in artificial and biological neural systems. *arXiv preprint arXiv:2310.11431*, 2023. Cited on page 1.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. ImageNet Classification with Deep Convolutional Neural Networks. In Bartlett, P. L., Pereira, F. C. N., Burges, C. J. C., Bottou, L., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a Meeting Held December 3-6, 2012, Lake Tahoe, Nevada, United States*, pp. 1106–1114, 2012. Cited on pages 1 and 13.
- Leavitt, M. L. and Morcos, A. S. Towards falsifiable interpretability research. *CoRR*, abs/2010.12016, 2020. URL <https://arxiv.org/abs/2010.12016>. Cited on page 2.
- Mahendran, A. and Vedaldi, A. Understanding deep image representations by inverting them. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5188–5196, Los Alamitos, CA, USA, jun 2015. IEEE Computer Society. doi: 10.1109/CVPR.2015.7299155. URL <https://doi.ieeecomputersociety.org/10.1109/CVPR.2015.7299155>. Cited on pages 1 and 2.
- Morcos, A. S., Barrett, D. G., Rabinowitz, N. C., and Botvinick, M. On the importance of single directions for generalization. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rliuQjxCZ>. Cited on page 2.
- Mordvintsev, A., Olah, C., and Tyka, M. Inceptionism: Going deeper into neural networks, 2015. URL <https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>. Cited on page 2.
- Nanda, N., Chan, L., Lieberum, T., Smith, J., and Steinhart, J. Progress measures for grokking via mechanistic interpretability, 2023. Cited on page 1.
- Nguyen, A., Yosinski, J., and Clune, J. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 12 2014. Cited on page 2.
- Nguyen, A., Clune, J., Bengio, Y., Dosovitskiy, A., and Yosinski, J. Plug & play generative networks: Conditional iterative generation of images in latent space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2017. Cited on page 2.
- Oikarinen, T. and Weng, T.-W. Clip-dissect: Automatic description of neuron representations in deep vision networks. In *The Eleventh International Conference on Learning Representations*, 2022. Cited on page 3.
- Olah, C. Mechanistic interpretability, variables, and the importance of interpretable bases, 2022. URL <https://transformer-circuits.pub/2022/mech-interp-essay/index.html>. Cited on pages 1 and 2.
- Olah, C., Mordvintsev, A., and Schubert, L. Feature visualization. *Distill*, 2017. doi: 10.23915/distill.00007. <https://distill.pub/2017/feature-visualization>. Cited on pages 1 and 2.
- Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., and Carter, S. Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001. <https://distill.pub/2020/circuits/zoom-in>. Cited on page 7.

- 550 Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma,  
551 N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen,  
552 A., Conerly, T., Drain, D., Ganguli, D., Hatfield-Dodds,  
553 Z., Hernandez, D., Johnston, S., Jones, A., Kernion, J.,  
554 Lovitt, L., Ndousse, K., Amodei, D., Brown, T., Clark, J.,  
555 Kaplan, J., McCandlish, S., and Olah, C. In-context  
556 learning and induction heads. *Transformer Circuits*  
557 *Thread*, 2022. [https://transformer-circuits.pub/2022/in-](https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html)  
558 [context-learning-and-induction-heads/index.html](https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html). Cited  
559 on page 2.
- 560 Quiroga, R. Q., Reddy, L., Kreiman, G., Koch, C., and Fried,  
561 I. Invariant visual representation by single neurons in the  
562 human brain. *Nature*, 435(7045):1102–1107, 2005. Cited  
563 on page 2.
- 564 Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S.,  
565 Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein,  
566 M., Berg, A. C., and Fei-Fei, L. ImageNet Large Scale  
567 Visual Recognition Challenge. *International Journal of*  
568 *Computer Vision (IJCV)*, 115(3):211–252, 2015. doi:  
569 10.1007/s11263-015-0816-y. Cited on page 4.
- 570 Schwettmann, S., Shaham, T. R., Materzynska, J., Chowd-  
571 hury, N., Li, S., Andreas, J., Bau, D., and Torralba, A.  
572 FIND: A Function Description Benchmark for Evaluat-  
573 ing Interpretability Methods, December 2023. Cited on  
574 page 3.
- 575 Simonyan, K. and Zisserman, A. Very Deep Convolutional  
576 Networks for Large-Scale Image Recognition. In Bengio,  
577 Y. and LeCun, Y. (eds.), *3rd International Conference*  
578 *on Learning Representations, ICLR 2015, San Diego,*  
579 *CA, USA, May 7-9, 2015, Conference Track Proceedings,*  
580 *2015*. Cited on page 13.
- 581 Wightman, R. Pytorch image models. [https://github.](https://github.com/rwightman/pytorch-image-models)  
582 [com/rwightman/pytorch-image-models,](https://github.com/rwightman/pytorch-image-models)  
583 2019. Cited on pages 5 and 16.
- 584 Wightman, R., Touvron, H., and Jégou, H. Resnet strikes  
585 back: An improved training procedure in timm. *arXiv*  
586 *preprint arXiv:2110.00476*, 2021. Cited on page 7.
- 587 Yosinski, J., Clune, J., Nguyen, A., Fuchs, T., and Lipson,  
588 H. Understanding neural networks through deep visu-  
589 alization. In *Deep Learning Workshop, International*  
590 *Conference on Machine Learning (ICML)*, 2015. Cited  
591 on page 2.
- 592 Zhang, R., Isola, P., Efros, A. A., Shechtman, E., and Wang,  
593 O. The Unreasonable Effectiveness of Deep Features  
594 as a Perceptual Metric. In *2018 IEEE Conference on*  
595 *Computer Vision and Pattern Recognition, CVPR 2018,*  
596 *Salt Lake City, UT, USA, June 18-22, 2018*, pp. 586–595.  
597 Computer Vision Foundation / IEEE Computer Society,  
598 2018. doi: 10.1109/CVPR.2018.00068. Cited on pages 3  
599 and 13.
- 600 Zhou, B., Sun, Y., Bau, D., and Torralba, A. Revisiting  
601 the importance of individual units in cnns via ablation.  
602 *CoRR*, abs/1806.02891, 2018. URL [http://arxiv.](http://arxiv.org/abs/1806.02891)  
603 [org/abs/1806.02891](http://arxiv.org/abs/1806.02891). Cited on page 2.
- 604 Zimmermann, R. S., Borowski, J., Geirhos, R., Bethge,  
M., Wallis, T., and Brendel, W. How well do fea-  
ture visualizations support causal understanding of  
cnn activations? In Ranzato, M., Beygelzimer, A.,  
Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.),  
*Advances in Neural Information Processing Systems*,  
volume 34, pp. 11730–11744. Curran Associates, Inc.,  
2021. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2021/file/618faa1728eb2ef6e3733645273ab145-Paper.pdf)  
[cc/paper\\_files/paper/2021/file/](https://proceedings.neurips.cc/paper_files/paper/2021/file/618faa1728eb2ef6e3733645273ab145-Paper.pdf)  
[618faa1728eb2ef6e3733645273ab145-Paper.](https://proceedings.neurips.cc/paper_files/paper/2021/file/618faa1728eb2ef6e3733645273ab145-Paper.pdf)  
pdf. Cited on pages 2, 4, and 12.
- Zimmermann, R. S., Klein, T., and Brendel, W. Scale  
alone does not improve mechanistic inter-  
pretability in vision models. In *Thirty-seventh Con-*  
*ference on Neural Information Processing Systems,*  
2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=OZ7aImD4uQ)  
[id=OZ7aImD4uQ](https://openreview.net/forum?id=OZ7aImD4uQ). Cited on pages 1, 2, 3, 4, 5, 6, 7, 8,  
12, and 16.

## A. Description of the 2-AFC Task

### A.1. Task Design

Our proposed MIS builds on the 2-AFC task designed by [Borowski et al. \(2021\)](#) to conduct human psychophysics experiments. An example of such a task is given in Fig. 13.

This task aims to probe how well (human) participants can detect the sensitivity of a unit of a neural network based on visual explanations of it. Understanding the unit’s sensitivity should allow participants to distinguish between a stimulus eliciting highly activating from one yielding low activation. Therefore, the task shows the participants two such images, called query images, and asks them to pick the image eliciting higher activation. To solve the task, participants also see two sets of visual explanations: Positive explanations describe the patterns the unit activates strongly for, while negative activations show patterns the unit weakly responds to. For solving this task, there are two potential strategies: Participants can either recognize a common pattern of the positive explanations in one of the query images, making this the correct choice. Or they detect a common pattern of the negative explanations in a query image, making the other one the right choice. See [Borowski et al. \(2021\)](#); [Zimmermann et al. \(2021\)](#) or [Zimmermann et al. \(2023\)](#) for alternative descriptions and visualizations of the task.



Fig. 13: **Examples of the 2-AFC Task.** For two different units of GoogLeNet one task each is shown. Every task contains a set of negative (left) and positive (right) visual explanations describing which visual feature the unit is sensitive to. In the center, two query images in the form of strongly and weakly activating dataset examples are shown, respectively. This means that each one of the two query images corresponds to the positive and the other to the negative explanations. The task is now to choose which query image corresponds to the positive ones.

### A.2. Task Construction

For constructing tasks, we follow [Zimmermann et al. \(2023\)](#). Specifically, this means that we use  $K = 9$  (positive and negative) explanations in each task. We restrict explanations to natural dataset examples to reduce complexity but note that the same setup can also be applied to other visual explanations, such as feature visualizations. To choose query images and explanations, we proceed as follows: For each unit, we determine the  $N \cdot (K + 1)$  most and least activating images, respectively. Out of these, the  $N \cdot K$  most extreme images are used as explanations, the others as query images. The  $N \cdot K$  potential explanation images are uniformly distributed across tasks according to their elicited activation level (see [Borowski et al., 2021](#); [Zimmermann et al., 2023](#) for more details).

## B. Influence of the Underlying Perceptual Similarity on the Machine Interpretability Score

As stated in Sec. 3, we used DreamSim (Fu et al., 2023) as the underlying perceptual similarity  $f$  for all experiments shown so far. We now repeat the experiments on IMI in Sec. 4.1.1 with two alternative similarity measures: LPIPS (Zhang et al., 2018) and DISTS (Ding et al., 2022). While all three measures are based on learned image features, DreamSim leverages an ensemble of modern vision models trained on larger datasets compared to LPIPS and DISTS, which use AlexNet (Krizhevsky et al., 2012) and VGG16 (Simonyan & Zisserman, 2015) trained on ImageNet, respectively. According to Fu et al. (2023), DreamSim clearly outperforms LPIPS and DISTS on image similarity benchmarks.

When comparing MIS based on DreamSim with one based on LPIPS and DISTS on a per-model level (see Fig. 14) one sees very similar results and strong correlations between each MIS and HIS. This might suggest that the choice of the similarity function to use has little influence on the quality of MIS. The picture, however, changes when zooming in and looking at per-unit interpretability (see Fig. 16). Now, it becomes evident that the MIS based on DreamSim outperforms that based on LPIPS and DISTS, indicated by the higher correlation and smaller spread of the point cloud. We, therefore, conclude that DreamSim is the best perceptual similarity available for computing machine interpretability scores.

To put the difference in performance between the perceptual similarities on a per-unit level into context, we estimate the noise ceiling of the data: As the HIS for a single unit is a (potentially) noisy estimate over (up to 30) human decisions, it has some uncertainty. To take this into account, we run a statistical simulation, in which we model individual human responses as binary decisions from a Bernoulli distribution whose mean equals the unit’s HIS. We can now simulate human decisions by sampling from the distribution. Then, we compute the correlation between MIS and simulated HIS and repeat the process 1 000 times. The resulting *noise ceiling* is compared to the correlations obtained when using LPIPS, DISTS, and DreamSim in Fig. 15. We see that DreamSim’s performance is very close to the noise ceiling for estimating the per-unit human interpretability.

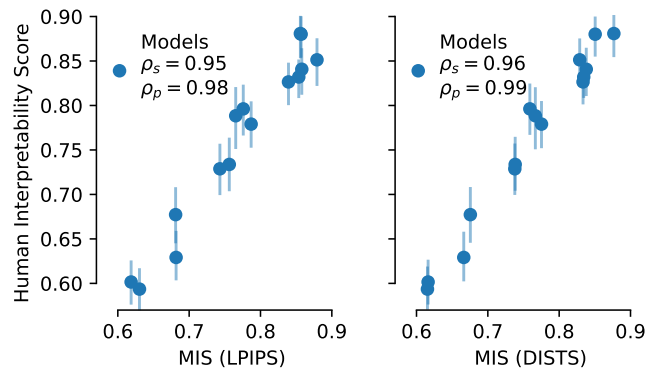


Fig. 14: **LPIPS and DISTS Perform Similarly as DreamSim when Comparing Models.** We compare DreamSim with two earlier perceptual similarity metrics, LPIPS and DISTS. All three lead to similar results on IMI (cf. Fig. 2). See Fig. 16 for comparing these similarity functions on a per-unit level. standard deviation.

## C. Sensitivity of the MIS on the Number of Tasks

As described in Sec. 3, we compute the MIS by averaging over  $N = 20$  tasks. This choice was initially motivated by previous work by Borowski et al. (2021). We investigate now how this choice influences the MIS. For this, we perform two experiments for GoogLeNet (see Fig. 17). First, we use the method for constructing tasks described before in Appx. A.2 to create 20 tasks per unit and then compute how the MIS changes when only using the first  $i = 1, \dots, 19$  tasks compared to all 20. While this setting is straightforward to analyze, it does not reflect how the number of tasks influences the MIS computation in practice: Using the task creation above, the chosen number of tasks influences the creation of all tasks, e.g., adding one more task changes which images are used for previous tasks. Therefore, in the second experiment, we again measure how the MIS changes when using  $i = 1, \dots, 19$  tasks compared to 20, but recreate all tasks when increasing their number. For both settings, we see that the residual converges to zero, with a slower convergence in the more realistic setting.

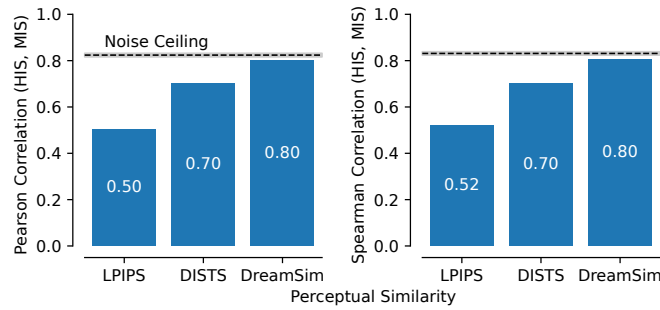


Fig. 15: **Best Perceptual Similarity Approaches Noise Ceiling.** Considering the noise ceiling, caused by the inherent uncertainty of the HIS, the best perceptual similarity (DreamSim) shows an almost perfect performance. The black bar and shaded area show the mean correlation and standard deviation over 1 000 simulations, respectively.

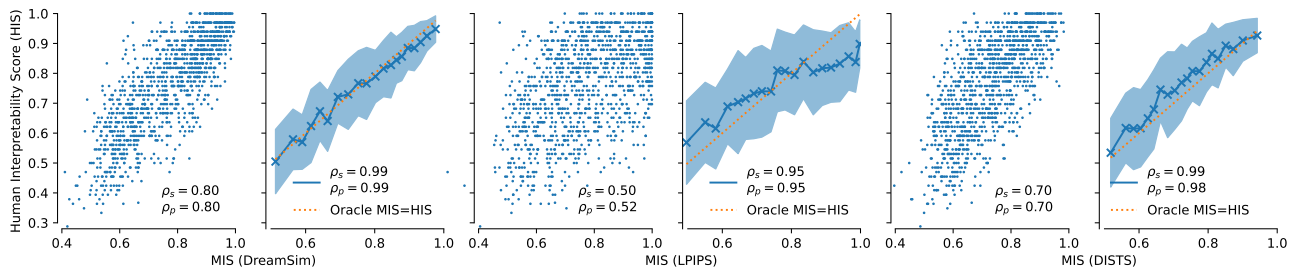


Fig. 16: **LPIPS and DISTS Perform Worse than DreamSim when Comparing Individual Units.** We compare DreamSim with two earlier perceptual similarity metrics, LPIPS and DISTS. While LPIPS and DISTS perform similarly to DreamSim on a per-model level of IMI (cf. Fig. 16), they lead to worse performance on a per-unit level.

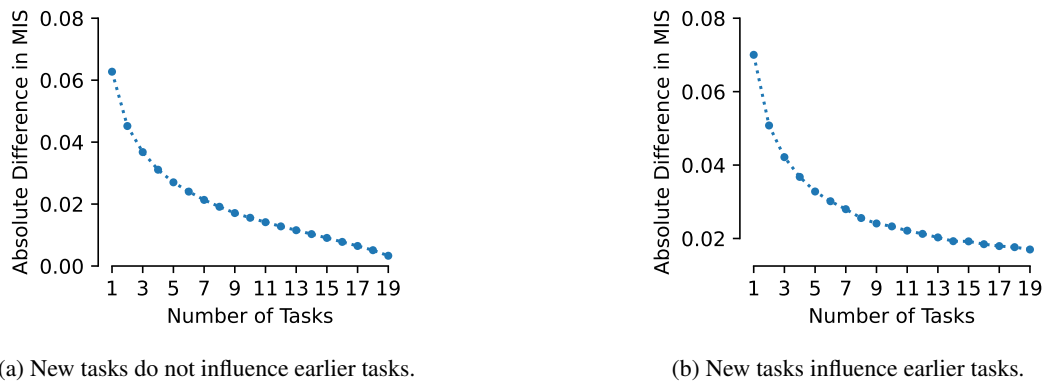


Fig. 17: **Convergence of MIS.** We investigate how MIS changes depending on the number of tasks  $N$  that it is computed over. Here, we distinguish between two settings. In (a), we simulate that adding another task does not change the selection of query images and explanations in earlier tasks; in (b), this is not the case. While the former is easier to analyze due to a reduced level of randomness, note that the latter is the more relevant setting in practice. For both cases, we visualize the average absolute difference in MIS estimated for  $< 20$  and  $N = 20$  tasks.

## D. Applying MIS for Different Explanation Methods

The experiments in Sec. 4 compute the MIS for one type of explanation, namely strongly activating dataset examples. We now demonstrate that the same approach easily generalizes to other visual explanations: feature visualizations. We do not tune any hyperparameters but re-use the same as presented in Sec. 3 for dataset examples as explanations. In Fig. 18 we repeat the experiment from Fig. 2 and again see a strong correlation between MIS and HIS.

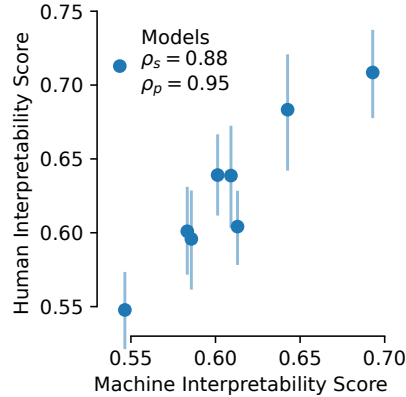


Fig. 18: **MIS Generalizes Well to Other Explanation Types.** We find a high correlation between MIS and HIS for other explanation types (feature visualizations). See Fig. 2 for the corresponding results for using natural dataset examples as explanations.

## E. Analysis of Constant Units

After training a network, it might happen that some of its units effectively become non-active/constant for any relevant image. We here call a unit *constant* if the difference between maximally and minimally elicited activation by the entire ImageNet-2012 training set is less than  $10^{-8}$ . As mentioned at the beginning of Sec. 4, we excluded those units in our analysis, as they do not present any interesting behavior that is worth understanding. Note that this does not mean that it will not be interesting to understand why such units exist. In Fig. 19, we display the ratio of constant units for each model. For most models, we see a low number of constant units: Specifically, we see that out of the 835 models investigated, 256 do not contain any constant units, 89 contain more than 1% and 22 more than 5%. Note that we here used the same notion of units as in the rest of the paper, meaning that we take the spatial mean of feature maps with spatial dimensions (e.g., for convolutional layers).

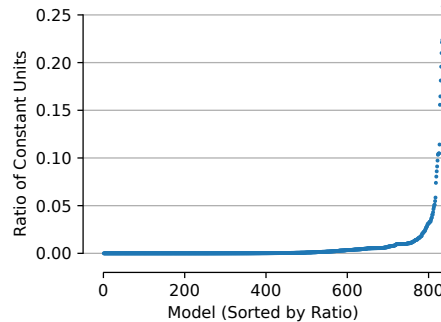


Fig. 19: **Ratio of Constant Units.** We compute the ratio of units constant with respect to the input (over the training set of ImageNet-2012) for all models considered. While the ratio is low for most models, it becomes large for a few models.

## F. Details on Models

In addition to the 9 models investigated by Zimmermann et al. (2023) (GoogLeNet, ResNet-50, Clip ResNet-50, Robust (L2) ResNet-50, DenseNet-101, WideResNet-50, Clip ViT-B32, ViT-B32), we include one more model suggested by them (Robust (L2) ResNet-50) and 825 models from timm (Wightman, 2019):

xcit\_tiny\_12\_p16\_224\_fb.in1k, vit\_tiny\_patch16\_384\_augreg\_in21k\_ft.in1k, pit\_xs\_224.in1k, repghostnet\_111.in1k, regnetz\_c16\_evos.ch.in1k, poolformer\_m48\_sail.in1k, repghostnet\_080.in1k, volo\_d3\_448.sail.in1k, vit\_base\_patch16\_224\_augreg\_in21k\_ft.in1k, regnety\_320.tv2.in1k, densenet121.ra.in1k, mobilenetv3\_large\_100.ra.in1k, repghostnet\_150.in1k, seresnext26ts.ch.in1k, regnety\_160.swag\_ft.in1k, hrnet\_w40.ms.in1k, convnext\_small.in12k\_ft.in1k, vit\_base\_patch16\_224\_sam.in1k, seresnextaa101d\_32x8d.sw.in12k\_ft.in1k\_288, vit\_tiny\_r\_s16\_p8\_384\_augreg\_in21k\_ft.in1k, regnety\_320.pycls.in1k, cs3darknet\_m.c2ns.in1k, vit\_tiny\_patch16\_224\_augreg\_in21k\_ft.in1k, resnet101c.gluon.in1k, convnextv2\_atto.fcmae\_ft.in1k, flexivit\_base.600ep.in1k, xcit\_small\_12\_p16\_384\_fb\_dist.in1k, mobilenetv2.050.lamb.in1k, flexivit\_base.300ep.in1k, resnext50\_32x4d.tv.in1k, resnet152.tv.in1k, seresnext26d\_32x4d.bt.in1k, fbnetv3\_g.ra2.in1k, poolformer\_s36.sail.in1k, resnext101\_32x8d.tv.in1k, rexnet\_130.nav.in1k, efficientvit\_b2\_r224.in1k, convnext\_small.fb.in22k\_ft.in1k\_384, resnet50\_gn.a1h.in1k, eva02\_small\_patch14\_336.mim.in22k\_ft.in1k, regnety\_032.ra.in1k, res2net50d.in1k, convit\_small.fb.in1k, regnetz\_160.pycls.in1k, convnextv2\_large.fcmae\_ft.in22k.in1k\_384, tf\_efficientnet\_b0.ns.jft.in1k, pit\_ti\_224.in1k, volo\_d1\_384.sail.in1k, xcit\_small\_12\_p8\_384\_fb\_dist.in1k, dpn131.mx.in1k, resnext101\_64x4d.gluon.in1k, densenet169.tv.in1k, resnet101d.ra2.in1k, repghostnet\_200.in1k, resnet18.a2.in1k, xcit\_small\_12\_p16\_224\_fb.in1k, pvt\_v2\_b3.in1k, dm\_nfnct\_fl.dm.in1k, vit\_large\_patch32\_384\_orig.in21k\_ft.in1k, convnextv2\_tiny.fcmae\_ft.in22k.in1k\_384, gresnet50.ra2.in1k, nf\_regnet\_b1.ra2.in1k, volo\_d1\_224.sail.in1k, resnet50\_ram.in1k, hrnet\_w18\_small.v2.ms.in1k, convnext\_base.clip\_laion2b\_augreg\_ft.in1k, regnetx\_160.tv2.in1k, sequencer2d.l.in1k, convnext\_large.fb.in22k\_ft.in1k, botnet26t\_256.c1.in1k, gc\_efficientnetv2\_rw.t.agc.in1k, wide\_resnet50\_2\_racm.in1k, halonet50ts.a1h.in1k, csresnext50.ra.in1k, resnetv2\_50d\_evos.a1h.in1k, tf\_efficientnetv2\_b3.in21k\_ft.in1k, resnet152.gluon.in1k, lambda\_resnet26prt\_256.c1.in1k, fastvit\_sa24.apple\_dist.in1k, xcit\_medium\_24\_p8\_384\_fb\_dist.in1k, revpit\_m0.9.dist.450e.in1k, regnetz\_320.pycls.in1k, seresnextaa101d\_32x8d.sw.in12k\_ft.in1k, efficientvit\_b2\_r288.in1k, convnext\_tiny.in12k\_ft.in1k, xcit\_large\_24\_p16\_384\_fb\_dist.in1k, resnetv2\_50.a1h.in1k, coatnet\_0\_rw\_224.sw.in1k, efficientnet.es.pruned.in1k, dla60\_res2net.in1k, efficientformer\_l7.snap\_dist.in1k, cait\_xs24\_224\_fb\_dist.in1k, vit\_small\_patch16\_224\_augreg\_in21k\_ft.in1k, tf\_efficientnet\_cc\_b1\_8e.in1k, efficientvit\_b1\_r288.in1k, halonet26t.a1h.in1k, mixnet\_m.ft.in1k, hrnet\_w44.ms.in1k, regnety\_160.tv2.in1k, xcit\_nano\_12\_p8\_384\_fb\_dist.in1k, seresnext101\_32x8d.a1h.in1k, efficientvit\_b2\_r256.in1k, vit\_base\_patch16\_clip\_224\_laion2b\_ft.in12k.in1k, tf\_efficientnet\_lite2.in1k, deit3\_small\_patch16\_224\_fb.in1k, hrnet\_w18\_ssl.paddle.in1k, tf\_efficientnet\_b2\_aa.in1k, crossvit\_15\_dagger\_240.in1k, deit3\_small\_patch16\_224\_fb.in22k\_ft.in1k, haloregnetz\_b.ra3.in1k, tf\_efficientnetv2\_b0.in1k, eca\_nfnct\_10.ra2.in1k, twins\_pcpvt\_small.in1k, ecaresnet50.ra2.in1k, fastvit\_sa12.apple\_dist.in1k, skresnext50\_32x4d.ra.in1k, resnet50d.a2.in1k, vit\_base\_patch32\_clip\_224\_laion2b\_ft.in1k, resnetblur50.bt.in1k, vit\_base\_patch16\_224\_orig.in21k\_ft.in1k, resnet50.a1h.in1k, hardcorenas\_e.mil\_green.in1k, coatnext\_nano\_rw\_224.sw.in1k, convnext\_base.clip\_laion2b\_augreg\_ft.in1k\_384, tresnet\_m.mil.in1k\_448, resnet101.c3.in1k, poolformer\_m2\_m48.sail.in1k, tf\_efficientnet\_b1\_aa.in1k, edgenext\_base.usi.in1k, tf\_efficientnet.es.in1k, tresnet\_l.mil.in1k\_448, resnet152.a1h.in1k, mixnet\_s.ft.in1k, resnet50.am.in1k, rexnet\_100.nav.in1k, xcit\_large\_24\_p8\_224\_fb\_dist.in1k, deit3\_base\_patch16\_224\_fb.in22k\_ft.in1k, xcit\_tiny\_24\_p8\_384\_fb\_dist.in1k, coat\_lite\_medium\_384.in1k, focalnet\_small\_srf.ms.in1k, vit\_base\_patch8\_224\_augreg\_in21k\_ft.in1k, convnext\_tiny\_hnf.a2h.in1k, visformer\_small.in1k, vit\_small\_r26\_s32\_384\_augreg\_in21k\_ft.in1k, vgg16.bn.tv.in1k, eca\_nfnct\_l1.ra2.in1k, xcit\_small\_12\_p8\_224\_fb.in1k, beitv2\_base\_patch16\_224.in1k\_ft.in22k.in1k, cs3edgenet\_x.c2.in1k, vit\_base\_patch16\_clip\_384\_laion2b\_ft.in12k.in1k, xcit\_small\_12\_p16\_224\_fb\_dist.in1k, convformer\_b36.sail.in1k\_384, bat\_resnet26ts.ch.in1k, caformer\_b36.sail.in1k, dla34.in1k, crossvit\_18\_dagger\_240.in1k, tf\_efficientnet\_b2.in21k\_ft.in1k, focalnet\_base\_srf.ms.in1k, convformer\_b36.sail.in22k\_ft.in1k\_384, resmlp\_24\_224\_fb\_distilled.in1k, convnext\_base.clip\_laion2b\_augreg\_ft.in12k.in1k, caformer\_s18.sail.in1k\_384, resnetaa50.a1h.in1k, beitv2\_base\_patch16\_224.in1k\_ft.in1k, convformer\_m36.sail.in22k\_ft.in1k, inception\_resnet\_v2\_tf\_ens\_adv.in1k, mobilenetv2\_l110d.ra.in1k, resnext101\_32x4d.fb\_sws\_lig1b\_ft.in1k, regnetx\_008.tv2.in1k, convnext\_small.in12k\_ft.in1k\_384, levit\_conv\_128.fb\_dist.in1k, volo\_d3\_224.sail.in1k, nest\_tiny\_jx\_goog.in1k, mobileone\_s2.apple.in1k, fastvit\_t8.apple\_dist.in1k, halo2botnet50ts\_256.a1h.in1k, mobilenetv2\_140.ra.in1k, caformer\_m36.sail.in1k, seresnet50.ra2.in1k, hardcorenas\_d.mil\_green.in1k, convformer\_b36.sail.in1k, regnety\_320.swag\_ft.in1k, volo\_d4\_448.sail.in1k, tf\_efficientnet\_b2\_ns.jft.in1k, sebotnet33ts\_256.a1h.in1k, vit\_small\_patch32\_224\_augreg\_in21k\_ft.in1k, vit\_base\_patch32\_224\_sam.in1k, resnetv2\_50d\_gn.a1h.in1k, mobileone\_s4.apple.in1k, coat\_small.in1k, tf\_mixnet\_l.in1k, resnet34.a2.in1k, regnetz\_032.pycls.in1k, resnetaa101d.sw.in12k\_ft.in1k, lenet\_100.ra2.in1k, repvgg\_b1\_rvgg.in1k, crossvit\_15\_240.in1k, edgenext\_x.small.in1k, revpit\_m1\_5.dist.300e.in1k, hardcorenas\_a.mil\_green.in1k, efficientformer\_l1.snap\_dist.in1k, tf\_mobilenetv3\_large\_075.in1k, hrnet\_w18\_small.ms.in1k, tf\_efficientnet\_b2.in1k, ghostnetv2\_130.in1k, ecaresnet26t.ra2.in1k, fastvit\_s12.apple.in1k, xcit\_tiny\_12\_p8\_224\_fb\_dist.in1k, tresnet\_m.mil.in21k\_ft.in1k, fastvit\_sa24.apple.in1k, resnetrs200.tf.in1k, convnextv2\_nano.fcmae\_ft.in1k, resnet50.ra.in1k, resnet34.bt.in1k, regnety\_002.pycls.in1k, focalnet\_base\_lrf.ms.in1k, dla102.in1k, regnetz\_e8.ra3.in1k, pvt\_v2\_b0.in1k, xcit\_medium\_24\_p8\_224\_fb.in1k, regnety\_640.seer.ft.in1k, resnet200d.ra2.in1k, caformer\_s36.sail.in1k\_384, deit3\_small\_patch16\_384\_fb.in22k\_ft.in1k, eca\_resnext26ts.ch.in1k, vgg13.tv.in1k, tf\_efficientnet\_lite0.in1k, resnet50\_b1k.in1k, dla60\_res2next.in1k, revpit\_m1\_l1.dist.300e.in1k, convnext\_base.fb.in22k\_ft.in1k, tf\_efficientnet\_cc\_b0\_4e.in1k, ese\_vovnet19b.dw.ra.in1k, resnetv2\_152x2\_bit\_goog\_teacher.in12k\_ft.in1k, deit\_base\_distilled\_patch16\_384\_fb.in1k, resnet101d.gluon.in1k, convnext\_large.fb.in22k\_ft.in1k\_384, darknet53.c2ns.in1k, poolformerv2\_s36.sail.in1k, convformer\_m36.sail.in22k\_ft.in1k\_384, gmp\_s16\_224.ra3.in1k, convformer\_s18.sail.in1k, efficientnet.em.ra2.in1k, inception\_v3.gluon.in1k, resmlp\_12\_224\_fb.in1k, tresnet\_l.mil.in1k, ecaresnet101d.pruned.mil.in1k, resnet152.a2.in1k, vit\_small\_patch32\_384\_augreg\_in21k\_ft.in1k, inception\_v3\_tf\_adv.in1k, repghostnet\_130.in1k, levit\_conv\_384\_fb\_dist.in1k, revpit\_m1\_5.dist.450e.in1k, efficientnet.el.ra.in1k, seresnet50.a2.in1k, pit\_s\_distilled\_224.in1k, cs3darknet53.ra.in1k, tf\_efficientnet\_cc\_b0\_8e.in1k, densenet201.tv.in1k, resnext50\_32x4d.a1.in1k, cs3darknet\_l.c2ns.in1k, cait\_s24\_384\_fb\_dist.in1k, spnasnet\_100.rmsp.in1k, res2net50\_14w\_8s.in1k, repvgg\_d2se.rvgg.in1k, regnetz\_032.tv2.in1k, crossvit\_18\_dagger\_408.in1k, pit\_b\_distilled\_224.in1k, cs3darknet\_focus\_l.c2ns.in1k, resnet50.bt.in1k, vgg11.tv.in1k, convnextv2\_femto.fcmae\_ft.in1k, convnext\_nano.in12k\_ft.in1k, resnext101\_64x4d.tv.in1k, convnext\_nano.d1h.in1k, csresnet50.ra.in1k, tf\_mixnet\_m.in1k, xcit\_tiny\_12\_p16\_384\_fb\_dist.in1k, seresnet50.a1.in1k, efficientnetv2\_rw.ra2.in1k, resnet152d.gluon.in1k, regnety\_032.tv2.in1k, inception\_resnet\_v2\_tf.in1k, eva\_large\_patch14\_196.in22k\_ft.in1k, pvt\_v2\_b1.in1k, convformer\_m36.sail.in1k\_384, densenet161.tv.in1k, dla102x.in1k, edgenext\_small\_rw.sw.in1k, regnety\_016.tv2.in1k, convnextv2\_base.fcmae\_ft.in1k, vit\_large\_patch14\_clip\_336\_laion2b\_ft.in12k.in1k, levit\_conv\_128s.fb\_dist.in1k, hrnet\_w48.ms.in1k, resnet101.a1h.in1k, xcit\_medium\_24\_p8\_224\_fb\_dist.in1k, resnetrs152.ft.in1k, convnextv2\_nano.fcmae\_ft.in22k.in1k, convnextv2\_tiny.fcmae\_ft.in22k.in1k, resnext50d\_32x4d.bt.in1k, gernet\_s.idstvc.in1k, seletcls42b.in1k, revpit\_m3.dist.in1k, resnet50d\_l.s4x24d.in1k, dpn98.mx.in1k, xcit\_nano\_12\_p16\_224\_fb.in1k, regnetx\_016.pycls.in1k, xcit\_medium\_24\_p16\_224\_fb.in1k, caformer\_s18.sail.in1k, sehalonet33ts.ra2.in1k, tinynet.c.in1k, xcit\_tiny\_24\_p16\_224\_fb\_dist.in1k, flexivit\_small\_300ep.in1k, resnext101\_32x8d.tv2.in1k, convnextv2\_base.fcmae\_ft.in22k.in1k\_384, semnasnet\_075.rmsp.in1k, res2net50\_26w\_4s.in1k, cait\_xs24\_384\_fb\_dist.in1k, mobilenetv2\_l120d.ra.in1k, seresnext26t\_32x4d.bt.in1k, flexivit\_base.1200ep.in1k, res2net50\_26w\_6s.in1k, vit\_base\_patch16\_clip\_384\_openai\_ft.in12k.in1k, nest\_base\_jx\_goog.in1k, ecaresnetlight.mil.in1k, repvgg\_b0.rvgg.in1k, ecaresnet50t.a1.in1k, inception\_next\_tiny.sail.in1k, regnety\_032.pycls.in1k, mixer\_b16\_224.mil.in12k\_ft.in1k, poolformer\_s12.sail.in1k, vit\_base\_patch32\_clip\_384\_openai\_ft.in12k.in1k, vit\_base\_patch32\_384\_augreg\_in21k\_ft.in1k, efficientvit\_b1\_r224.in1k, vit\_base\_patch16\_clip\_384\_laion2b\_ft.in1k, deit\_small\_distilled\_patch16\_224\_fb.in1k, efficientvit\_b0\_r224.in1k, resnet50d.in1k, regnety\_120.pycls.in1k, semnasnet\_100.rmsp.in1k, wide\_resnet50\_2.tv.in1k, xcit\_small\_24\_p16\_224\_fb.in1k, resnet101.a3.in1k, fastvit\_t12.apple.in1k, tf\_efficientnet\_lite1.in1k, tinynet.a.in1k, resmlp\_big\_24\_224\_fb\_distilled.in1k, cs3se.edgenet\_x.c2ns.in1k, resnetv2\_152x2\_bit\_goog\_teacher.in12k\_ft.in1k\_384, resnext50\_32x4d.tv2.in1k, efficientnet\_b2.ra.in1k, convformer\_s18.sail.in22k\_ft.in1k\_384, caformer\_s18.sail.in22k\_ft.in1k\_384, deit3\_base\_patch16\_224\_fb.in1k, vit\_base\_patch32\_clip\_384\_laion2b\_ft.in12k.in1k, vit\_medium\_patch16\_gap\_384.sw.in12k\_ft.in1k, sequencer2d.s.in1k, mobileone\_s0.apple.in1k, edgenext\_base.in21k\_ft.in1k, deit3\_medium\_patch16\_224\_fb.in1k, efficientformerv2\_l1.snap\_dist.in1k, lambda\_resnet50ts.a1h.in1k, xception4p.ra.in1k, resnext50\_32x4d.a3.in1k, crossvit\_small\_240.in1k, repvgg\_a1.rvgg.in1k, resnet51q.ra.in1k, xcit\_small\_12\_p16\_384\_fb\_dist.in1k, vit\_base\_patch32\_clip\_224\_openai\_ft.in1k, flexivit\_large\_300ep.in1k, repvgg\_b3g4.rvgg.in1k, resnext50\_32x4d.a1h.in1k, coat\_lite\_medium.in1k, vit\_base\_patch32\_clip\_448\_laion2b\_ft.in12k.in1k, resnext50\_32x4d.gluon.in1k, repvgg\_b2.rvgg.in1k, vit\_base\_patch16\_rpn\_224.sw.in1k, mixer\_b16\_224\_goog.in21k\_ft.in1k, resnet50.c2.in1k, lamhalobotnet50ts\_256.a1h.in1k, tiny\_vit\_21m\_512.dist.in22k\_ft.in1k, xcit\_large\_24\_p16\_224\_fb\_dist.in1k, repvgg\_a2.rvgg.in1k, gernet\_l.idstvc.in1k, mobilevitv2\_050.cvnets.in1k, convnextv2\_base.fcmae\_ft.in22k.in1k, resnet18.a3.in1k, ecaresnet50d.mil.in1k, coat\_lite\_small.in1k, convnext\_xlarge.fb.in22k\_ft.in1k, mobilevitv2\_075.cvnets.in1k, cait\_s36\_384\_fb\_dist.in1k, efficientformerv2\_s1.snap\_dist.in1k, resnet18.fb\_sws\_lig1b\_ft.in1k, mobileone\_s1.apple.in1k, resnet61q.ra2.in1k, tf\_efficientnetv2\_b3.in1k, mobilevitv2\_175.cvnets.in1k, convnext\_tiny.fb.in22k\_ft.in1k\_384, crossvit\_tiny\_240.in1k, caformer\_b36.sail.in22k\_ft.in1k\_384, resnet152d.ra2.in1k, convit\_base.fb.in1k, tinynet.b.in1k, deit3\_large\_patch16\_384\_fb.in22k\_ft.in1k, regnetz\_004.tv2.in1k, cait\_xs36\_384\_fb\_dist.in1k, convnext\_nano\_ols.d1h.in1k, efficientnet\_lite0.ra.in1k, inception\_v4.tf.in1k, hrnet\_w18.ms.in1k, gernet\_m.idstvc.in1k, convformer\_s36.sail.in22k\_ft.in1k\_384, deit\_tiny\_distilled\_patch16\_224\_fb.in1k, deit\_small\_patch16\_224\_fb.in1k, vit\_large\_patch14\_clip\_336\_laion2b\_ft.in1k, crossvit\_18\_240.in1k, resnet26.bt.in1k, resnet18.a1.in1k, deit3\_base\_patch16\_384\_fb.in22k\_ft.in1k, convformer\_s36.sail.in1k, convnext\_small.fb.in22k\_ft.in1k, seletcls60b.in1k, efficientnet\_b0.ra.in1k, focalnet\_tiny\_srf.ms.in1k, ecaresnet101d.mil.in1k, regnetx\_080.tv2.in1k, mobileone\_s3.apple.in1k, mobilenetv3\_rw.rmsp.in1k, poolformerv2\_m36.sail.in1k, seresnextaa101d\_32x8d.a1h.in1k, levit\_conv\_192.fb\_dist.in1k, focalnet\_tiny\_lrf.ms.in1k, regnety\_320.swag\_lc.in1k, tresnet\_v2\_l.mil.in21k\_ft.in1k,



## Measuring Mechanistic Interpretability at Scale Without Humans

880 seresnet50.a3.in1k, dla46x.c.in1k, cs3darknet.x.c2ns.in1k, tf.efficientnet\_b0.ap.in1k, vit\_base\_patch16\_224.augreg2.in21k.ft.in1k, resnext101\_32x8d.fb.ssl.yfcc100m.ft.in1k,  
881 xcit\_large\_24\_p8.384.fb.dist.in1k, tinynet.e.in1k, cait\_xs24.384.fb.dist.in1k, fastvit\_sai12.apple.in1k, hrnet\_w64.ms.in1k, regnety\_016.pycls.in1k, wide\_resnet101\_2.tv.in1k,  
882 beitv2\_large\_patch16\_224.in1k.ft.in22k.in1k, hrnet\_w30.ms.in1k, resnet101.tv.in1k, repvit\_m2.dist.in1k, coatnet.nano\_rw\_224.sw.in1k, flexivit\_small.1200ep.in1k,  
883 tf.efficientnet\_b0.in1k, tf.efficientnet\_b1.in1k, efficientformer\_l3.snap.dist.in1k, vit\_base\_patch16\_384.augreg.in21k.ft.in1k, xcit\_tiny\_24\_p8.224.fb.dist.in1k, dla102x2.in1k,  
884 hardcorenas.f.miil.green.in1k, regnety\_064.ra3.in1k, resnext101\_32x4d.gluon.in1k, tf.efficientnetv2\_b2.in1k, resnet32ts.ra2.in1k, xcit\_tiny\_12\_p8.384.fb.dist.in1k,  
885 inception\_v3.tv.in1k, xcit\_large\_24\_p16.224.fb.in1k, ecaresnet50t.a3.in1k, repvit\_m2\_3.dist.450e.in1k, fbnetv3\_b.ra2.in1k, vit\_base\_patch8\_224.augreg2.in21k.ft.in1k,  
886 cs3darknet.l.c2ns.in1k, convnext\_base.clip.laion2b.augreg.ft.in12k.in1k.384, regnety\_160.deit.in1k, regnety\_160.pycls.in1k, dla60x.in1k, xcit\_tiny\_24\_p16.384.fb.dist.in1k,  
887 eva02\_tiny\_patch14\_336.mim.in1k, flexivit\_small.600ep.in1k, visformer\_tiny.in1k, resnet50.a1.in1k, dla60.in1k, regnetz\_d32.ra3.in1k, senet154.gluon.in1k, efficient-  
888 netv2\_rw\_s.ra2.in1k, focalnet\_small\_lrf.ms.in1k, seresnet33ts.ra2.in1k, fbnetc\_100.rmsp.in1k, resnet18d.ra2.in1k, resnet34.a3.in1k, dla60x.c.in1k, efficient-  
889 net\_b1\_pruned.in1k, efficientformerv2\_s2.snap.dist.in1k, resnet50s.gluon.in1k, resnet101.a2.in1k, regnety\_040.ra3.in1k, convmixer\_1536.20.in1k, regnety\_008.tv.tv2.in1k,  
890 resnet152.a1.in1k, mixnet\_l.ft.in1k, gresnext26ts.ch.in1k, vit\_base\_patch16\_clip\_224.openai.ft.in1k, fastvit\_ma36.apple.in1k, vgg\_16.tv.in1k, gresnext50ts.ch.in1k,  
891 xcit\_tiny\_12\_p16.224.fb.dist.in1k, regnety\_008.pycls.in1k, resmlp\_36\_224.fb.distilled.in1k, regnetz\_040\_h.ra3.in1k, inception\_next\_base.sail.in1k, dm\_nfnnet\_f0.dm.in1k,  
892 resnet50\_d.in1k, efficientformerv2\_b2\_pruned.in1k, resnet18.tv.in1k, rexnet\_150.nav.in1k, convnext\_large\_mlp.clip.laion2b\_soup.ft.in12k.in1k.320, ghostnetv2\_160.in1k,  
893 vit\_small\_patch16\_384.augreg.in21k.ft.in1k, convnext\_xlarge.fb.in22k.ft.in1k.384, mobilenetv3\_small\_075.lamb.in1k, regnetz\_d8.evos.ch.in1k, dm\_nfnnet\_f3.dm.in1k,  
894 repvgg\_b3.rvvg.in1k, convnext\_large\_mlp.clip.laion2b.augreg.ft.in1k.384, dpn68b.mx.in1k, resnext101\_32x8d.fb.wsl.ig1b.ft.in1k, deit3\_large\_patch16\_384.fb.in1k,  
895 convformer\_s18.sail.in1k.384, repghostnet\_058.in1k, fastvit\_sa36.apple.dist.in1k, resnext50\_32x4d.a2.in1k, regnetx\_040.pycls.in1k, vit\_base\_r50\_s16.384.orig.in21k.ft.in1k,  
896 vit\_base\_patch16\_clip\_224.laion2b.ft.in1k, deit3\_base\_patch16\_224.fb.in1k, tf.efficientnetv2\_s.in1k, ecaresnet50.a2.in1k, resnetrs50.ft.in1k, resnet24\_224.ra3.in1k, resne-  
897 taa50d.sw.in12k.ft.in1k, tresnet\_xl.miil.in1k, resnet101e.in1k, regnetx\_004.pycls.in1k, mnasnet\_small.lamb.in1k, repvgg\_a0.rvvg.in1k, resnetv2\_50x1\_bit.goog.in21k.ft.in1k,  
898 cait\_s24\_224.fb.dist.in1k, regnety\_004.tv2.in1k, convnext\_base.fb.in22k.ft.in1k.384, convnext\_tiny.fb.in22k.ft.in1k, convnext\_tiny.in12k.ft.in1k.384, eca\_halonext26ts.c1.in1k,  
899 resnet18.gluon.in1k, fastvit\_s12.apple.dist.in1k, deit\_base\_patch16\_224.in1k, hrnet\_w18.ms.aug.in1k, resnet33ts.ra2.in1k, seresnext101\_64x4d.gluon.in1k, conv-  
900 next\_small.fb.in1k, convformer\_s36.sail.in1k.384, pit\_ti.distilled.224.in1k, resnet50.tv2.in1k, nest\_small\_jx.goog.in1k, resmlp\_36\_224.fb.in1k, hrnet\_w18\_small.gluon.in1k,  
901 vit\_base\_patch16\_384.augreg.in1k, resnet50.fb.swsl.ig1b.ft.in1k, poolformer\_m36.sail.in1k, tf.mobilenetv3\_small\_100.in1k, regnety\_040.pycls.in1k, gresnext33ts.ra2.in1k,  
902 resnet101s.gluon.in1k, darknetaa53.c2ns.in1k, poolformerv2\_s12.sail.in1k, resnext50\_32x4d.fb.ssl.yfcc100m.ft.in1k, poolformerv2\_s24.sail.in1k, eca\_resnet33ts.ra2.in1k,  
903 repvit\_m2\_3.dist.300e.in1k, nf\_resnet50.ra2.in1k, convnext\_pico\_ols.d1.in1k, caformer\_s36.sail.in1k, regnetz\_040.ra3.in1k, vit\_small\_r26\_s32\_224.augreg.in21k.ft.in1k,  
904 resnext26ts.ra2.in1k, mixnet\_xl.ra.in1k, deit\_base\_patch16\_384.fb.in1k, repvit\_m1\_0.dist.450e.in1k, convmixer\_1024\_20\_ks9\_p14.in1k, regnety\_064.pycls.in1k,  
905 resnet34.gluon.in1k, res2net101\_26w\_4s.in1k, nfnnet\_j0.ra2.in1k, resnet34d.ra2.in1k, convnextv2\_nano.fcmae.ft.in22k.ft.in1k.384, twins\_pcpvt\_base.in1k,  
906 resnetv2\_101.a1h.in1k, xcit\_nano\_12\_p8.224.fb.dist.in1k, xcit\_small\_24\_p8.224.fb.dist.in1k, resnet50.b2k.in1k, deit3\_small\_patch16\_384.fb.in1k, hardcorenas.c.miil.green.in1k,  
907 coat\_lite\_mini.in1k, resnet152.tv2.in1k, densenetblur121d.ra.in1k, hrnet\_w18\_small\_v2.gluon.in1k, vit\_base\_patch16\_384.orig.in21k.ft.in1k, xcit\_small\_12\_p8.224.fb.dist.in1k,  
908 convformer\_m36.sail.in1k, convformer\_m36.sail.in1k, xcit\_nano\_12\_p16.384.fb.dist.in1k, resnet34.a1.in1k, convnext\_atto\_ols.a2.in1k, resnet14t.c3.in1k, twins\_pcpvt\_large.in1k,  
909 resnet26d.gluon.in1k, mobilenetv3\_small\_100.lamb.in1k, efficientnet\_b3\_pruned.in1k, vit\_small\_patch16\_224.augreg.in1k, convnext\_tiny.fb.in1k, resnet50d.a3.in1k, mobilevitv2\_175.cvnets.in22k.ft.in1k.384,  
910 deit3\_medium\_patch16\_224.fb.in22k.ft.in1k, seresnext101\_32x4d.gluon.in1k, hardcorenas\_b.miil.green.in1k, caformer\_m36.sail.in22k.ft.in1k, ghostnetv2\_100.in1k,  
911 ecaresnet50d\_pruned.miil.in1k, caformer\_s36.sail.in22k.ft.in1k.384, deit\_tiny\_patch16\_224.fb.in1k, fastvit\_sa36.apple.in1k, regnety\_320.seer.ft.in1k, edgenext\_small.usi.in1k,  
912 resmlp\_big\_24\_224.fb.in22k.ft.in1k, regnety\_160.lion.in12k.ft.in1k, regnety\_160.sw.in12k.ft.in1k, tf.efficientnet\_b1.ap.in1k, res2net50\_48w\_2s.in1k,  
913 eca\_botnext26ts\_256.c1.in1k, xcit\_small\_24\_p8.224.fb.in1k, crossvit\_9\_dagger\_240.in1k, coat\_lite\_tiny.in1k, wide\_resnet50\_2.tv2.in1k, vit\_base\_patch16\_clip\_224.openai.ft.in12k.in1k,  
914 skresnet34.ra.in1k, repvgg\_b1g4.rvvg.in1k, vgg19.bn.tv.in1k, repghostnet\_100.in1k, regnetv\_064.ra3.in1k, mobilenetv2\_100.ra.in1k, convnext\_femto.d1.in1k,  
915 resnet26t.ra2.in1k, regnetv\_040.ra3.in1k, skresnet18.ra.in1k, caformer\_m36.sail.in22k.ft.in1k.384, vit\_base\_patch32\_384.augreg.in1k, regnetz\_b16.ra3.in1k, hrnet\_w48\_ssl.paddle.in1k, resnet50d\_4s2x40d.in1k, cait\_xxs36\_224.fb.dist.in1k, regnetx\_016.tv2.in1k,  
916 xcit\_small\_24\_p8.384.fb.dist.in1k, vit\_tiny\_r\_s16\_p8.224.augreg.in21k.ft.in1k, coat\_mini.in1k, xcit\_small\_24\_p16.224.fb.dist.in1k, caformer\_s36.sail.in22k.ft.in1k, poolformer\_s24.sail.in1k,  
917 resmlp\_big\_24\_224.fb.in1k, regnetx\_120.pycls.in1k, regnetz\_d8.ra3.in1k, resnet50d.ra2.in1k, repvit\_m1.dist.in1k, eca\_nfnnet\_l2.ra3.in1k, resnet50d.gluon.in1k, seresnext50\_32x4d.racm.in1k,  
918 vit\_small\_patch16\_384.augreg.in1k, coat\_tiny.in1k, xcit\_nano\_12\_p8.224.fb.in1k, crossvit\_base\_240.in1k, resnet50d.a1.in1k, convformer\_s36.sail.in22k.ft.in1k, convnextv2\_large.fcmae.ft.in22k.in1k, resnet50.tv.in1k, resnet50.c1.in1k, pit\_xs.distilled.224.in1k, efficient-  
919 net\_b1.ft.in1k, tf.efficientnet\_el.in1k, hrnet\_w32.ms.in1k, vit\_base\_patch16\_224.miil.in21k.ft.in1k, cs3darknet.x.c2ns.in1k, dpn68b.ra.in1k, tf.efficientnetv2\_b1.in1k, regnety\_004.pycls.in1k,  
920 tf.mobilenetv3\_large\_minimal\_100.in1k, resnetrs101.ft.in1k, ese\_vovnet39b.ra.in1k, mixer\_l16\_224.goog.in21k.ft.in1k, repghostnet\_050.in1k, repvgg\_b2g4.rvvg.in1k, repvit\_m1\_l1.dist.450e.in1k,  
921 vit\_base\_patch32\_224.augreg.in21k.ft.in1k, tf.mobilenetv3\_large\_100.in1k, pit\_s.224.in1k, caformer\_s18.sail.in22k.ft.in1k, wide\_resnet101\_2.tv.in1k, fastvit\_l12.apple.dist.in1k, convmixer\_768.32.in1k, vit\_base\_patch32\_224.augreg.in1k, efficientformerv2\_s0.snap.dist.in1k, resnet200e.in1k,  
922 levit\_conv\_256.fb.dist.in1k, resnet18.fb.ssl.yfcc100m.ft.in1k, vgg13.bn.tv.in1k, resnet152c.gluon.in1k, dla169.in1k, pvt\_v2\_b4.in1k, crossvit\_15.dagger.408.in1k, conv-  
923 next\_femto\_ols.d1.in1k, convnext\_large.fb.in1k, regnetx\_064.pycls.in1k, fastvit\_t8.apple.in1k, seresnet152d.ra2.in1k, vgg19.tv.in1k, vgg11.bn.tv.in1k, dm\_nfnnet\_f2.dm.in1k, seresnext101d\_32x8d.ah.in1k,  
924 inception\_next\_base.sail.in1k.384, lambda\_resnet26t.c1.in1k, resnetv2\_152x2\_bit.goog.in21k.ft.in1k, fastvit\_ma36.apple.dist.in1k, regnety\_006.pycls.in1k, regnety\_080.pycls.in1k, resnet50.fb.ssl.yfcc100m.ft.in1k, tf.mobilenetv3\_small\_075.in1k, regnetz\_c16.ra3.in1k, edgenext\_xx\_small.in1k, crossvit\_9\_240.in1k,  
925 xcit\_tiny\_24\_p8.224.fb.in1k, regnety\_080.ra3.in1k, efficientvit\_b1\_r256.in1k, tinynet.d.in1k, caformer\_b36.sail.in1k.384, pvt\_v2\_b2.in1k, resnet26d.bt.in1k, convnext\_pico.d1.in1k,  
926 pit\_b.224.in1k, convnextv2\_pico.fcmae.ft.in1k, fbnetv3\_d.ra2.in1k, flexivit\_large.1200ep.in1k, resnet50c.gluon.in1k, regnetx\_080.pycls.in1k, convnext\_base.fb.in1k, tf.efficientnet\_em.in1k, vit\_base\_patch16\_224.augreg.in1k, convit\_tiny.fb.in1k, resnext50\_32x4d.fb.swsl.ig1b.ft.in1k, dm\_nfnnet\_f4.dm.in1k, resnet50.a3.in1k,  
927 convnext\_atto.d2.in1k, efficientnet\_el\_pruned.in1k, volo\_d2.384.sail.in1k, resnext101\_32x4d.fb.ssl.yfcc100m.ft.in1k, repvit\_m0\_9.dist.300e.in1k, regnety\_120.sw.in12k.ft.in1k, beit\_base\_patch16\_384.in22k.ft.in22k.in1k, mobilenetv3\_large\_100.miil.in21k.ft.in1k, tf.efficientnet\_b0.a1.in1k, inception\_next\_small.sail.in1k,  
928 deit\_base\_distilled\_patch16\_224.fb.in1k, lcnets\_075.ra2.in1k, xcit\_tiny\_12\_p8.224.fb.in1k, resnet101.gluon.in1k, dpn92.mx.in1k, resnet101.a1.in1k, sececls60.in1k, beit\_base\_patch16\_224.in22k.ft.in22k.in1k, convnextv2\_tiny.fcmae.ft.in1k, res2net50\_26w\_8s.in1k, sequencer2d.m.in1k, vit\_medium\_patch16\_gap\_256.sw.in12k.ft.in1k, regnetx\_008.pycls.in1k,  
929 resnet50.a2.in1k, resnet101d.in1k, vit\_large\_patch16\_384.augreg.in21k.ft.in1k, pvt\_v2\_b2.li.in1k, regnetx\_006.pycls.in1k, xcit\_tiny\_24\_p16.224.fb.in1k, pvt\_v2\_b5.in1k, resnext50\_32x4d.ra.in1k, resnet14d.gluon.in1k, caformer\_m36.sail.in1k.384, resnet50.gluon.in1k, resnet152s.gluon.in1k, flexivit\_large.600ep.in1k, resnetv2\_50x1\_bit.goog.distilled.in1k, resmlp\_24\_224.fb.in1k, deit3\_large\_patch16\_224.fb.in1k, seresnext50\_32x4d.gluon.in1k, densenet121.tv.in1k, resnet152.a3.in1k, ghostnet\_100.in1k, tf.efficientnet\_b2.ap.in1k, regnetx\_002.pycls.in1k.

G. Additional Results

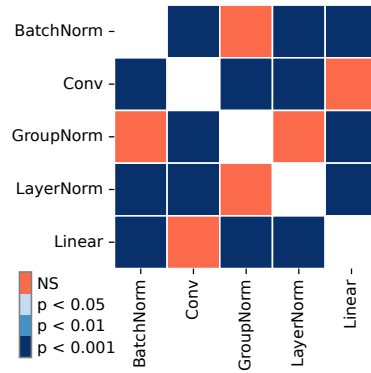


Fig. 20: **Differences Between Layer Types are Significant.** We analyze and test for statistical significances in the differences in MIS between different layer types (see Fig. 8). The reported significance levels were computed using Conover’s test over the per-model and per-layer-type means with Holm’s correction for multiple comparisons.

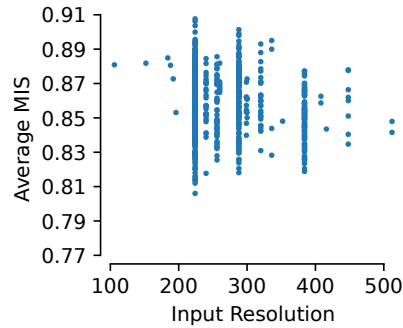
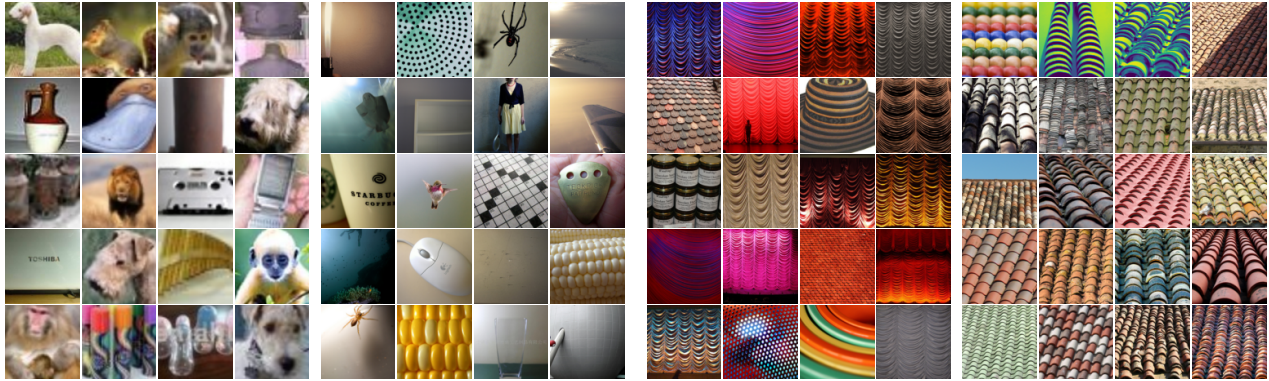


Fig. 21: **Influence of Input Resolution of MIS.** We show the average MIS per model as a function of the model’s input resolution. No trend is apparent; models with the same resolution yield different interpretability levels.

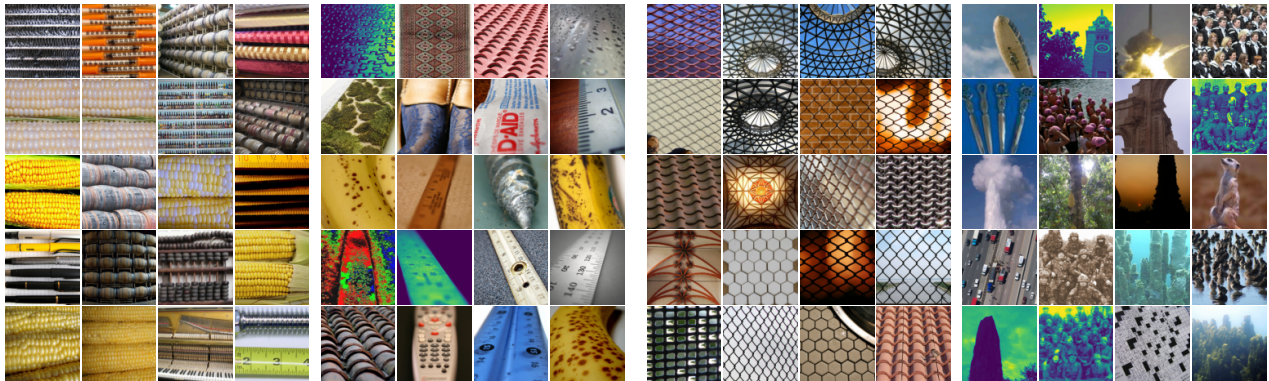


1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099



(a) Clip ViT-B32, resblocks\_3\_mlp\_c\_proj, unit 573

(b) DenseNet201, block3\_layer10\_norm1, unit 138



(c) DenseNet201, block3\_layer29\_conv1, unit 39

(d) DenseNet201, block3\_layer35\_norm2, unit 123

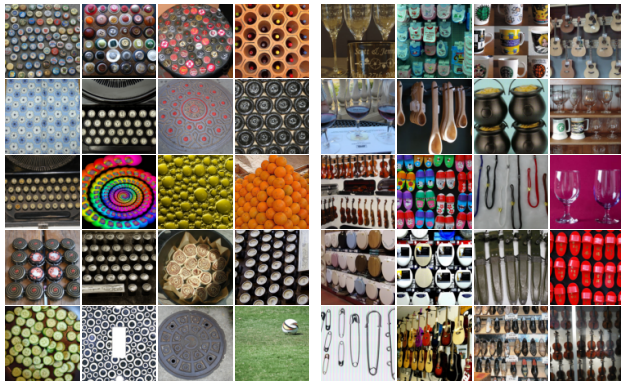
Fig. 24: **Visualization of Units for which MIS overestimates HIS.** To showcase the shortcomings of the MIS, we visualize four units for which the MIS predicts an interpretability that is higher than the measured HIS in Fig. 3. See Fig. 25 for the opposite direction. For each unit, we show the 20 most (right) and 20 least (left) activating dataset exemplars.

1100  
1101  
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  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154

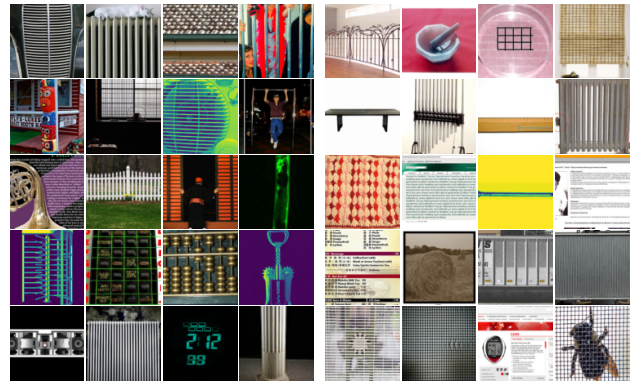


(a) Clip ResNet-50, layer4\_1\_conv2, unit 430

(b) DenseNet201, block3\_layer48\_conv2, unit 21



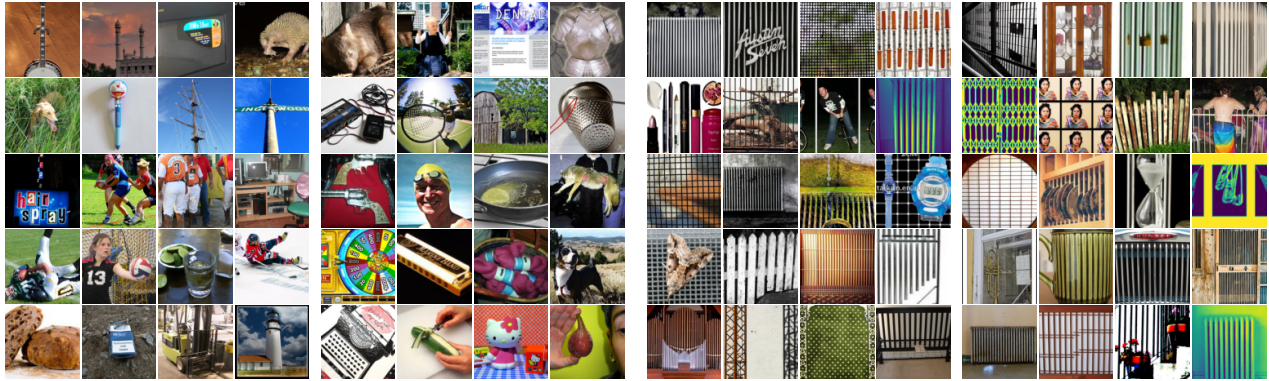
(c) DenseNet201, block3\_layer48\_norm1, unit 1369



(d) ViT-B32, block0\_norm2, unit 358

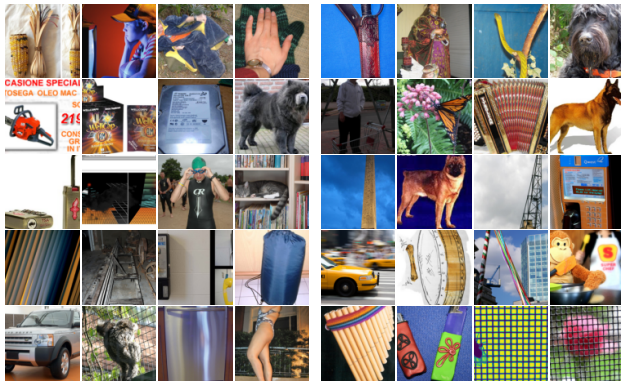
Fig. 25: **Visualization of Units for which MIS underestimates HIS.** To showcase the shortcomings of the MIS, we visualize four units for which the MIS predicts an interpretability that is lower than the measured HIS in Fig. 3. See Fig. 24 for the opposite direction. For each unit, we show the 20 most (right) and 20 least (left) activating dataset exemplars.

1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209

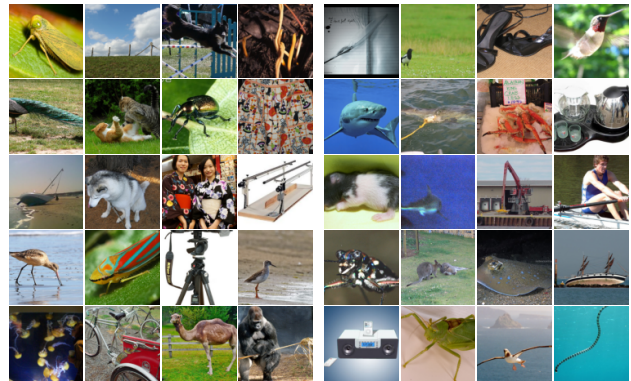


(a) ResMLP-36, blocks\_10\_linear\_tokens, unit 61

(b) ResMLP-24, blocks\_0\_mlp\_channels\_fc1, unit 110



(c) GMixer-24, blocks\_5\_mlp\_tokens\_fc1, unit 166



(d) ResMLP-12, blocks\_7\_linear\_tokens, unit 127

**Fig. 26: Visualization of Hard Units from Models with High Variability.** For the four models with the highest variability in MIS (see Fig. 6), we visualize one of the units with the lowest MIS each. For each unit, we show the 20 most (right) and 20 least (left) activating dataset exemplars.