

Explaining Local, Global, And Higher-Order Interactions In Deep Learning

Sam Lerman
University of Rochester
Rochester, NY
slerman@ur.rochester.edu

Charles Venuto
University of Rochester Medical Center
Rochester, NY
Charles.Venuto@chert.rochester.edu

Chenliang Xu
University of Rochester
Rochester, NY
chenliang.xu@rochester.edu

Henry Kautz
National Science Foundation
Washington, D.C.
hkautz@nsf.gov

Abstract

We present a simple yet highly generalizable method for explaining interacting parts within a neural network’s reasoning process. First, we design an algorithm based on cross derivatives for computing statistical interaction effects between individual features, which is generalized to both 2-way and higher-order (3-way or more) interactions. We present results side by side with a weight-based attribution technique, corroborating that cross derivatives are a superior metric for both 2-way and higher-order interaction detection. Moreover, we extend the use of cross derivatives as an explanatory device in neural networks to the computer vision setting by expanding Grad-CAM, a popular gradient-based explanatory tool for CNNs, to the higher order. While Grad-CAM can only explain the importance of individual objects in images, our method, which we call Taylor-CAM, can explain a neural network’s relational reasoning across multiple objects. We show the success of our explanations both qualitatively and quantitatively, including with a user study. We will release all code as a tool package to facilitate explainable deep learning.

1. Introduction

The universe is made up of myriad interacting parts. To truly understand complex systems and processes, it is not enough to view their functions as an amalgamation of independent contributors. Rather, they are a complex web of inter-operating influences. For much of the past, explainable deep learning has concerned itself with identifying important features, feature vectors, and isolated concepts. However, in the real world, humans intuitively understand that decisions are consequences of complex relations, not merely extrapolated from rankings of singular phenomena.



Figure 1: An automated driver decides whether to “stop” or “go.” Here, the decision cannot be explained by individual factors alone, but by the interaction between the yield sign and the passing car. Taylor-CAM identifies interactions by considering how changing one object affects the significance of another, such as how changing a passing car into an empty road would change the meaning of the yield sign from “stop” to “go.”

For example, upon seeing a yield sign, it is natural to look to see if there are also passing cars. If not, the yield sign may be safely dismissed and one could keep driving without stopping. If there is a passing car, the law is to yield to the other car. If an intelligent agent made the decision to stop upon approaching a yield sign and a passing car, explaining their actions with precision would require an explanation of this interaction. As far as individual factors go, perhaps a nearby pedestrian is also present, but without an interactional interpretation, one would not be able to distinguish the independence of the yield sign and passing

car from the pedestrian, and one would not be privy to the knowledge of the salient interaction. Furthermore, a naive observer might think that yield signs always indicate “stop” without realizing that the agent’s response to the yield sign would depend on the presence of a passing car.

Similarly, explaining an agent’s strategies in any task — be it computer vision, natural language processing, biomedicine, reinforcement learning, or future forecasting — is imprecise without an interactional approach. However, interactional strategies are not always summarizable by heatmaps [5, 32, 33, 47, 48] or ordered rankings [12, 29, 37]; and they often require an understanding of many dependencies — complex dependencies, such as those between higher-level concepts (*e.g.* vector representations in deep neural networks [2, 30, 31, 46]) — not just single-dimensional features as typically explored in the statistical interaction effects literature [9, 17, 40, 41]. In light of all of this, we propose a number of contributions towards explaining interactions in deep learning:

T-NID, an algorithm for statistical interaction effects that outperforms recent state-of-the-art baselines with both pairwise and higher-order interactions. Interaction effects are a fundamental notion in statistics [45]. We make this computation tractable by translating local interaction effects into global interaction effects via representative samples and employing a simple subsampling heuristic.

Taylor-CAM, an explanatory tool that extends Grad-CAM [32], which assigns importances to feature vectors based on input gradients, by generalizing it to the 2-way and higher-order setting using the same formalism of interaction effects as for T-NID. This method is demonstrated on multi-object detection and relational reasoning in visual question-answering (VQA).

Visualizations of Taylor-CAM’s explanations that enable a human cohort to reverse engineer questions in relational VQA without knowing the answers and interpret relational reasoning better than with existing explanatory tools like Grad-CAM and GLIDER [40] from just a convolutional neural network’s (CNN) feature maps.

2. Related Work

In Deep Learning Recently, there have been several attempts to compute statistical interactions with deep learning. Neural Interaction Detection (NID) [41] used neural network weights to interpret interactions, observing that interactions occur at nonlinear activations in the first hidden layer of an MLP. Like our approach T-NID, [6] used gradient information to compute statistical interaction effects. However, they relied on Bayesian neural networks, required averaging a high number of Hessians, and only computed global interaction effects, not focusing on local or higher-order interactions. [9] used cross derivatives between single features to explain interactions in deep similarity models,

whereas we use an adaptation of Grad-CAM to demonstrate explainability in a more general computer vision setting. [35] relied on self attention [42] to compute a measure analogous to non-emergent interaction effects and apply this to an analysis in the biomedical domain.

Cui et al. [6] applied their approach to a toy MNIST dataset consisting of a fixed set of feature vectors such that they could compute global interaction effects, but they mapped those feature vectors to single neurons and computed standard interaction effects between those mapped neurons. The limitation of this approach is that it cannot be used to explain local phenomena, which is traditionally what is of interest in computer vision, NLP, and other areas where multidimensional feature vectors are used.

[17] and [40], like our substitution of ReLU with GELU, substitute ReLU with Sofplus in order to induce differentiability. The latter, like our work, translate local interaction effects to global interaction effects by aggregating across representative samples. While they use a random batch, we use a small subset of common aggregates. While our Taylor-CAM formulation is expressly adapted from Grad-CAM for intuitively explaining feature vectors in CNNs, [17] derive their formulation from integrated gradients and [40] directly uses cross partials.

Individual Importances [12, 29, 37] used input gradients to explain the reasoning of a neural network. [48] did so with class activation maps. Grad-CAM [32] and Grad-CAM++ [5] combined both approaches to localize important feature vectors in computer vision with class activation maps and gradients, visualized by heatmaps. Similar to us, [23] used Taylor decomposition to explain neural network decisions, but only for main effects, not interactions.

Relational Reasoning We also connect interaction effects with relational reasoning, which has received increased attention in deep learning [2, 30, 31, 46], and use Taylor-CAM to interpret the reasoning process of Relation Networks [31]. While most past works have mainly focused on explaining individual factors of a neural network’s predictions, the weights in multi-head dot product attention [42] could be interpreted as relational explanations for neural networks that include MHDPA in their architecture [35]. In contrast, Taylor-CAM is architecture agnostic and can explain decisions unique to each output dimension directly from gradient information.

Unlike other works, we expressly derived Taylor-CAM for the purpose of explaining interactions between higher level representations, such as feature maps from a CNN, which standardly represent objects in computer vision (rather than using raw RGB pixels). As Grad-CAM is built on projected feature vectors in addition to gradients, so is our higher-order extension w.r.t. cross derivatives to explain interactions rather than isolated phenomena.

3. Statistical Interaction Effects

We define statistical interaction effect analogous to [1]:

Definition 3.1. Interaction Effect An interaction effect $\text{IE}_{1,\dots,\ell}$ between variables $x_1, \dots, x_\ell \in \mathbf{x}$ on a function $F(\mathbf{x})$ with inputs \mathbf{x} is measured as:

$$\text{IE}_{1,\dots,\ell} = \frac{\partial^\ell F(\mathbf{x})}{\partial x_1 \cdots \partial x_\ell}. \quad (1)$$

In plain English, an interaction effect is how much the meaning of one variable changes for a unit change in another variable. This change is reflected by the cross partial derivative. “Change” is an intuitive measure for interaction. From the earlier example, given a representation of a yield sign and an oncoming car, *changing* the representation of the oncoming car into a representation of an empty road also changes the meaning of the yield sign from “stop” to “go.” For a more formal example, consider $F(\mathbf{x}) = x_1 \sin(x_2) + \cos(x_3)$. F consists of an interaction between x_1 and x_2 for some \mathbf{x} since $\partial^2 F(\mathbf{x}) / (\partial x_1 \partial x_2)$ is nonzero. However, x_3 does not belong to an interaction since any cross derivative w.r.t. x_3 is zero.

Adapt to Neural Networks Substituting F with a trained neural network, we can compute the local interaction effects for a datapoint up to order ℓ as long as the neural network F is ℓ -times differentiable. In classification, softmax ensures this to be the case. In regression, we substitute ReLUs with Gaussian-error rectified linear units (GELUs), which have been shown comparable in performance [13]. Otherwise, Definition 3.1 affords the computation of interaction effects for arbitrary neural network architectures.

Translate Local Effects to Global Effects Often in statistics, there is greater interest in computing global interaction effects, statistics that generalize across all datapoints. Similarly, this need may be found in analyzing scene graphs, object co-currency, and contextual information [28, 34, 43]. In tandem with our work, [6] converted local pairwise interaction effects to global pairwise interaction effects by averaging a set of representative samples retrieved via k-means clustering, in effect dividing the dataset by Euclidean distance and computing the global average from the centroids. We will similarly average representative local interaction effects in order to compute a global summary, but we will use a simpler and more efficient technique. In our case, efficiency is of more concern because computing higher-order interaction effects requires the computation of higher-order derivatives, which for many samples can become intractable.

To translate local interaction effects into global ones at any order, we sample representative samples that have a wide range over the dataset and that are potentially meaningful. We choose the samples that are closest to a subset of common aggregates, including mean, median, min, max,

and mode. As well as a random sample for good measure. Likewise, we used L2 distance to measure closeness. In addition, we considered different ways to aggregate the interaction effects of these samples. Again, namely mean, median, min, max, or mode. We ran a wide sweep of the complete power set of these potential samples and aggregates to find which combination performed best on a wide array of synthetic datasets distinct from those we trained on selected from prior works [15, 22, 36, 41], chosen to test for various types of interactions. Results of this power sweep are reported in the *Appendix*. We ended up using the mean interaction effect of the samples closest to the mean, minimum, and mode of all samples, as well as a random sample.

Improve Efficiency Another heuristic for efficiency that we employed was subsampling the interactions that would be computed. Naturally, testing for every combination up to order ℓ would be very expensive. Every double, every triple, every quadruple, etc. — the problem grows combinatorially. We were able to mitigate this to a degree by taking advantage of the property of statistical interaction effects that *an ℓ -way interaction can only exist if all its corresponding $(\ell - 1)$ -interactions exist* [36]. In turn, we were able to reduce the search space by only selecting non-redundant combinations of the k interactions from the previous order whose interaction effects were highest, beginning with using every combination up to order o and then subsampling the top k for every order thereafter.

Our complete algorithm, which we call Taylor-Neural Interaction Detection (T-NID) due to the higher-order derivatives, is described in pseudocode in the *Appendix*.

Finally, we need to make a point about the sign of the resulting cross partial derivatives. A positive value indicates change in the positive direction; negative, negative. Since in regression we are interested in the overall effect of an interaction and are agnostic to the direction, we take the squared value of the cross-partial as our measure of interaction effect. In contrast, for classification, we use the sign — positive or negative — corresponding to the class of interest. And for multi-class classification, we take F to be the network corresponding to the class output of interest, and use its squared cross partial derivatives.

4. Taylor-CAM

To this point, we have generalized our computation of interaction effects to the local, global, and higher-order setting, but we have not yet considered the case where features are multidimensional, as is the case in higher-level deep neural network representations.

Explaining the influence of feature vectors is common in computer vision and interpreting CNNs. However, we have illustrated with multiple examples why a precise explanation of a model’s decisions requires an explanation of its interacting components, not just singular entities.

4.1. Intuition

For arbitrary objects in the computer vision setting, a cross derivative alone is not sufficient. Besides the obvious reason that such objects are not represented by singular features but by multidimensional feature vectors learned by a CNN, it is also because fundamentally a cross derivative measures changes of changes. More formally, a cross derivative $\frac{\partial^2 F}{\partial x \partial y}$ measures the effect of a unit change of x on the effect on F of a unit change of y . When reasoning about visual relations, it is convenient to think of dependencies between objects that inform a decision, such as the dependency between a yield sign and a passing car in informing an automated driver’s decision to “stop” or “go.” *Changing* the passing car into another object, such as merely an empty road, would on its own change the neural network’s interpretation of the yield sign from meaning “stop” to meaning “go,” even while keeping the yield sign fixed and unchanged — yet a cross derivative only measures the effect of changing both. To account for this, instead of naively using cross derivatives, we measure how much changing one object would change the *importance* of another object to a neural network’s decision, *e.g.*, how changing the yield sign into a speed limit sign would change the passing car’s importance or how changing the passing car into a gush of leaves would change the yield sign’s importance with regards to the decision of whether to “stop” or “go” — even when not necessarily both are changed.

Given car C , yield sign Y , and binary decision “go” G , this intuition may be summarized mathematically as:

$$S_{Y,C} = \partial \text{IMP}(Y, G) / \partial C, \quad (2)$$

where $S_{Y,C}$ represents the interaction salience between the yield sign and passing car, and $\text{IMP}(Y, G)$ represents the importance of the yield sign to the neural network’s decision to go or stop. Fortunately, the importance of individual objects in computer vision is the characteristic problem of the explanatory tool Grad-CAM [5, 32, 48], which we use to derive our method. We use the term *interaction salience* due to deviation from interaction effects in Definition 3.1.

4.2. Methodology

Suppose we have an ℓ -times differentiable function $F : \mathbb{R}^{n,d} \rightarrow \mathbb{R}$, which will stand for our neural network, where $\ell \geq 2$. F takes in matrix \mathbf{x} consisting of n feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ of dimension d . So $\mathbf{x}_1, \dots, \mathbf{x}_n$ are just feature vectors produced by a CNN and each one is associated with an image region. F is the portion of the network downstream of these feature vectors.

Quantify Importance To fill IMP in Equation 2, we turn to class activation maps (CAMs) [48]. However, as observed by the solution of [32], to find out how a class activation map increases the class’s likelihood, we would

like to know how its features contribute to the output, which we can do with their gradients. We can estimate the global effect by summing the gradient of each feature vector \mathbf{x}_k and weighing the sum to each CAM. This amounts exactly to Grad-CAM [32]:

$$\begin{aligned} \text{IMP}(\mathbf{x}_i, F(\mathbf{x})) &= \text{GradCAM}(\mathbf{x}_i, F(\mathbf{x})) \\ &= \sum_p \mathbf{x}_{ip} \sum_k \frac{\partial F(\mathbf{x})}{\partial \mathbf{x}_{kp}} . \end{aligned} \quad (3)$$

Generalize Grad-CAM to Compute Interactions Now that we have the importance of a feature vector (via essentially Grad-CAM), we can formulate S_{ij} , the interaction salience between feature vectors \mathbf{x}_i and \mathbf{x}_j , by substituting Equation 3 into 2 and summing the dimensions as follows:

$$S_{ij} = \sum_m \partial \left[\sum_p \mathbf{x}_{ip} \sum_k \frac{\partial F(\mathbf{x})}{\partial \mathbf{x}_{kp}} \right] / \partial \mathbf{x}_{jm}. \quad (4)$$

Merge with Statistical Interaction Effects Finally, we bring this to an easy-to-compute form by realizing that the partial derivative in the denominator $\partial \mathbf{x}_j$ can be computed together with the partial derivative in the numerator. We also square the salience because a change of importance in either direction would be significant. We note that the following is a generalization of Grad-CAM that reduces elegantly to a modified interaction effects Definition 3.1:

$$\begin{aligned} S_{ij}^2 &= \left(\sum_m \sum_p \mathbf{x}_{ip} \sum_k \frac{\partial^2 F(\mathbf{x})}{\partial \mathbf{x}_{kp} \partial \mathbf{x}_{jm}} \right)^2 \\ &= \left(\sum_{m,p,k} \mathbf{x}_{ip} \mathbf{IE}_{kp,jm} \right)^2 . \end{aligned} \quad (5)$$

In tests, we found setting $k = i$ in Equations 3 - 5 without the global sum over k to perform just as well and often better, perhaps because the local gradients in Equation 3 more precisely correspond to features. We call Equation 5 Hessian-CAM. Hessian-CAM may be further differentiated with respect to a cross partial $\partial \mathbf{x}_q$ to get a 3-way interaction salience, and that can be further differentiated up to any order ℓ . Thus, we name this Taylor-CAM, a higher-order generalization of Grad-CAM, where Grad-CAM (or a close variant) is the special case $\ell = 1$ and Hessian-CAM is the special case $\ell = 2$.

Note that interaction saliences are conditional. The interaction salience of feature \mathbf{x}_i on feature \mathbf{x}_j is not necessarily the same as that of \mathbf{x}_j on \mathbf{x}_i . Interaction salience S_{ij} represents the influence of \mathbf{x}_i on the importance of \mathbf{x}_j . Interaction salience $S_{ijk\dots}$ represents the influence of \mathbf{x}_i on the interaction salience of interaction $\mathbf{x}_j, \mathbf{x}_k, \dots$. To address this, we sum the mutual pairs, *e.g.*, $S_{ij} + S_{ji}$, although we note that we did so only to make the presentation clearer and

not because it is required. For many interpretation tasks, understanding that the meaning of the yield sign depends on the car, but the meaning of the car does not depend on the yield sign is crucial to getting the most precise understanding. Computing the mutual pairs does not require re-computation of any derivatives, and can be achieved easily by permuting the resulting interaction saliences and summing them. Lastly, we zero out the diagonals and redundant grid cells of the resulting interaction saliences to only consider interactions between non-redundant feature vectors.

4.3. Limitations

One limitation of Taylor-CAM, much like Grad-CAM, is that “importance” is based on contribution to the output, so if two different objects have the same contribution to the output, then changing one into the other would be considered meaningless, and so the interactions might not be identified. Suppose we have the setup from Sort-Of-CLEVR [18], a relational reasoning task. Here, we have an image with an assortment of shapes of different colors and a relational question related to that image. An example of this limitation is when an agent is asked, “What is the color of the circle furthest from the pink square?” If the furthest circle is blue, and the second furthest is also blue, then changing the furthest into a square does not meaningfully impact the pink square’s contribution to the output, as determined by Grad-CAM, since the answer to the question would be unchanged (blue). Grad-CAM++ [5] may hold an insight as to how to address this, via even-higher order derivatives.

Another limitation is that “change” is measured locally, as derivatives do not account for non-local rates of change. This means that Taylor-CAM, like other deep learning explanatory tools, depends on local regions of representations.

Lastly, of course, is the time complexity of computing higher-order derivatives. Higher-order differentiation has become increasingly more accessible with Taylor-mode autograd methods like JAX [3] and libraries like the new Pytorch functional autograd API [24], yet remains a challenge as the order grows. For Hessian-CAM, we had no trouble computing 2nd-order derivatives of Relation Networks using Pytorch and CPU memory. None of our individual explanations required more than a few minutes to compute on a CPU, excluding neural network training.

5. Experiments

5.1. Statistical Interaction Effects

We evaluate T-NID’s ability to rank interactions on the suite of synthetic functions proposed by [15, 22, 36, 41], which were “designed to have a mixture of pairwise and higher-order interactions, with varying order, strength, non-linearity, and overlap” [41]. These are available to see in the *Appendix* and in Table 1 of [41].

Pairwise Interactions For pairwise interaction effects (see Table 1), we report or reproduce the experiments of [41] verbatim, measuring AUC scores between predicted interaction rankings and ground truths. A pair x_i, x_j is considered an interaction either by itself or when it is a subset of a higher-order interaction, as in [22, 36]. Included for comparison are benchmarks from various statistical and machine learning methods [36, 38, 40, 41, 45]. NID [41] uses an interpretation of the weights from a standard MLP to detect interactions, whereas NID + MLP-M uses an MLP with additional univariate networks summed at the output to discourage modeling of main effects and false spurious interactions. GLIDER [40] is a recent cross-partial method that induces higher-order differentiability with Softplus.

In contrast, T-NID uses only a standard MLP and GELU activations. GELU demonstrably performs better. Unlike NID, we found no benefit from MLP-M or sparsity regularization. Despite the simpler architecture, T-NID is immune to some of the deficits of NID and NID + MLP-M. T-NID is able to distinguish main effects and spurious interactions in F_2 and F_4 , and while NID + MLP-M modeled spurious main effects in the $\{8, 9, 10\}$ interaction of F_6 and GLIDER appears to struggle with this as well, T-NID recognizes it as an interaction. All around, T-NID performs on par or better than NID and GLIDER at computing pairwise statistical interaction effects on these synthetic tasks.

Higher-Order Interactions For higher-order interactions, we do not report AUC scores against the full ground truth, as that would grow combinatorially more expensive with higher orders. Since NID also extracts interactions one order at a time, we compare the AUC scores of NID and T-NID one order at a time and use ground truths from the union of their discovered interactions. That way, they can be assessed relative to one another, albeit not universally. In addition to the results reported in Table 2, we tested many variants of architectures and report results with NID + MLP-M in the *Appendix*. In all cases, the relative results were largely the same, with T-NID achieving the highest scores, except less so at 4-way interactions when equipped with its own main effects network (MLP-M). Since any-order NID tends to find supersets much better than subsets, at 3-way interactions, NID misses nearly all present interactions, whereas T-NID fares relatively well. Along with recent works [6], we have shown that cross derivatives are a promising metric for interaction attribution in DNNs.

5.2. Object Detection

We ran two qualitative assessments of Taylor-CAM in multi-object detection. In both, the task was to identify whether a pair of objects were present in tandem. We tested the objects “car” and “person” in the COCO annotated-image dataset [21], and we designed our own toy dataset consisting of cars (rectangles), signs (triangles), and a yield

Table 1: AUC scores for pairwise interaction effects. Top-1 scores are bolded.

	ANOVA	HierLasso	RuleFit	AG	NID [41]	NID MLP-M [41]	GLIDER [40]	T-NID
$F_1(\mathbf{x})$	0.992	1.00	0.754	1	0.970	$0.995 \pm 4.4e - 3$	0.973 ± 0.01	0.962 ± 0.022
$F_2(\mathbf{x})$	0.468	0.636	0.698	0.88	0.79	$0.85 \pm 3.9e - 2$	0.84 ± 0.097	0.885 ± 0.039
$F_3(\mathbf{x})$	0.657	0.556	0.815	1	0.999	1 ± 0.0	0.919 ± 0.075	0.999 ± 0.001
$F_4(\mathbf{x})$	0.563	0.634	0.689	0.999	0.85	$0.996 \pm 4.7e - 3$	0.951 ± 0.073	0.998 ± 0.003
$F_5(\mathbf{x})$	0.544	0.625	0.797	0.67	1	1 ± 0.0	0.997 ± 0.008	0.991 ± 0.016
$F_6(\mathbf{x})$	0.780	0.730	0.811	0.64	0.98	$0.70 \pm 4.8e - 2$	0.767 ± 0.033	0.954 ± 0.026
$F_7(\mathbf{x})$	0.726	0.571	0.666	0.81	0.84	$0.82 \pm 2.2e - 2$	0.751 ± 0.207	0.98 ± 0.021
$F_8(\mathbf{x})$	0.929	0.958	0.946	0.937	0.989	$0.989 \pm 4.5e - 3$	0.998 ± 0.005	1.0 ± 0.0
$F_9(\mathbf{x})$	0.783	0.681	0.584	0.808	0.83	$0.83 \pm 3.7e - 2$	0.754 ± 0.098	0.98 ± 0.023
$F_{10}(\mathbf{x})$	0.765	0.583	0.876	1	0.995	$0.99 \pm 2.1e - 2$	0.974 ± 0.027	1.0 ± 0.0
Average	0.721	0.698	0.764	0.87	0.92	$0.92 \pm 1.8e - 2$	0.892 ± 0.063	0.975 ± 0.015

Table 2: AUC scores for higher-order n -way interaction effects

	3-Way Interactions		4-Way Interactions		5-Way Interactions	
	NID [41]	T-NID	NID [41]	T-NID	NID [41]	T-NID
Average	0.08 ± 0.013	0.76 ± 0.07	0.75 ± 0.13	0.78 ± 0.11	0.92 ± 0.06	0.97 ± 0.05

sign (red triangle) with labels “go” or “stop.” The COCO task suffered from model overfitting and lower test accuracy due to the limited pairwise data, but we still observed sensible explanations. Figure 2a) shows such interactions assigned the highest interaction salience by Taylor-CAM.

In the Yield-or-Go task, Taylor-CAM revealed two prediction strategies. The first is expected: the model interacts the yield sign (red triangle) with a car (rectangle), as seen in Figure 2b), then predicts “stop” accordingly. In the second, the model interacts one car with all of the other cars. One would expect it to relate the car and the yield sign, but the model discovered that the problem can be solved by checking if (1) a car is present, and (2) a red car is not present. Since each object has a different color, (2) implies that a yield sign is present and thus to “stop.” Demystifying such reasoning strategies is a unique benefit of Taylor-CAM.

However, when the correct label is “go,” *i.e.*, a car and yield sign are not present together, Taylor-CAM finds that the model rarely interacts anything, but rather either all interaction saliences are zero or objects interact with themselves (immediately adjacent regions) (Figure 2c)). This self-interacting is an intuitive and convenient interpretation that Taylor-CAM provides in the lack of salient interactions.

5.3. Relational Reasoning

Sort-Of-CLEVR is a toy dataset for relational reasoning proposed by [31]. It is a less-computationally expensive 2D form of the CLEVR VQA dataset [18] with a focus on relational questions. In our setup, these questions include distance and compare-&-count tasks. To test Taylor-CAM’s capacity to explain a neural network’s relational reasoning, we train a Relation Network [31] on Sort-Of-CLEVR and

visualize its top interactions in Figure 3. Relation Networks are simple modules augmented to CNNs that enable relational reasoning between image regions.

In Figure 3, interacting regions are indicated by two bounding boxes, and the top 4 interactions discovered by Taylor-CAM are shown per image. The input is an image of objects and a question about a particular *object of interest* and its relation to another object, and the output is the answer to that question. Since these questions are relational in nature, this problem requires relational reasoning, which we hope Taylor-CAM can be suited to explain. We invite the reader to use the discovered interactions in Figure 3 (as visualized by the bounding boxes) to try to deduce the objects of interest and questions for themselves before looking at the captions. For example, if the top 4 interactions each consist of objects that are close to each other and if each interaction includes the pink square, one might guess that the question is “Which shape is closest to the pink square?”

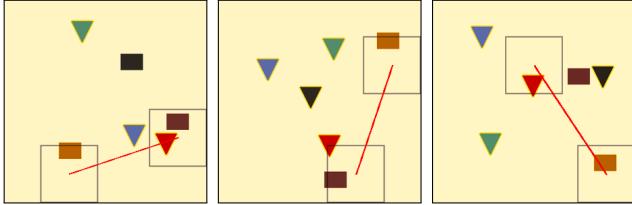
The 6 objects are “blue”, “purple”, “pink”, “yellow”, “orange”, and “green” and the 3 questions are (1) “Which shape is closest to the object of interest?”, (2) “Which shape is furthest from the object of interest?”, and (3) “How many objects have the same shape as the object of interest?”

While decisions are frequently relational [2], Grad-CAM is only designed to explain the importance of individual objects in isolation. We observed that Taylor-CAM affords much clearer explanations when decisions are relational.

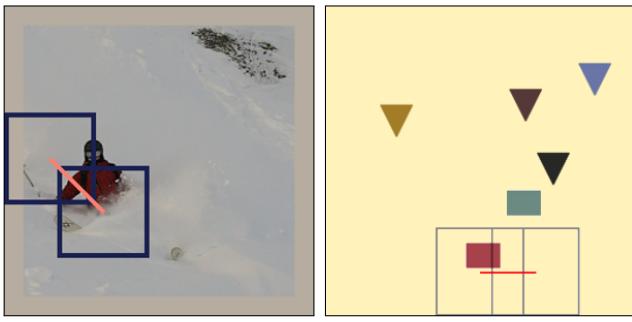
Quantitative Performance To assess quantitatively, 20 images per question that were classified correctly by the model were randomly selected and annotated with their question’s object of interest and answer-relevant objects. For example, for the question, “What is the shape of the ob-



a) Objects “person” and “car” are interacted to produce the output classification of whether both are present in the image in tandem.



b) Taylor-CAM interacts the yield sign (red triangle) with any present car (rectangle).



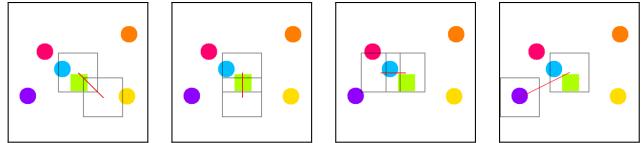
c) When no interactions present, Taylor-CAM’s interactions intuitively are 0 or occur primarily between adjacent regions as above.

Figure 2: Top-1 bounding boxes generated by Taylor-CAM representing simple interactions in multi-object detection.

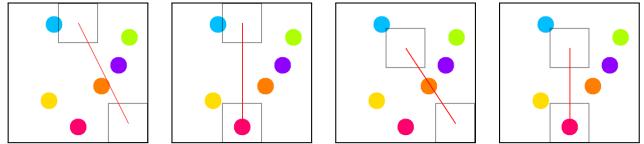
Table 3: Quantitative analysis on Sort-Of-CLEVR (%)

	Taylor-CAM	Grad-CAM* [32]	GLIDER [40]
Ques 1	90%	35%	60%
Ques 2	55%	50%	35%
Ques 3	60%	40%	45%

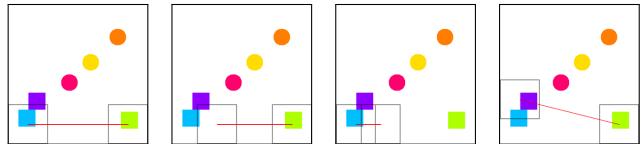
ject closest to the green square?” the green square and the object that is closest to it are annotated. If Taylor-CAM’s top-1 interaction (a pair of bounding boxes) intersects with the annotated pair, then it is counted as accurate for that image. Same with GLIDER. If Grad-CAM’s top-2 saliences include the annotated pair, then it is counted as accurate for that image. Since Grad-CAM does not provide relational interpretations, we refer to this relational interpretation of Grad-CAM’s saliences as Grad-CAM*. The bounding boxes in Figure 4 exemplify what a single salience looks like for Taylor-CAM and Grad-CAM respectively. Results of the quantitative analysis are reported in Table 3.



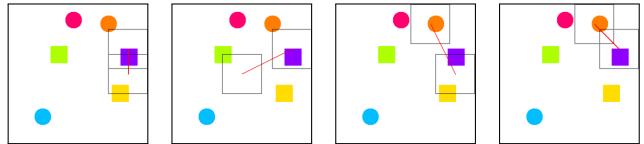
a) Q: “Which shape is closest to the green square?”



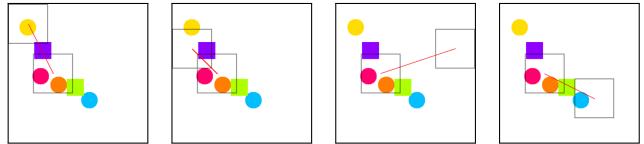
b) Q: “Which shape is furthest from the blue circle?”



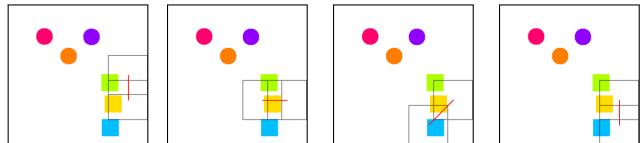
c) Q: “How many objects have shape of green object?”



d) Q: “Which shape is closest to the purple square?”



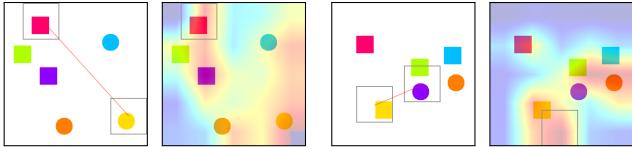
e) Q: “Which shape is furthest from the pink circle?”



f) Q: “How many objects have shape of yellow object?”

Figure 3: Shown are the top 4 interactions identified from a Relation Network’s predictions on 6 visual question-answering samples. The bounding boxes proposed by Taylor-CAM may be interpreted as indicating a relation. We recommend testing yourself to see if you can guess (1) the object of interest and (2) the question being asked (*closest*, *furthest*, or *same shape*), without looking at the caption.

Qualitative Performance To measure Taylor-CAM’s qualitative explainability, we selected a random batch of 15 samples and their ordered interaction saliences, and conducted a small user study ($n = 10$), asking each individual to guess (1) the object of interest and (2) the question being asked, from just looking at the top-4 ranked interaction vi-



a) Q: “Which shape is furthest from the pink square?” b) Q: “Which shape is closest to the yellow square?”

Figure 4: The bounding boxes show what a top-1 salience looks like for Taylor-CAM (on the left) and Grad-CAM (on the right) respectively. Taylor-CAM offers interpretable relational explanations from a single top-1 salience, whereas Grad-CAM depends on all saliences to produce a non-relational heatmap.

Table 4: User study **object of interest** accuracy (%)

	Grad-CAM [32]	GLIDER [40]	Taylor-CAM
Green	13.3%	33.3%	40%
Pink	30%	10%	46.7%
Blue	10%	22.2%	40%
Purple	N/A	15%	10%
Orange	3.3%	10%	15%
Yellow	25%	16.7%	33.3%

Table 5: User study **question** accuracy (%)

	Grad-CAM [32]	GLIDER [40]	Taylor-CAM
Ques 1	44%	38.9%	76%
Ques 2	14%	38.9%	55%
Ques 3	30%	23.8%	48.3%

suals. Taylor-CAM achieves strong explainability with better guess-accuracy than Grad-CAM and the recent GLIDER [40]. With Taylor-CAM, participants were able to reverse engineer questions in relational VQA from just looking at the visualized interactions. We report a wide range of explainability across different colors and questions in Tables 4 and 5. Due to random sampling, none of the 15 sampled images for Grad-CAM included a purple object of interest, so it is marked “N/A” in Table 4.

While some Grad-CAM colors strongly outperform random guessing (pink and yellow), on average, people struggled guessing the object of interest with Grad-CAM. This is because Grad-CAM only explains which individual objects contribute to the output, which in relational VQA, is all of them with an equal importance assigned to the object of interest and any objects that are included in the question-answer, such as the furthest or nearest object. This results in uninterpretable and sometimes misleading visualizations, making it very hard to guess an object of interest from the visual only. Without knowing the object of interest, it is consequently much harder to guess the question asked.

Grad-CAM, GLIDER, and Taylor-CAM all did relatively well on question 1. Closeness is easier to interpret with all three explanatory tools, since it is usually more visually apparent. However, we found question 2 (furthest distance) to be harder to interpret for Grad-CAM, perhaps because it is unclear what the object of interest is, with multiple “far away” objects of different relative proximity being ranked highly. For example, two objects that are far away from the object of interest might be close to each other, creating the false impression that the question is asking about closeness. Thus, without confidence regarding the object of interest and the interacting parts, we found ranked importances alone to be unintuitive and even misleading.

5.4. Biomedical Application

We also applied T-NID to determine interactions in the PPMI study dataset (www.ppmi-info.org). Our analysis suggests that various measures previously thought to be unrelated should be considered together when predicting faster cognitive progression in Parkinson’s disease. Please see Appendix for details in this domain.

6. Architecture Configurations

T-NID For T-NID, we trained a GELU-activated multi-layer perceptron with hidden layer sizes 140, 100, 60, and 20 for 200 epochs with a learning rate of 0.003 using early stopping [4] with a patience of 10. Results were averaged across 10 trials. Input data was normalized by standard deviation. T-NID hyperparameters were set as $\ell = 5$, $o = 2$, $k = 10$.

Taylor-CAM For our COCO [21] task, we used Pytorch’s ResNet-50 [11] pretrained on ImageNet [7], except we replaced the global average pooling layer with an additional convolutional layer composed of 1024 out-channels, size 2 kernel, 2 stride, and 2 padding, followed by 3 hidden linear layers of size 512, 256, 64, because global average pooling yields no higher-order derivatives. For our Relation Network, we used an open source reference implementation, which can be found here: <https://github.com/kimhc6028/relational-networks>, since [31] did not release code to the public. We trained for 50 epochs.

7. Conclusion

With T-NID and Taylor-CAM, we have shown that input cross derivatives, combined with a few simple heuristics and intuitions, are a powerful tool for explaining interactions in deep learning. T-NID, using GELU activations, representative samples, and interaction subsampling, successfully ranks statistical interactions, outperforming NID. Meanwhile, Taylor-CAM generalizes Grad-CAM to the higher order and effectively explains interactions in object detection and relational reasoning, affording a user cohort

the insight to guess questions in VQA from only seeing the top discovered visual interactions. Future work may explore localizing multi-modal interactions such as in audio-visual tasks, an agent’s interactions in RL and robotics, and interactions between word embeddings in NLP. By making our code publicly available, we hope that these simple explanatory tools can be used and built upon to better explain the complex interoperating factors underlying neural network reasoning and the world.

8. Acknowledgements

Research reported in this publication was supported by the National Institute Of Neurological Disorders And Stroke of the National Institutes of Health under Award Number P50NS108676. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

References

- [1] Chunrong Ai and Edward C Norton. Interaction terms in logit and probit models. *Economics letters*, 80(1):123–129, 2003. 3
- [2] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018. 2, 6
- [3] Jesse Bettencourt, Matthew J. Johnson, and David Duvenaud. Taylor-mode automatic differentiation for higher-order derivatives in JAX. In *Advances in neural information processing systems, Workshop Program Transformations*, 2019. 5
- [4] Rich Caruana, Steve Lawrence, and C Lee Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Advances in neural information processing systems*, pages 402–408, 2001. 8
- [5] A. Chattopadhyay, A. Sarkar, P. Howlader, and V. N. Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 839–847, 2018. 2, 4, 5
- [6] Tianyu Cui, Pekka Marttinen, and Samuel Kaski. Recovering pairwise interactions using neural networks. In *Advances in neural information processing systems, Bayesian Deep Learning workshop*, 2019. 2, 3, 5
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 8
- [8] Dennis W Dickson. Parkinson’s disease and parkinsonism: neuropathology. *Cold Spring Harbor perspectives in medicine*, 2(8):a009258, 2012. 14
- [9] Oliver Eberle, Jochen Büttner, Florian Kräutli, Klaus-Robert Müller, Matteo Valleriani, and Grégoire Montavon. Building and interpreting deep similarity models. *arXiv preprint arXiv:2003.05431*, 2020. 2
- [10] Jerome H Friedman, Bogdan E Popescu, et al. Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2(3):916–954, 2008. 11
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 8
- [12] Yotam Hechtlinger. Interpretation of prediction models using the input gradient. *ArXiv*, abs/1611.07634, 2016. 2
- [13] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. *CoRR*, abs/1606.08415, 2016. 3
- [14] Elke Heremans, Evelien Nackaerts, Sanne Broeder, Griet Vervoort, Stephan P Swinnen, and Alice Nieuwboer. Handwriting impairments in people with parkinson’s disease and freezing of gait. *Neurorehabilitation and neural repair*, 30(10):911–919, 2016. 18
- [15] Giles Hooker. Discovering additive structure in black box functions. In *Proceedings of the 2004 ACM SIGKDD international conference on Knowledge discovery and data mining - KDD ’04*, page 575. ACM Press. 3, 5
- [16] Sungjae Hwang, Peter Agada, Stephen Grill, Tim Kiemel, and John J Jeka. A central processing sensory deficit with parkinson’s disease. *Experimental brain research*, 234(8):2369–2379, 2016. 18
- [17] Joseph D Janizek, Pascal Sturmels, and Su-In Lee. Explaining explanations: Axiomatic feature interactions for deep networks. *arXiv preprint arXiv:2002.04138*, 2020. 2
- [18] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2901–2910, 2017. 5, 6
- [19] Natalia Jozwiak, Ronald B Postuma, Jacques Montplaisir, Véronique Latreille, Michel Panisset, Sylvain Chouinard, Pierre-Alexandre Bourgouin, and Jean-François Gagnon. Rem sleep behavior disorder and cognitive impairment in parkinson’s disease. *Sleep*, 40(8), 2017. 14
- [20] VE Kelly, CO Johnson, EL McGough, A Shumway-Cook, FB Horak, KA Chung, AJ Espay, FJ Revilla, J Devoto, C Wood-Siverio, et al. Association of cognitive domains with postural instability/gait disturbance in parkinson’s disease. *Parkinsonism & related disorders*, 21(7):692–697, 2015. 14
- [21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 5, 8
- [22] Yin Lou, Rich Caruana, Johannes Gehrke, and Giles Hooker. Accurate intelligible models with pairwise interactions. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD ’13*, page 623. ACM Press. 3, 5

- [23] Gr  oire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert M  ller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017. [2](#)
- [24] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *Advances in neural information processing systems*, 2017. [5](#)
- [25] Werner Poewe. Dysautonomia and cognitive dysfunction in parkinson’s disease. *Movement disorders: official journal of the Movement Disorder Society*, 22(S17):S374–S378, 2007. [14](#)
- [26] Mirthe M Ponsen, Andreas Daffertshofer, Erik Ch Wolters, Peter J Beek, and Henk W Berendse. Impairment of complex upper limb motor function in de novo parkinson’s disease. *Parkinsonism & Related Disorders*, 14(3):199–204, 2008. [18](#)
- [27] AH Rajput, A Voll, ML Rajput, CA Robinson, and A Rajput. Course in parkinson disease subtypes: a 39-year clinicopathologic study. *Neurology*, 73(3):206–212, 2009. [14](#)
- [28] Amir Rosenfeld, Richard S. Zemel, and John K. Tsotsos. The elephant in the room. *CoRR*, abs/1808.03305, 2018. [3](#)
- [29] Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. Right for the right reasons: training differentiable models by constraining their explanations. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence, IJCAI’17*, pages 2662–2670. AAAI Press. [2](#)
- [30] Adam Santoro, Ryan Faulkner, David Raposo, Jack Rae, Mike Chrzanowski, Theophane Weber, Daan Wierstra, Oriol Vinyals, Razvan Pascanu, and Timothy Lillicrap. Relational recurrent neural networks. In *Advances in neural information processing systems*, pages 7299–7310, 2018. [2, 11](#)
- [31] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning. In *Advances in neural information processing systems*, pages 4967–4976, 2017. [2, 6, 8, 11](#)
- [32] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 618–626, 2017. [2, 4, 7, 8](#)
- [33] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013. [2](#)
- [34] Krishna Kumar Singh, Dhruv Mahajan, Kristen Grauman, Yong Jae Lee, Matt Feiszli, and Deepti Ghadiyaram. Don’t judge an object by its context: Learning to overcome contextual bias. *CoRR*, abs/2001.03152, 2020. [3](#)
- [35] Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. Autoint: Automatic feature interaction learning via self-attentive neural networks. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, pages 1161–1170, 2019. [2](#)
- [36] Daria Sorokina, Rich Caruana, Mirek Riedewald, and Daniel Fink. Detecting statistical interactions with additive groves of trees. In *Proceedings of the 25th international conference on Machine learning - ICML ’08*, pages 1000–1007. ACM Press. [3, 5](#)
- [37] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML’17, page 3319–3328. JMLR.org, 2017. [2](#)
- [38] Robert Tibshirani. Regression shrinkage and selection via the lasso: a retrospective. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(3):273–282, 2011. [5](#)
- [39] AI Troster, AM Paolo, KE Lyons, SL Glatt, JP Hubble, and WC Koller. The influence of depression on cognition in parkinson’s disease: a pattern of impairment distinguishable from alzheimer’s disease. *Neurology*, 45(4):672–676, 1995. [14](#)
- [40] Michael Tsang, Dehua Cheng, Hanpeng Liu, Xue Feng, Eric Zhou, and Yan Liu. Feature interaction interpretability: A case for explaining ad-recommendation systems via neural interaction detection. In *International Conference on Learning Representations*, 2020. [2, 5, 6, 7, 8](#)
- [41] Michael Tsang, Dehua Cheng, and Yan Liu. Detecting statistical interactions from neural network weights. In *International Conference on Learning Representations*, 2018. [2, 3, 5, 6, 12, 13](#)
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. [2](#)
- [43] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. Adversarial removal of gender from deep image representations. *CoRR*, abs/1811.08489, 2018. [3](#)
- [44] Akira Wiberg, Michael Ng, Yasser Al Omran, Fidel Alfaro-Almagro, Paul McCarthy, Jonathan Marchini, David L Bennett, Stephen Smith, Gwenaelle Douaud, and Dominic Furniss. Handedness, language areas and neuropsychiatric diseases: insights from brain imaging and genetics. *Brain*, 142(10):2938–2947, 2019. [14](#)
- [45] T.H. Wonnacott and R.J. Wonnacott. *Introductory statistics*. Wiley series in probability and mathematical statistics. Wiley, 1977. [2, 5](#)
- [46] Vinicius Zambaldi, David Raposo, Adam Santoro, Victor Bapst, Yujia Li, Igor Babuschkin, Karl Tuyls, David Reichert, Timothy Lillicrap, Edward Lockhart, Murray Shanahan, Victoria Langston, Razvan Pascanu, Matthew Botvinick, Oriol Vinyals, and Peter Battaglia. Deep reinforcement learning with relational inductive biases. In *International Conference on Learning Representations*, 2019. [2, 11](#)
- [47] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014. [2](#)
- [48] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016. [2, 4](#)

Appendix

A Note On Terminology	11
B Representative Samples & Aggregations	11
C T-NID Algorithm	12
D Test Suite Of Synthetic Functions	12
E Additional Architectures For N-Way Interactions	12
F COCO Multi-Object Detection	12
G Grad-CAM On Relational Reasoning	12
H Interactional Relation Network (IRN)	12
I Taylor-CAM Pipeline	13
J Biomedical Analysis	13

A. Note On Terminology

In colloquial terms, two things are said to interact when they depend on each other in some way. Similar to [10], this can be formalized as follows:

Definition A.1. Entity Interaction Given an entity e_1 with attributes (a_1, \dots, a_n) , an interaction exists with another entity e_2 with attributes (b_1, \dots, b_n) if some a_i depends on some b_j or some b_j depends on some a_i .

Now we will define mathematical relation.

Definition A.2. Relation Given sets A and B , the binary relation from A to B is a subset of the Cartesian product $A \times B$.

We would like to unify our colloquial understanding of interaction in Definition A.1, our mathematical definition of relation in Definition A.2, and our definition for statistical interaction effects in Definition 1 of the main paper.

To connect this to Definition 1, we will reframe features as entities with the following theorem:

Theorem 1. *Given a function $F(\mathbf{x})$ and feature x_i , let entity e_i consist of attributes $(x_i, F/\partial x_i)$. An interaction exists between e_1 and e_2 if there is a nonzero interaction effect between x_1 and x_2 .*

Proof. If there is a nonzero interaction effect between x_1 and x_2 , then $\partial^2 F(\mathbf{x})/(\partial x_1 \partial x_2) \neq 0$ for some input \mathbf{x} . Then $F/\partial x_1$ depends on x_2 and consequently, there exists an interaction between entities e_1 and e_2 . \square

We have shown that our statistical interaction implies an interaction according to our colloquial understanding. An interaction exists between e_1 and e_2 if (but not only if, since the change need not be local) $F(\mathbf{x})/(\partial x_1 \partial x_2) \neq 0$, meaning $F/\partial x_1$ depends on x_2 . This is considered a binary relation between the two attributes, as all functions are relations, though not all relations are functions. Formally: given a function $F(\mathbf{x})$, a feature x_i , and entity e_i consisting of attributes $(x_i, F/\partial x_i)$, if there is a nonzero interaction effect between x_1 and x_2 , then a relation exists between the attributes of the two entities.

We have shown that, under this framing, an interaction effect is a relation, and if the interaction effect is nonzero, there must be a dependency/interaction between those entities. Since feature vectors in CNNs could be treated as entities [31, 46, 30], and if one interprets their gradients on the output to be implicit attributes, computing interaction effects between CNN feature vectors is equivalent to identifying the colloquial interactions and relations described in this formulation.

This is trivially generalized to interactions/relations of higher orders.

To summarize, a mathematical relation is implied by a colloquial interaction is implied by a statistical interaction, and this hierarchy can be formalized by regarding a feature x_i as an entity whose attributes include its gradients with respect to the function of interest. Thus, we offer a simple, formal connection between our statistical interaction effects definition and mathematical relations, as well as an integration of both into the colloquial understanding of “interaction” as merely a dependency between two “things.”

B. Representative Samples & Aggregations

Aggregation Of Representative Samples	AUC Score
Mean Of Mean-Min-Mode-Rand	0.61825
Mean Of Med-Min-Mode-Rand	0.61825
Mean Of Mean-Med-Min-Mode-Rand	0.61775
Mean Of Mean-Min-Max-Mode-Rand	0.6155
Mean Of Med-Min-Max-Mode-Rand	0.6155
Med Of Mean-Min-Mode-Rand	0.61525
Med Of Med-Min-Mode-Rand	0.61525
Mean Of Mean-Med-Min-Max-Mode-Rand	0.61525
Mean Of Mean-Min-Rand	0.614
Mean Of Med-Min-Rand	0.614

Table 6: Top average (across all orders) AUC scores for different aggregations of representative samples

Table 6 displays the top 10 aggregations and representative samples discovered via our power sweep.

C. T-NID Algorithm

Our complete T-NID is described in Algorithm 1. Note that each derivation of interaction effect using Definition 1 of the main paper for an interaction $I = \hat{I} \cup j$ of size ℓ where $|\hat{I}| = \hat{\ell} - 1$ for sample x can be derived as a single-order partial derivative $\partial IE_{\hat{I}} / \partial x_j$ and does not need to be recomputed from the ground up.

Algorithm 1 T-NID algorithm in pseudocode

Inputs ℓ -times differentiable trained neural network F , dataset \mathbf{X} with i th sample features $\mathbf{X}_{i1}, \dots, \mathbf{X}_{in}$, order ℓ , orders without subsampling o , subsampling size k .
Outputs Interaction effects IE_I for top estimated interactions $I \subseteq \{1, \dots, n\}$, where $|I| \leq \ell$.

Get representative samples:

For j th aggregation \in mean, minimum, mode, random
 $c = \operatorname{argmin}_i \| \mathbf{X}_i - \text{aggregation}(\mathbf{X}, \text{axis} = 0) \|$
 $r_j = \mathbf{X}_c$

For each representative sample:

For $r_j \in r$
Compute all non-redundant partial derivatives up to order o :
For $I \subseteq \{1, \dots, n\}$, where $|I| \leq o$
 $I = \text{sort}(I)$
If $IE_I^{(j)}$ uninitiated
Initiate $IE_I^{(j)}$ according to Definition 1 of the main paper

Compute remaining partial derivatives up to order ℓ by subsampling top k from previous orders:

For $\hat{\ell} \in o + 1, \dots, \ell$
For $\hat{I} \in \text{top } k \operatorname{argmax}$ of $IE_I^{(j)}$, where $|I| = \hat{\ell} - 1$
For $I \subseteq \{1, \dots, n\}$, where $|I| = \ell$ and $\hat{I} \subset I$
If $IE_I^{(j)}$ uninitiated
Initiate $IE_I^{(j)}$ according to Definition 1

of the main paper

Take the mean interaction effects across representative samples:

For $I \subseteq \{1, \dots, n\}$ if $IE_I^{(j)}$ initiated for some j
 $IE_I = \text{mean}(IE_I^{(j)})$ for all j where $IE_I^{(j)}$ initiated
Return IE

D. Test Suite Of Synthetic Functions

The test-suite of synthetic functions used to evaluate T-NID may be found in Table 7, courtesy of [41].

E. Additional Architectures For N -Way Interactions

Table 8 shows results for T-NID + MLP-M (T-NID using a neural network equipped with a main effects network as well as trained with sparsity regularization) and NID + MLP-M, the architecture used in [41].

F. COCO Multi-Object Detection

The task is to identify whether a pair of objects are each present in tandem. If only one is present, then the class label is negative. We tested this on the objects “car” and “person” in the COCO annotated-image dataset. We configured the frequency of the labels such that an even amount of positive and negative samples were in the training set. We found the COCO task to be somewhat inconclusive, because of model overfitting and rather low test accuracy, but still observed reasonable explanations, as seen in Figure 6. Taylor-CAM often prioritizes the car-person interaction correctly.

G. Grad-CAM On Relational Reasoning

Grad-CAM is a first-order explanatory tool that ranks different image regions and produces a heatmap of saliences. As shown qualitatively in Figure 7, Grad-CAM’s heatmaps are much harder to interpret and to reverse engineer questions and objects from compared to the results obtained from Taylor-CAM, shown in Figure 3 of the main paper, as corroborated quantitatively with our human study.

H. Interactional Relation Network (IRN)

A standard RN pools a set of feature vectors $O = \{o_1, \dots, o_n\}$, their corresponding positional encodings $C = \{c_1, \dots, c_n\}$, and a question q as follows:

$$\text{RN}(O, C, q) = f_\phi \left(\sum_{i,j} g_\theta(o_i, o_j, c_i, c_j, q) \right), \quad (6)$$

where f and g are modeled by neural networks parameterized by ϕ and θ respectively.

We observed through Taylor-CAM that many of the top interactions in the RN’s reasoning were between individual regions and themselves, even when we zeroed out diagonals, such as in Figure 2c) of the main paper. We found that we could mitigate this by making a simple modification to the RN architecture which we found to yield better test accuracy:

$$\begin{aligned} \text{IRN}(O, C, q) = \\ f_\phi \left(\sum_{i,j} g_\theta(h_\psi(o_i, c_i, q), h_\psi(o_j, c_j, q), c_i, c_j, q) \right), \end{aligned} \quad (7)$$

$F_1(\mathbf{x})$	$\pi^{x_1 x_2} \sqrt{2x_3} - \sin^{-1}(x_4) + \log(x_3 + x_5) - \frac{x_9}{x_{10}} \sqrt{\frac{x_7}{x_8}} - x_2 x_7$
$F_2(\mathbf{x})$	$\pi^{x_1 x_2} \sqrt{2 x_3 } - \sin^{-1}(0.5x_4) + \log(x_3 + x_5 + 1) + \frac{x_9}{1+ x_{10} } \sqrt{\frac{ x_7 }{1+ x_8 }} - x_2 x_7$
$F_3(\mathbf{x})$	$\exp x_1 - x_2 + x_2 x_3 - x_3^{2 x_4 } + \log(x_4^2 + x_5^2 + x_7^2 + x_8^2) + x_9 + \frac{1}{1+x_{10}^2}$
$F_4(\mathbf{x})$	$\exp x_1 - x_2 + x_2 x_3 - x_3^{2 x_4 } + (x_1 x_4)^2 + \log(x_4^2 + x_5^2 + x_7^2 + x_8^2) + x_9 + \frac{1}{1+x_{10}^2}$
$F_5(\mathbf{x})$	$\frac{1}{1+x_1^2+x_2^2+x_3^2} + \sqrt{ x_4 + x_5 + x_6 + x_7 + x_8 x_9 x_{10}}$
$F_6(\mathbf{x})$	$\exp(x_1 x_2 + 1) - \exp(x_3 + x_4 + 1) + \cos(x_5 + x_6 - x_8) + \sqrt{x_8^2 + x_9^2 + x_{10}^2}$
$F_7(\mathbf{x})$	$(\arctan(x_1) + \arctan(x_2))^2 + \max(x_3 x_4 + x_6, 0) - \frac{1}{1+(x_4 x_5 x_6 x_7 x_8)^2} + \left(\frac{ x_7 }{1+ x_9 }\right)^5 + \sum_{i=1}^{10} x_i$
$F_8(\mathbf{x})$	$x_1 x_2 + 2^{x_3+x_5+x_6} + 2^{x_3+x_4+x_5+x_7} + \sin(x_7 \sin(x_8 + x_9)) + \arccos(0.9x_{10})$
$F_9(\mathbf{x})$	$\tanh(x_1 x_2 + x_3 x_4) \sqrt{ x_5 } + \exp(x_5 + x_6) + \log((x_6 x_7 x_8)^2 + 1) + x_9 x_{10} + \frac{1}{1+ x_{10} }$
$F_{10}(\mathbf{x})$	$\sinh(x_1 + x_2) + \arccos(\tanh(x_3 + x_5 + x_7)) + \cos(x_4 + x_5) + \sec(x_7 x_9)$

Table 7: Synthetic test-suite functions

	T-NID 3-Way	NID 3-Way	T-NID 4-Way	NID 4-Way	T-NID 5-Way	NID 5-Way
$F_1(\mathbf{x})$	0.831 ± 0.064	0.122 ± 0.028	0.777 ± 0.389	0.555 ± 0.456	N/A	N/A
$F_2(\mathbf{x})$	0.629 ± 0.165	0.07 ± 0.011	0.032 ± 0.065	0.185 ± 0.37	N/A	N/A
$F_3(\mathbf{x})$	0.991 ± 0.013	0.095 ± 0.008	0.997 ± 0.006	1.0 ± 0.0	N/A	N/A
$F_4(\mathbf{x})$	0.993 ± 0.007	0.09 ± 0.028	0.96 ± 0.08	0.996 ± 0.009	N/A	N/A
$F_5(\mathbf{x})$	0.493 ± 0.009	0.035 ± 0.005	N/A	N/A	N/A	N/A
$F_6(\mathbf{x})$	0.103 ± 0.025	0.034 ± 0.005	N/A	N/A	N/A	N/A
$F_7(\mathbf{x})$	0.417 ± 0.264	0.156 ± 0.031	0.363 ± 0.322	0.711 ± 0.046	0.303 ± 0.367	0.536 ± 0.453
$F_8(\mathbf{x})$	1.0 ± 0.0	0.141 ± 0.008	1.0 ± 0.0	0.994 ± 0.008	N/A	N/A
$F_9(\mathbf{x})$	0.838 ± 0.146	0.113 ± 0.01	0.859 ± 0.068	0.618 ± 0.084	0.988 ± 0.024	0.549 ± 0.452
$F_{10}(\mathbf{x})$	1.0 ± 0.0	0.03 ± 0.002	N/A	N/A	N/A	N/A
average	0.73 ± 0.069	0.089 ± 0.014	0.713 ± 0.133	0.723 ± 0.139	0.646 ± 0.200	0.543 ± 0.453

Table 8: N-Way AUC scores for T-NID + MLP-M and NID + MLP-M, both using a main effects network and sparsity regularization, as described in [41]

where h is an MLP parameterized by ψ .

We refer to this as Interactional Relation Network (IRN) since it explicitly separates within its architecture the concerns of reasoning about interactions from reasoning about individual objects. While IRN does not relate to Taylor-CAM directly, it does highlight how visualizing a network’s relational reasoning can inspire potential ideas for improvement to a network’s architecture.

I. Taylor-CAM Pipeline

In Figure 5, we illustrate the full pipeline of Taylor-CAM. A model consisting of a CNN and Relation Network predicts answers based on images and questions. Taylor-CAM intercepts the model’s gradients, reverse engineers the question asked, and visualizes for a human observer.

Given three possible question categories (*closest*, *furthest*, and *same shape*), the user is able to interpret which the model is reasoning about by looking at Taylor-CAM’s proposed interactions, as shown quantitatively in Tables 4 and 5 of the main paper.

J. Biomedical Analysis

We applied these techniques to the Parkinson’s Progression Marker Initiative (PPMI) study (<http://www.ppmi-info.org/>) dataset, which follows persons living with early-stage Parkinson’s disease for up to approximately eight years collecting clinical and biological data from participants. Parkinson’s disease (PD) is a neurodegenerative progressive disease, characterized clinically by motor (e.g., tremor, rigidity) and non-motor (e.g., cognition and autonomic dysfunction) symptoms that vary over time within and between patients. Progression of motor and non-motor symptoms are likely not independent of each other. Instead, collateral damage may be inflicted multilaterally with non-motor and motor pathological features progressing interdependently. As an example, depressive symptoms in Parkinson’s disease are common and may perpetuate motor and cognitive deficits, which could impact function, and ultimately diminish quality of life. Therefore, it is necessary to take as comprehensive of an approach as possible in unraveling the clinical progression of Parkinson’s disease.

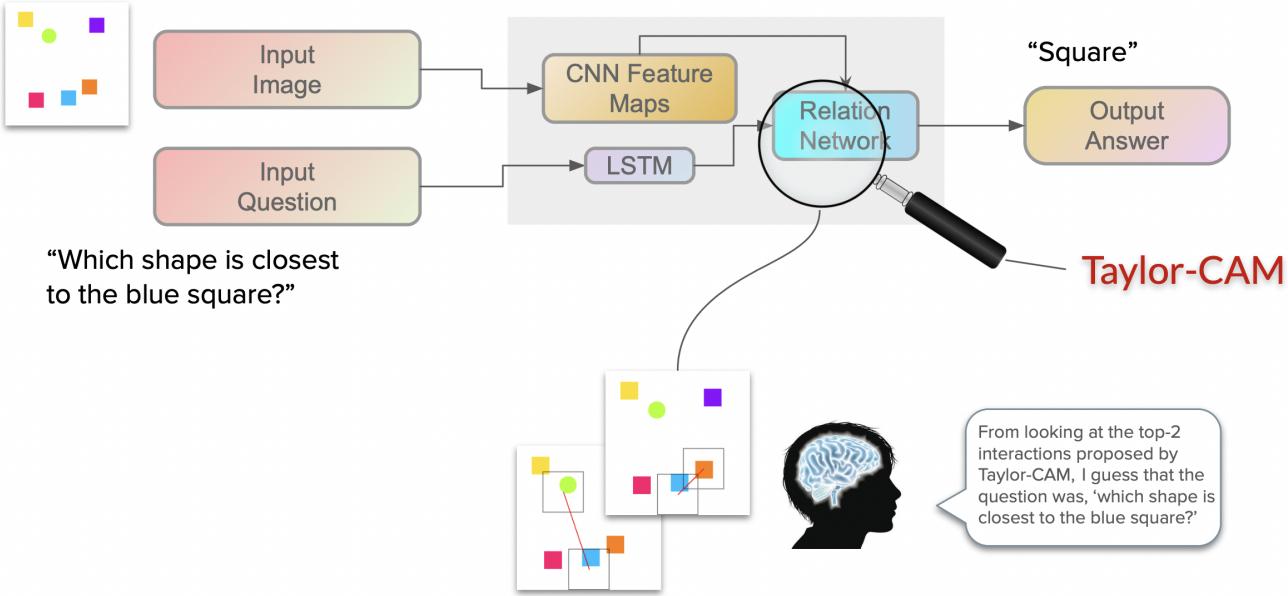


Figure 5: Full Taylor-CAM pipeline on the problem of relational visual question answering. A CNN + RN take in an image and question, and output an answer. Taylor-CAM is able to visualize the model’s reasoning from just its gradients such that a human can interpret what and how the model reasoned.

As PD progresses, cognitive impairment leading to dementia may affect up to 80% of patients, ultimately impairing one’s functional independence. Within the PPMI study, we tested 2-, 3-, 4- and 5-way interactions to understand multivariable features at baseline that distinguish patients with a more severe progression in decline of cognitive function (“fast progressors”) compared to those with a more benign course of cognitive changes, as measured by the Montreal Cognitive Assessment (MoCA) scale. Top 2-way interaction effects identified (Figure 8) among “fast progressors” included feature interactions between handedness (handed) and severity of rigidity in the neck (np3rign); presence of resting tremor at disease diagnosis (dxtremor) and severity of rigidity in the lower extremities (np3rigll); and, severity of tremor (np2trmr) and alternating trail making test from the MoCA scale (mcaalttm) – which ultimately is a measure of processing speed, mental flexibility, ability to sequence, and visuo-motor skills. Each of these features individually have some established associations with cognitive dysfunction or neuropsychological disorders; however, their interactions together have not been previously considered. For example, handedness, has been significantly associated with functional connectivity between language networks, as well as specific genetic loci implicated in the pathogenesis of neurologic disorders including Parkinson’s disease [44]. More severe rigidity symptoms in Parkinson’s disease are also associated with faster cognitive decline [27]. Our analysis, for the first time, suggests that measures of both handedness and rigidity severity together

are important to consider when predicting faster cognitive progression in Parkinson’s disease. As shown in Figure 8, we provide 3-, 4-, and 5-way interactions between features.

When broadening the interactions to 3-, 4-, and 5-way interactions between features that predict fast cognitive decline, we observe additional features with some consistency (Figure 9 provides 3-, 4-, and 5-way interactions between features). Broadly speaking, some of the most important interactions occurred between symptoms of autonomic dysfunction: urinary (np1urin, scau11, scau13) and constipation issues (np1cnst, scau5), problems tolerating cold/heat [scau20]; mood and sleep disturbances: depression (np1dprs), apathy (np1apat), and restless sleep (np1slpn, slplmbmv); postural instability and balance issues (np2walk, np3pstbl); overall severity of Parkinson’s disease (nhy); and, memory impairment: delayed recall (mcarec2, mcarec4). Each of these symptoms, singularly, have been thought to be associated with cognitive impairment [25, 19, 39, 20]. It is novel, yet biologically plausible to consider these symptoms interacting, as the neuropathology underlying Parkinson’s disease involves multiple areas of the brain and nervous system beyond the nigrostriatal dopamine pathway. For instance, Lewy Body pathology affects the limbic cortex and frontal neocortical areas, sympathetic ganglia and even the peripheral autonomic nervous system including the myenteric plexus [8].

We also performed our analysis to predict fast progression of ambulatory impairment which stems from worsening progression of motor symptoms and is a major source



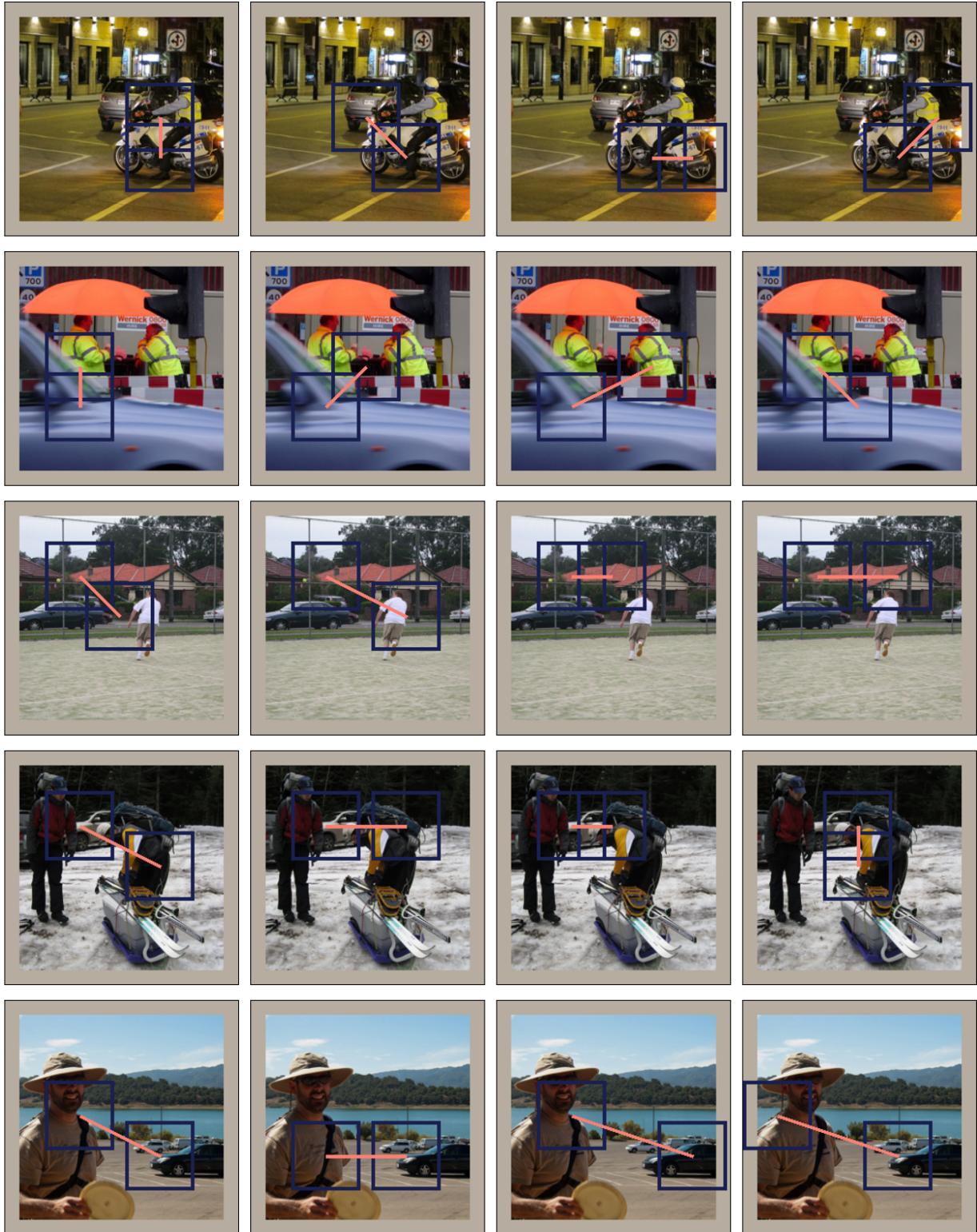
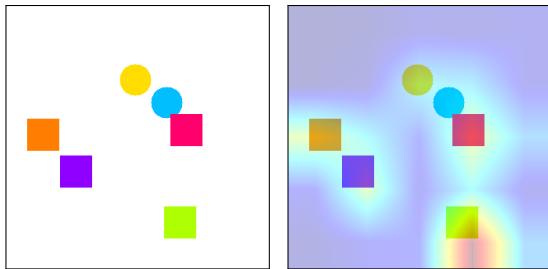
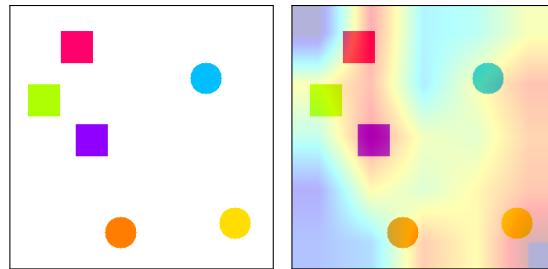


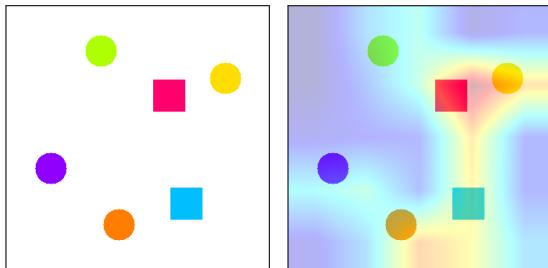
Figure 6: Taylor-CAM interacts “car” and “person” in the custom COCO multi-object detection task. Taylor-CAM’s top 4 discovered interactions per image are shown. As discussed in the main paper, in cases where the interaction is not present (as in 4th example), Taylor-CAM interacts immediately adjacent regions around whichever of the two objects is present.



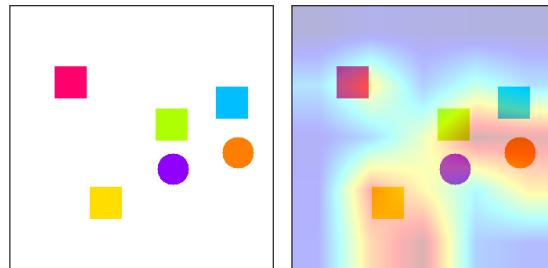
a) Q: "How many objects have shape of purple object?"



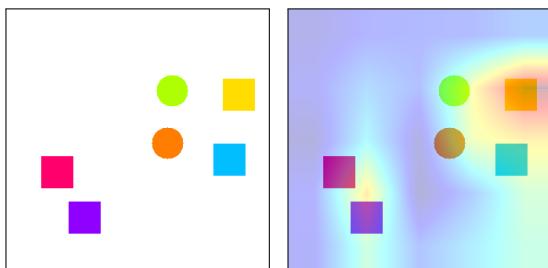
b) Q: "Which shape is furthest from the pink square?"



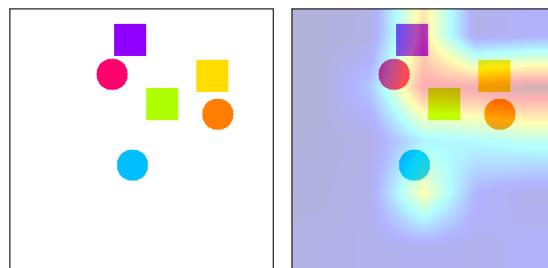
c) Q: "How many objects have shape of orange object?"



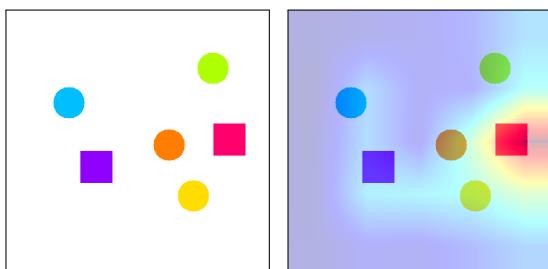
d) Q: "Which shape is closest to the yellow square?"



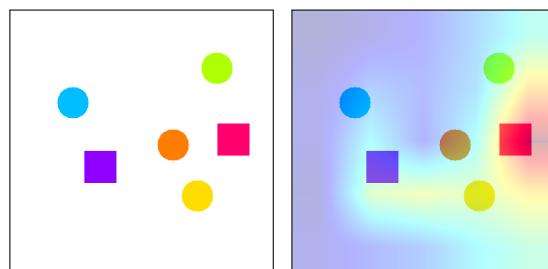
e) Q: "Which shape is closest to the purple square?"



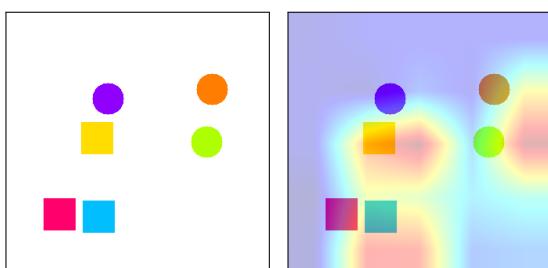
f) Q: "How many objects have shape of pink object?"



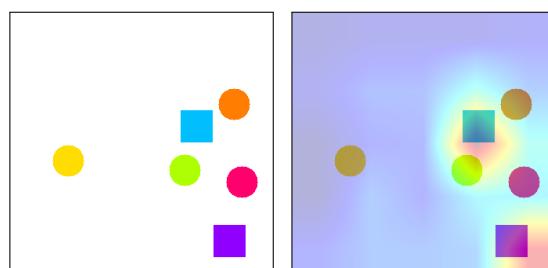
g) Q: "Which shape is furthest from the green circle?"



h) Q: "Which shape is closest to the blue circle?"

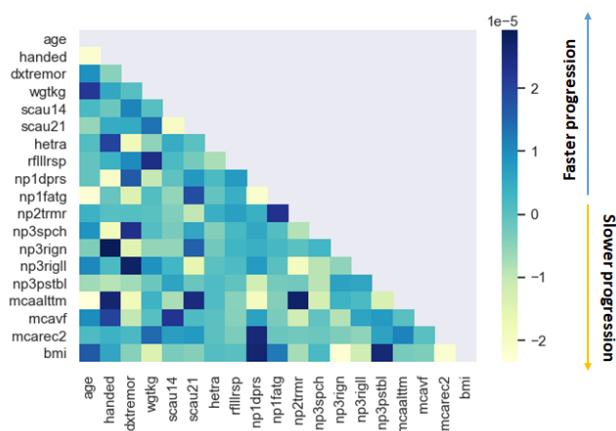


i) Q: "Which shape is furthest from the green circle?"



j) Q: "Which shape is furthest from the pink circle?"

Figure 7: Grad-CAM's salience heatmaps per image/question. *Left*, the raw image. *Right*, Grad-CAM's explanatory heatmap.

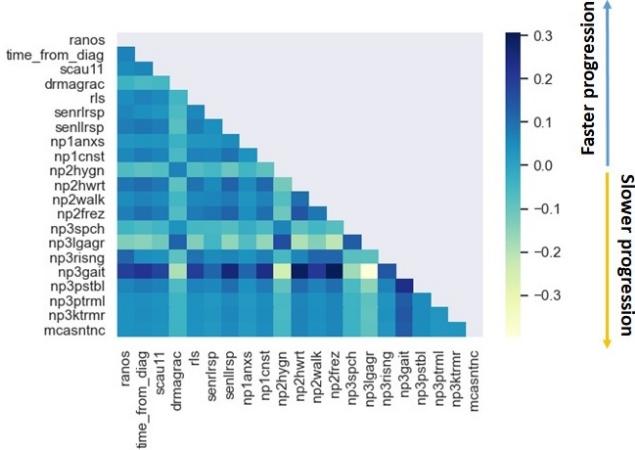


a) Top 2-way interactions for MoCA fast progression

N-Way Interaction	Strength
id_num, scau20, mcarec4	4.77E-06
id_num, drmagrac, mcarec4	4.69E-06
educyrs, np1apat, bmi	4.56E-06
np1dprs, np2walk, np3pstbl	4.27E-06
scau13, np1slpn, np1cnst, nhyl	6.00E-07
scau11, scau13, scau20, bmi	5.66E-07
scau11, scau13, np1slpn, nhyl	5.64E-07
scau13, scau20, np1urin, nhyl	5.43E-07
slplmbmv, np1dprs, np2walk, np3rigr, np3pstbl	1.23E-07
slplmbmv, np1dprs, np2walk, np3rign, np3pstbl	1.22E-07
scau5, np1dprs, np2walk, np3rigr, np3pstbl	1.19E-07
slplmbmv, np1dprs, np2walk, np3pstbl, mcarec2	1.18E-07

b) Top 3-way, 4-way, and 5-way interactions for MoCA fast progression

Figure 8: Interaction effects for classifying fast clinical progression of MoCA scores from baseline



a) Top 2-way interactions for uMCA fast progression

N-Way Interaction	Strength
scau1, np3lgagr, np3risng	1.97E+00
scau1, scau9, np3lgagr	1.41E+00
scau1, np2hwrt, np3lgagr	1.31E+00
scau1, senllrsp, np3lgagr	1.19E+00
time_from_diag, scau1, np2frez, np3lgagr	5.39E+00
scau1, np2walk, np2frez, np3lgagr	5.28E+00
ranos, scau1, np2frez, np3lgagr	5.25E+00
scau1, rls, np2frez, np3lgagr	5.21E+00
scau1, np3lgagr, np3risng, np3gait, np3tarl	1.93E+01
scau1, np2hygn, np3lgagr, np3risng, np3gait	1.80E+01
scau1, mslarsp, np3lgagr, np3risng, np3gait	1.68E+01
dxrigid, scau1, np3lgagr, np3risng, np3gait	1.47E+01

b) Top 3-way, 4-way, and 5-way interactions for uMCA fast progression

Figure 9: Interaction effects for classifying fast clinical progression of uMCA scores from baseline

of disability for patients with Parkinson’s disease. Severity of ambulatory impairment was measured by an ambulatory capacity score derived from sum of scores of the MDS-UPDRS items 2.13 (freezing), 2.12 (walking and balance), 3.10 (gait), 3.12 (postural stability), and 3.11 (freezing of gait). The top 2-way interaction among “fast progressors” of ambulatory capacity was between severity of freezing (np2frez) and gait (np3gait), which is unsurprising as both are components of the ambulatory capacity scale score. Interestingly, however, the next top 2-way interactions were between handwriting (np2hwrt) and gait (np3gait); and, sensory of legs (senllrsp) and gait (np3gait). Worsening of handwriting is often reported as an initial symptom of Parkinson’s disease and is reported to be more problematic

in people with Parkinson’s disease who experience freezing of gait [26, 14]. Periphery sensory defects in the lower limbs are also commonly noted in people with Parkinson’s disease, and could be a main contributor to balance control issues and postural instability [16]. As interactions were expanded to 3-, 4-, and 5-ways, other items that were consistently identified were difficulty in swallowing and chewing (scau1), and leg agility (np3lgagr). While these items are not often considered as obvious predictors of ambulatory capacity by themselves, their interactions with some more apparent features (e.g., gait, freezing, arising from chair) provide new insights on how symptoms in Parkinson’s disease patients contribute to disease progression.