# Towards Best Practice in Explaining Neural Network Decisions with LRP

Maximilian Kohlbrenner[1], Alexander Bauer[2], Shinichi Nakajima[2], Alexander Binder[3],
Wojciech Samek[1,*] and Sebastian Lapuschkin[1,*]

[1]Dept. of Video Coding and Analytics, Fraunhofer Heinrich Hertz Institute, Berlin, Germany
[2]Dept. of Electrical Engineering and Computer Science, Technische Universität Berlin, Berlin, Germany
[3]ISTD Pillar, Singapore University of Technology and Design, Singapore, Singapore
[*]{wojciech.samek|sebastian.lapuschkin}@hhi.fraunhofer.de

*Abstract*—**Within the last decade, neural network based predictors have demonstrated impressive — and at times superhuman — capabilities. This performance is often paid for with an intransparent prediction process and thus has sparked numerous contributions in the novel field of *explainable artificial intelligence (XAI)*. In this paper, we focus on a popular and widely used method of XAI, the *Layer-wise Relevance Propagation (LRP)*. Since its initial proposition LRP has evolved as a method, and a *best practice* for applying the method has tacitly emerged, based however on humanly observed evidence alone. In this paper we investigate — and for the first time *quantify* — the effect of this current best practice on feedforward neural networks in a visual object detection setting. The results verify that the layer-dependent approach to LRP applied in recent literature better represents the model's reasoning, and at the same time increases the object localization and class discriminativity of LRP.**

*Index Terms*—**layer-wise relevance propagation, explainable artificial intelligence, neural networks, visual object recognition, quantitative evaluation**

## I. INTRODUCTION

In recent years, deep neural networks (*DNN*) have become the state of the art method in many different fields, but are mainly applied as black-box predictors. Since understanding the decisions of artificial intelligence systems is crucial in numerous scenarios and partially demanded by law[1], neural network interpretability has been established as an important and active research area. Consequently, many approaches to explaining neural network decisions have been proposed in recent years, *e.g.* [3]–[6]. The *Layer-wise Relevance Propagation (LRP)* [7] framework has proven successful at providing a meaningful intuition and measurable quantities describing a network's feature processing and decision making [8]–[10]. LRP attributes *relevance scores* $R_i$ to the model inputs or intermediate neurons $i$ by decomposing a model output of interest. The method follows the principles of *relevance conservation* and *proportional decomposition*. Therefore, attribu-

tions computed with LRP maintain a strong connection to the predictor output. While early applications of LRP administer a single decomposition rule uniformly to all layers of a model [7], [11], [12], more recent work describes a trend towards assigning specific decomposition rules purposely to layers wrt. function and position within the network [10], [13]–[16]. This trend has tacitly emerged and formulates a *best practice* for applying LRP. Under qualitative evaluation, the attribution maps resulting from this current approach seem to be more robust against the well-known effects of shattered gradients [12], [13], [17] and demonstrate an increased discriminativity between different target classes [13], [14] compared to the uniform application of a single rule.

However, recent literature applying LRP-rules in a layer-dependent manner do not justify the beneficial effects of this novel variant *quantitatively*, but only based on human observation. In this paper, we design and conduct a series of experiments in order to verify whether a layer-specific application of different decomposition rules actually constitutes an improvement above earlier descriptions and applications of LRP [11], [18]. That is, we measure and compare capabilities of various methods from explainable AI — with a focus on earlier and more recent approaches to LRP — to *precisely* localize the ground-truth objects in images via attribution of relevance scores. Our experiments are conducted on popular computer vision data sets with ground truth object localizations, the ImageNet [19] and PascalVOC [20] datasets, using different neural network models.

## II. FEEDFORWARD NEURAL NETWORKS AND LRP

Feedforward neural networks constitute a popular architecture type, ranging from simple multi-layer perceptrons and shallower convolutional architectures such as the LeNet-5 [21] to deeper and more complex Inception [22] and VGG-like architectures [23]. These types of neural network commonly use ReLU non-linearities and first pass information through a stack of convolution and pooling layers, followed by several fully connected layers. The good performance of feedforward architectures in numerous problem domains, and the availability as pre-trained models makes them a valuable standard architecture in neural network design.

[1]*e.g.* via the "right to explanation" proclaimed in the General Data Protection Regulation of the European Union [1], [2]
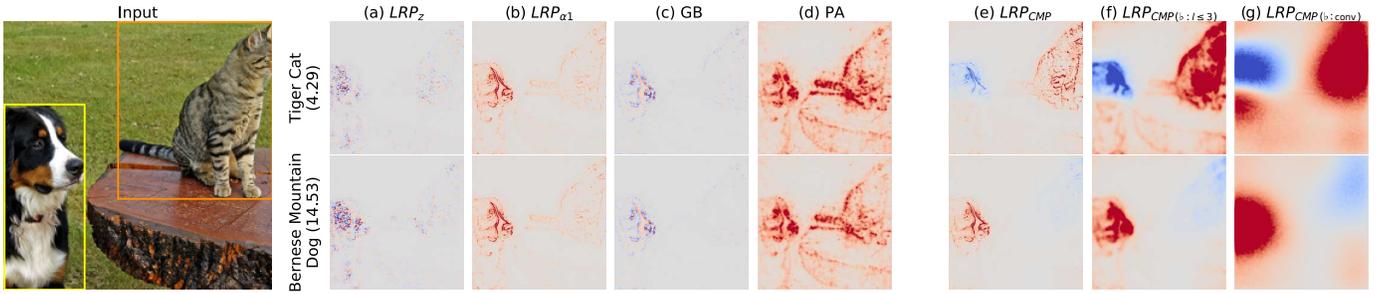
Fig. 1. Different attributions for the output classes "Tiger Cat" and "Bernese Mountain Dog" using the VGG-16 model. Network output strengths (logit) of the respective classes is given in parentheses. Network-widely applied rules in *(a)* - *(d)* ($LRP_z$, $LRP_{\alpha\beta}$, Guided Backprop and Pattern Attribution respectively), are not, or hardly class discriminative. An application of $LRP_z$ to every layer shows the effect of gradient shattering. Variants of $\boldsymbol{LRP}_{CMP}$ implementing a composite strategy of LRP rule application shown in *(e)* - *(g)* — here, from left to right, the $LRP_\flat$-rule is not applied at all, the three lowest convolution layers, and all convolution and pooling layers — are sensitive to class-specific information and highlight features on different levels of scale and conceptuality (*e.g.* highlighting the fur pattern activating "Tiger Cat" vs highlighting the general region showing a "Tiger Cat"). Attributions from $\boldsymbol{LRP}_{CMP}$ visualized in red/warm colors identify image regions contributing to the prediction of the target class, while regions marked in blue/cold hues provide contradictory evidence. Further examples can be found in the Appendix.

## A. Layer-wise Relevance Propagation

Consequently, feedforward networks have been subject to investigations in countless contributions towards neural network interpretability, including applications of LRP [7], [11], [18], which finds its mathematical foundation in *Deep Taylor Decomposition (DTD)* [24].

The most basic attribution rule of LRP (to which we will refer to as $LRP_z$) is defined as

$$R_i^{(l)} = \sum_j \frac{z_{ij}}{z_j} R_j^{(l+1)} \quad (1)$$

and performs a proportional decomposition of a given upper layer relevance value $R_j^{(l+1)}$ at some layer $(l+1)$ and neuron $j$ to obtain lower layer relevance scores $R_i^{(l)}$ for neurons $i$ at layer $(l)$, wrt. to the localized preactivations $z_{ij}$ and their respective aggregations $z_j$ at the layer output. Here, the localized preactivations $z_{ij}$ describe quantities propagated through the model during prediction time, *e.g.* $z_{ij} = x_i w_{ij}$ and $z_j = \sum_i z_{ij}$ within a neural network layer with learned weight parameters $w_{ij}$. Note that Eq. (1) is conservative between layers and in general maintains an equality $\sum_i R_i^{(l)} = f(x)$ at any layer $(l)$ of the model.

Further purposed LRP-rules beyond Eq. (1) are introduced in [7], which can be understood as advancements thereof:

So does the $LRP_\varepsilon$ decomposition rule [7] add a signed and small constant $\varepsilon$ to the denominator in order to prevent divisions by zero and to diminish the effect of recessive (*e.g.* weak and noisy) mappings $z_{ij}$ to the relevance decomposition.

$$R_i^{(l)} = \sum_j \frac{z_{ij}}{z_j + \varepsilon \cdot sign(z_j)} R_j^{(l+1)} \quad (2)$$

The $LRP_{\alpha\beta}$-rule [7] performs and then merges separate decompositions for the activatory ($z_{ij}^+$) and inhibitory ($z_{ij}^-$) parts of the forward pass

$$R_i^{(l)} = \sum_j \left( \alpha \frac{z_{ij}^+}{z_j^+} + \beta \frac{z_{ij}^-}{z_j^-} \right) R_j^{(l+1)} \quad (3)$$

where

$$z_{ij}^+ = \begin{cases} z_{ij} & ; z_{ij} > 0 \\ 0 & ; \text{else} \end{cases} \qquad z_{ij}^- = \begin{cases} 0 & ; \text{else} \\ z_{ij} & ; z_{ij} < 0 \end{cases} \quad (4)$$

Here, the non-negative $\alpha$ parameter permits a weighting of relevance distribution towards activations and inhibitions. The $\beta$ parameter is given implicitly s.t. $\alpha + \beta = 1$ in order to uphold conservativity of relevance between layers. The commonly used parameter $\alpha = 1$ can be derived from DTD and has been rediscovered in *ExcitationBackprop* [25].

Later work [14], [26] introduces $LRP_\flat$[2], a decomposition rule which spreads the relevance of a neuron uniformly across all its inputs. This rule assumes $z_{ij} = 1$ and $z_j = \sum_i 1$ in Eq. (1) *only* for backpropagating given relevance scores $R_j^{(l+1)}$ to lower layers $(l)$, and has seen application in the input layer(s) of neural networks. The $LRP_\flat$-rule provides invariance to the decomposition process wrt. to translations in the input domain and effectively propagates relevance scores of higher layer neurons — encoding "explanations" of more abstract concepts — towards the input via the neurons' *receptive fields*, without further transformation. Note that the $LRP_\flat$ decomposition rule is thus unsuitable for decomposing fully connected layers.

Earlier applications of LRP (*e.g.* [7], [11]) did use one single decomposition rule uniformly over the whole network, which often resulted in suboptimal "explanations" of model behavior [13]. So are network-wide applications of $LRP_z$ (in the following denoted as $\boldsymbol{LRP}_z$, in order to distinguish this specific *configuration* of LRP from the *rule* $LRP_z$) and network-wide applications of $LRP_\varepsilon$ (denoted as $\boldsymbol{LRP}_\varepsilon$) respectively identical and highly similar to *Gradient×Input (G×I)* in ReLU-activated DNNs [12]. $\boldsymbol{LRP}_z$ and $\boldsymbol{LRP}_\varepsilon$ demonstrate — albeit working well for shallower convolutional models [27], [28] such as the LeNet-5 [21] or simpler fully-connected networks [29] — the effect of gradient shattering as overly complex attributions for deeper models [12], [13]

---

[2]read: $\flat$ ="flat", as in the musical $\flat$.

(*cf.* Fig. 1*(a)*). A Network-wide application of $LRP_{\alpha\beta}$ (denoted as $\boldsymbol{LRP}_{\alpha\beta}$) demonstrates robustness against gradient shattering and produces visually pleasing attribution maps, however is lacking in class- or object discriminativity [13], [30]. By separately considering activatory and inhibitory mappings $z_{ij}$ during the decomposition process, $\boldsymbol{LRP}_{\alpha\beta}$ tends to attribute relevance to similar sets of input features activating sequences of neurons throughout the network, regardless of the output class chosen for relevance decomposition (*cf.* Fig. 1*(b)*). Further, $\boldsymbol{LRP}_{\alpha\beta}$ introduces the constraint of strictly positive layer activations [24], which is in general not guaranteed, especially at the (logit) output of a model. A dissatisfaction of this constraint may result in a sign inversion of all backpropagated relevance scores.

### B. A Current Best Practice for LRP

A recent trend among XAI researchers and practitioners employing LRP is the use of a *composite strategy* of rule applications for decomposing the prediction of a neural network [10], [13]–[16]. That is, different parts of the DNN are decomposed using purposed rules, which in combination are robust against gradient shattering while sustaining object discriminativity. Common among these works is the utilization of $LRP_{\varepsilon}$ with $\varepsilon \ll 1$ (or just $LRP_z$) to decompose fully connected layers close to the model output, followed by an application of $LRP_{\alpha\beta}$ to the underlying convolutional layers (usually with $\alpha \in \{1, 2\}$). Here, the separate decomposition of the positive and negative forward mappings complements the localized feature activation of convolutional filters activated by, and feeding into ReLUs. A final decomposition step within the convolution layers near the input uses the $LRP_{\flat}$-rule. Most commonly this rule (or alternatively the $DTD_{z^B}$-rule defined in context of Deep Taylor Decompositon [24]) is applied to the input layer only. In summary, we here describe this pattern of rule application as $\boldsymbol{LRP}_{CMP}$ (for *CoMP*osite). Fig. 1 provides a qualitative overview of the effect of $\boldsymbol{LRP}_{CMP}$ in contrast to other parameterizations and methods, which we will further discuss in Sec. IV-B. Note that the option to apply the $LRP_{\flat}$ decomposition to the first $n$ layers near the input (instead of only the first *one* layer) provides control over the local and semantic scale [26] of the computed attributions (see Fig. 1*(e)-(g)*). Previous works profit from this option for comparing DNNs of varying depth, and differently configured convolutional stacks [14], or by increasing readability of attributions maps aligned to the requirements of human investigators [15].

## III. METRIC AND ASSUMPTIONS

### A. Motivation

The declared purpose of LRP is to precisely and quantitatively inform about the (image or intermediate) features which contribute towards or against the decision of a model wrt. to a specific predictor output [7]. While the recent $\boldsymbol{LRP}_{CMP}$ exhibits improved properties above previous variants of LRP *by eyeballing*, an objective verification requires quantification. The visual object detection setting, as it is described by the Pascal VOC (PVOC) [20] or ImageNet [31] datasets — both
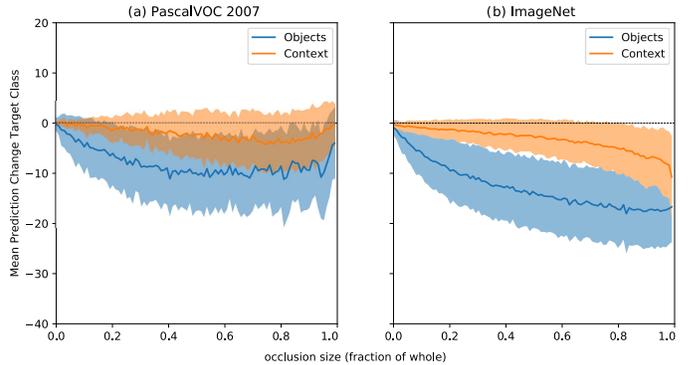


Fig. 2. Mean prediction changes $\Delta f(x)$ measured in the logit outputs of the true class as a function of the occluded area, when occluding the pixels within (object) and without (image context) the class-specific bounding boxes on PVOC 2007 (*left*) and ImageNet (*right*). Lower values indicate a stronger reaction of the model. Shaded areas show the standard deviation.

of which include object bounding box annotations — delivers an optimal experimental setting for this purpose.

An assumed ideal model would, in such a setting, exhibit true object understanding by only predicting based on the object itself. A good and *representative* attribution method should therefore reflect that object understanding of the model closely *i.e.* by marking (parts of) the shown object as relevant and disregarding visual features not representing the object itself. Similar to [11], we therefore rely on a measure based on localization of attribution scores. In the following, we will evaluate $\boldsymbol{LRP}_{CMP}$ against other methods and variants of LRP on ImageNet using a pre-trained VGG-16 network, and on PVOC 2007 using a pre-trained (on PVOC 2012) CaffeNet model [11]. Both models perform well on their respective task and have been obtained from https://modelzoo.co/.

### B. Verifying Object-centricity During Prediction

In practice, both datasets can not be assumed to be free from contextual biases (*cf.* [10], [32]), and in both settings models are trained to categorize images rather than localize objects. Still, we (necessarily) assume that the models we use dominantly base their decision on the target object, as opposed to the image context.

We verify our hypothesis in Fig. 2, by showing for both models and datasets the reaction of the corresponding predictor $f$ to the occlusion of the object area vs. the occlusion of the image background. That is, for each image $x$ of the respective dataset, we leverage the available bounding box annotations and compute partially occluded versions $x'$ where either the object area or class-specific image background (*i.e.* the non-object area) are replaced with mean color values per corresponding pixel and dataset. We then measure the $\Delta f(x) = f(x') - f(x)$ for the ground truth label(s) of $x$ based on the network's logit outputs, and plot this value as a function of relative bounding box size. Fig. 2 shows the average values and standard deviation for $\Delta f(x)$ per bounding box size (discretized into 100 uniform bins) when replacing

either the object (area within the bounding box) or the context (rest of the image).

Occluding the object area consistently leads to a sharper decrease in the output for the specific class. The trend is especially evident for smaller objects. This supports our claim that the networks base their decision mainly on the object itself.

### C. Attribution Localization as a Quantitative Measure

This gives us a performance criterion for attribution methods in object detection and classification. In order to track the fraction of the total amount of relevance that is attributed to the object, we use the inside-total relevance ratio $\mu$ without, and a weighted variant $\mu_w$ within consideration of the object size:

$$\mu = \frac{R_{\text{in}}}{R_{\text{tot}}} \qquad \mu_w = \mu \cdot \frac{S_{\text{tot}}}{S_{\text{in}}} \qquad (5)$$

While conceptually similar to the inside-outside ratio used in [11], $\mu$ and $\mu_w$ avoid numerical issues in edge cases wrt. bounding box size. Here, $R_{\text{in}}$ is the sum of positive relevance in the bounding box, $R_{\text{tot}}$ the total sum of positive relevance in the image and $S_{\text{in}}$ and $S_{\text{tot}}$ are the size of the bounding box and the image respectively, in pixels. The subscript $w$ signals the addition of a normalization factor in $\mu_w$ considering the size of image and object.

Correctly locating small objects is more difficult than locating image-sized objects. Since the ratio $S_{\text{tot}}/S_{\text{in}}$ is always greater than or equal to 1 and increases for smaller objects, $\mu_w$ puts additional emphasis on measuring the outcome for small bounding box sizes. In both cases, higher values indicate larger fractions of relevance attributed to the object area (and not background), and therefore are the desirable outcome.

## IV. EXPERIMENTS AND RESULTS

### A. Experimental Setup

We perform our experiments on both the ImageNet and the PVOC 2007 datasets, since both collections provide large numbers of ground truth object bounding boxes.

For PVOC, we compute attribution maps for all samples (approx. 10.000) from PVOC 2007, using a model which has been pre-trained on the multi label setting of PVOC 2012 [11], [20]. The respective model performs with a mean AP of 72.12 on PVOC 2007. Since PVOC describes a multi label setting, multiple classes can be present in the same image. We therefore evaluate $\mu$ and $\mu_w$ once for each unique existing pair of { class × sample }, yielding approximately 15.000 measurements. Images with a higher number of (smaller) bounding boxes thus effectively have a stronger impact on the results than images with larger (and fewer), image-filling objects, while at the same time describing a *more difficult* setting. Many of the objects shown in PVOC images are not centered. In order to use all available object information in our evaluation, we rescale the input images to the network's desired input shape, avoiding the (partial) cropping of objects.

On ImageNet [19] (2012 version), bounding box information does only exist for the 50.000 validation samples (displaying one class per image) and can be downloaded from the official website[3]. We evaluate a pre-trained VGG-16 model from the keras model zoo, obtained via the iNNvestigate [33] toolbox. The model performs with a 90.1% top-5 accuracy on the ImageNet test set. For all images the shortest side is rescaled to fit the model input and the longest side is center-cropped to obtain a quadratic input shape. Bounding box information is adjusted correspondingly.

For computing attribution maps, we make use of existing XAI software packages, depending on the models' formats. That is, for the VGG-16 model we use the Keras [34] and Tensorflow [35] based iNNvestigate [33] toolbox. For the PVOC data and the CaffeNet architecture, we compute attributions using the Caffe [36] based LRP Toolbox [27].

Both XAI packages support the same functionality regarding LRP, yet differ in the provided selection of other attribution methods. Our study, however, shall be focussed on the beneficial or detrimental effects between the variants of LRP used in literature.

We compute attribution maps and values for $\mu$ and $\mu_w$ for both models and different variants of LRP: $\boldsymbol{LRP}_z$, $\boldsymbol{LRP}_{\alpha\beta}$ (both for $\alpha = 1$ and $\alpha = 2$), and several parameterizations of $\boldsymbol{LRP}_{CMP}$. For the latter we distinguish parameter choices for $\alpha$ in a subscript when discussing quantitative results in Sec. IV-C. Additionally, in case $LRP_\flat$ is applied to the input layer, we add "$+\flat$" to the subscript, *e.g.* as "$\boldsymbol{LRP}_{CMP:\alpha1+\flat}$".

We complement the results with Guided Backprop (GB) [3] and for ImageNet with Pattern Attribution (PA) [4] only available in iNNvestigate. On both datasets, we evaluate attributions for the ground truth class labels, independent of the network prediction.

### B. Qualitative Observations

Fig. 1 exemplarily shows attribution maps computed with different methods based on the VGG-16 model, for two object classes present in the ImageNet labels and the input image; "Bernese Mountain Dog" and "Tiger Cat". Attributions in Figs. 1(a)-(d) result from uniform rule application to the whole network. Next to applications of $\boldsymbol{LRP}_z$ and $\boldsymbol{LRP}_{\alpha\beta}$ with $\alpha = 1$, this includes Guided Backprop [3] and Pattern Attribution [4]. Neither of these maps demonstrate class-discriminativeness and prominently attribute scores to the same areas, regardless of the target class chosen for attribution. $\boldsymbol{LRP}_z$ additionally shows the effects of gradient shattering in a highly complex attribution structure due to its equivalence to G×I. Such attributions would be difficult to use and juxtapose in further algorithmic or manual analyses of model behavior.

To the right, attribution maps in Figs. 1(e)-(g) correspond to variants of $\boldsymbol{LRP}_{CMP}$, which apply different decomposition rules depending on layer type and position. In Fig. 1(e), the $LRP_\flat$-rule is not applied at all, while in Fig. 1(f) it is used for the first three convolutional layers, and the whole
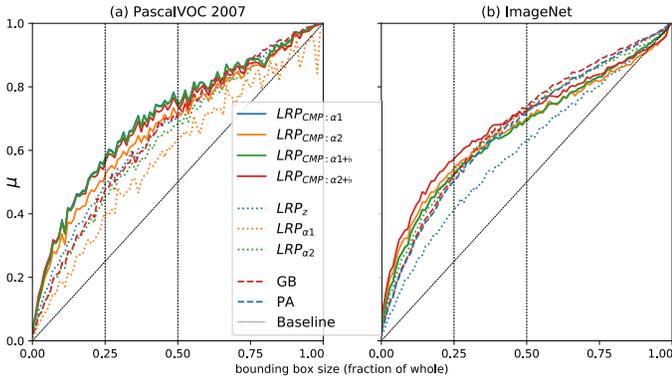
---

[3]http://www.image-net.org/challenges/LSVRC/2012

Fig. 3. Average in-total ratio $\mu$ as a function of bounding box size. Vertical lines mark thresholds of $25\%$ and $50\%$ covered image area. The baseline can be reached by uniformly attributing to all pixels of the image. Higher values are better.

TABLE I
AVERAGE CONTEXT ATTRIBUTION METRICS FOR DIFFERENT ANALYZERS AND DATASETS. ROW ORDER IS DETERMINED BY $\mu_w$. HIGHER $\mu_*$ ARE BETTER.

| Data | Analyzer | $\mu_w$ | $\mu_{<0.25}$ | $\mu_{<0.5}$ | $\mu$ |
|---|---|---|---|---|---|
| PVOC (CaffeNet) | $LRP_{CMP:\alpha2+\flat}$ | **2.716** | **0.307** | 0.421 | 0.532 |
| | $LRP_{CMP:\alpha1}$ | 2.664 | 0.306 | **0.426** | **0.539** |
| | $LRP_{CMP:\alpha1+\flat}$ | 2.598 | 0.301 | 0.421 | 0.535 |
| | $LRP_{CMP:\alpha2}$ | 2.475 | 0.276 | 0.388 | 0.504 |
| | $LRP_z$ | 2.128 | 0.236 | 0.353 | 0.480 |
| | $GB$ | 1.943 | 0.212 | 0.335 | 0.470 |
| | $LRP_{\alpha2}$ | 1.843 | 0.205 | 0.320 | 0.452 |
| | $LRP_{\alpha1}$ | 1.486 | 0.163 | 0.273 | 0.403 |
| | Baseline | 1.000 | 0.100 | 0.186 | 0.322 |
| ImageNet (VGG-16) | $LRP_{CMP:\alpha2+\flat}$ | **1.902** | **0.397** | **0.534** | **0.714** |
| | $LRP_{CMP:\alpha2}$ | 1.797 | 0.368 | 0.505 | 0.693 |
| | $LRP_{CMP:\alpha1}$ | 1.7044 | 0.3467 | 0.4887 | 0.6898 |
| | $LRP_{CMP:\alpha1+\flat}$ | 1.7043 | 0.3466 | 0.4886 | 0.6898 |
| | $LRP_{\alpha2}$ | 1.702 | 0.332 | 0.496 | 0.706 |
| | $GB$ | 1.640 | 0.312 | 0.485 | 0.710 |
| | $LRP_{\alpha1}$ | 1.609 | 0.306 | 0.475 | 0.699 |
| | $PA$ | 1.591 | 0.303 | 0.471 | 0.698 |
| | $LRP_z$ | 1.347 | 0.236 | 0.389 | 0.632 |
| | Baseline | 1.000 | 0.128 | 0.260 | 0.547 |

convolutional stack — including pooling layers — in Fig. 1(g). Both heatmaps in Fig. 1(e) and Fig. 1(f) use $\alpha = 1$. Here altogether, the visualized attribution maps correspond more to an "intuitive expectation" of how relevance should be attributed compared to Figs. 1(a)-(d), assuming a model predicts based on object understanding. Figs. 1(e)-(g) demonstrate the change in scale and semantic, from attributions to local features to a very coarse localization map, with changing placements of the $LRP_\flat$-rule. Further, it becomes clear that with an application of the $LRP_{\alpha\beta}$-rule in upper layers, object localization is lost (see Fig. 1(b) vs. Fig. 1(g)), while an application in lower layers avoids issues related to gradient shattering, as shown in Figs. 1(e)-(f) compared to Fig. 1(a).

Note that the special case shown in Fig. 1(g) is highly similar to an application of the Class Activation Mapping (CAM) [37] method in the fully connected part of the model, however replaces the upsampling over the model's convolutional stack of the CAM approach with the $LRP_\flat$ decomposition based approach of the LRP framework, and is thus naturally capable of distributing negative relevance scores.

Note that the VGG-16 network used here never has been trained in a multi-label setting. Despite only receiving one object category per input sample, it has learned to distinguish between different object types shown in the same image, *e.g.* that a dog is not a cat. This in turn reflects well in the attribution maps computed after the $LRP_{CMP}$ pattern.

Further examples akin to Fig. 1 are given in the Appendix.

*C. Quantitative Results*

Figs. 3(a) and (b) show the average in-total ratio $\mu$ as a function of bounding box size, discretized over 100 equally spaced intervals, for PVOC 2007 and ImageNet. Averages for $\mu$ and $\mu_w$ over the whole (and partial) datasets can be found in Tab. I. Large values indicate more precise attribution to the relevant object.

The inside-total relevance ratio highly depends on the size of the bounding box. In addition to the average $\mu$ and $\mu_w$ as an aggregate over all classes and images, we also report $\mu_{\leq0.25}$ and $\mu_{\leq0.5}$, the average values over all objects whose

bounding box does not span more than $0.25$ and $0.5$ times the area of the whole image respectively. The assumed *Baseline* is the uniform attribution of relevance over the whole image, which is outperformed by all methods.

$LRP_z$ performs noticeably worse on ImageNet than on PVOC, which we trace back to the significant difference in model depth (13 vs 21 layers) affecting gradient shattering. We omit $LRP_\varepsilon$ in Tab. I due to the identity in results to $LRP_z$. $LRP_{\alpha\beta}$ has the tendency to attribute to all shown objects (via generally neuron-activating features) and suffers from the multiple object classes per image in PVOC, where ImageNet shows only one class. Also, the similarity of attributions between *PA* and $LRP_{\alpha\beta}$ with $\alpha = 1$ observed in Fig. 1 seem consistent on ImageNet and result in close measurements in Tab. I.

Tab. I demonstrates that $LRP_{CMP}$ clearly outperforms other methods consistently on large datasets. That is, the increased precision in attribution to relevant objects is especially evident in the presence of smaller bounding boxes in $\mu_w$. This can also be seen in $\mu_{\leq0.25}$ and $\mu_{\leq0.5}$ in Tab. I and the left parts of Figs. 3(a) and (b), where a majority of the image shows contextual information or other classes. Once bounding boxes become (significantly) larger and cover over $50\%$ of the image, all methods converge towards perfect performance, as expected. In both settings, $LRP_{CMP:\alpha2+\flat}$ yields the best results, while overall the composite strategy is more effectful than a fine tuning of decomposition rule parameters.

*D. Conclusion*

In this study, we discuss a recent development in the application of Layer-wise Relevance Propagation. We summarize this emerging strategy of a composite application of multiple purposed decomposition rules as $LRP_{CMP}$ and juxtapose its effects to previous approaches to LRP and other methods, which uniformly apply a single decomposition rule to all

layers of the model. For the first time, our results show that $LRP_{CMP}$ does not only yield *measurably* more representative attribution maps, but also provides a solution against gradient shattering affecting previous approaches, and improves properties related to object localization and class discrimination via attribution. Moreover, $LRP_{CMP}$ is able to precisely attribute negative relevance scores to class-contradicting features while requiring only one modified backward pass though the model, using established tools from the LRP framework. The discussed beneficial effects are demonstrated qualitatively and verified quantitatively at hand of two large and widely used computer vision datasets.

## REFERENCES

[1] Parliament and Council of the European Union, "General data protection regulation," 2016.

[2] B. Goodman and S. R. Flaxman, "European Union regulations on algorithmic decision-making and a "right to explanation"," *AI Magazine*, vol. 38, no. 3, pp. 50–57, 2017.

[3] J.T. Springenberg, A. Dosovitskiy, T. Brox, and M.A. Riedmiller, "Striving for simplicity: The all convolutional net," in *Proc. of International Conference on Learning Representations (ICLR)*, 2015.

[4] P.-J. Kindermans, K.T. Schütt, M. Alber, K.-R. Müller, D. Erhan, B. Kim, and S. Dähne, "Learning how to explain neural networks: Patternnet and patternattribution," in *Proc. of International Conference on Learning Representations (ICLR)*, 2018.

[5] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *Proc. of International Conference on Machine Learning (ICML)*, 2017, pp. 3319–3328.

[6] D. Smilkov, N. Thorat, B. Kim, F.B. Viégas, and M. Wattenberg, "Smoothgrad: removing noise by adding noise," *CoRR*, vol. abs/1706.03825, 2017.

[7] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, "On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation," *PloS one*, vol. 10, no. 7, pp. e0130140, 2015.

[8] Y. Yang, V. Tresp, M. Wunderle, and P. A. Fasching, "Explaining therapy predictions with layer-wise relevance propagation in neural networks," in *Proc. of IEEE International Conference on Healthcare Informatics (ICHI)*, 2018, pp. 152–162.

[9] A W Thomas, H R Heekeren, K-R Müller, and W Samek, "Analyzing neuroimaging data through recurrent deep learning models," *Frontiers in Neuroscience*, vol. 13, pp. 1321, 2019.

[10] S. Lapuschkin, S. Wäldchen, A. Binder, G. Montavon, W. Samek, and K.-R. Müller, "Unmasking clever hans predictors and assessing what machines really learn," *Nature Communications*, vol. 10, no. 1, pp. 1096, 2019.

[11] S. Lapuschkin, A. Binder, G. Montavon, K.-R. Müller, and W. Samek, "Analyzing classifiers: Fisher vectors and deep neural networks," in *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2912–2920.

[12] M. Ancona, E. Ceolini, C. Öztireli, and M. Gross, "Gradient-based attribution methods," in *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pp. 169–191. Springer, 2019.

[13] G. Montavon, A. Binder, S. Lapuschkin, W. Samek, and K.-R. Müller, "Layer-wise relevance propagation: an overview," in *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pp. 193–209. Springer, 2019.

[14] S. Lapuschkin, A. Binder, K.-R. Müller, and W. Samek, "Understanding and comparing deep neural networks for age and gender classification," in *Proc. of IEEE International Conference on Computer Vision Workshops (ICCVW)*, 2017, pp. 1629–1638.

[15] M. Hägele, P. Seegerer, S. Lapuschkin, M. Bockmayr, W. Samek, F. Klauschen, K.-R. Müller, and A. Binder, "Resolving challenges in deep learning-based analyses of histopathological images using explanation methods," *CoRR*, vol. abs/1908.06943, 2019.

[16] L. Y. W. Hui and A. Binder, "Batchnorm decomposition for deep neural network interpretation," in *International Work-Conference on Artificial Neural Networks (IWANN)*, 2019, pp. 280–291.

[17] D. Balduzzi, M. Frean, L. Leary, J.P. Lewis, K. W.-D. Ma, and B. McWilliams, "The shattered gradients problem: If resnets are the answer, then what is the question?," in *Proc. of International Conference on Machine Learning (ICML)*, 2017, pp. 342–350.

[18] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, and K.-R. Müller, "Evaluating the visualization of what a deep neural network has learned," *IEEE Transactions on Neural Network Learning Systems*, vol. 28, no. 11, pp. 2660–2673, 2017.

[19] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, Ma S., Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.

[20] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results," "http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html".

[21] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, et al., "Gradient-based learning applied to document recognition," *Proc. of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[22] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2818–2826.

[23] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations (ICLR)*, 2015.

[24] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, and K.-R. Müller, "Explaining nonlinear classification decisions with deep taylor decomposition," *Pattern Recognition*, vol. 65, pp. 211–222, 2017.

[25] J. Zhang, S. A. Bargal, Z. Lin, J. Brandt, X. Shen, and S. Sclaroff, "Top-down neural attention by excitation backprop," *International Journal of Computer Vision (IJCV)*, vol. 126, no. 10, pp. 1084–1102, 2018.

[26] S. Bach, A. Binder, K.-R. Müller, and W. Samek, "Controlling explanatory heatmap resolution and semantics via decomposition depth," in *Proc. of IEEE International Conference on Image Processing (ICIP)*. IEEE, 2016, pp. 2271–2275.

[27] S. Lapuschkin, A. Binder, G. Montavon, K.-R. Müller, and S. Samek, "The LRP toolbox for artificial neural networks," *Journal of Machine Learning Research (JMLR)*, vol. 17, pp. 114:1–114:5, 2016.

[28] F. Horst, S. Lapuschkin, W. Samek, K.-R. Müller, and W. I. Schöllhorn, "Explaining the unique nature of individual gait patterns with deep learning," *Scientific Reports*, vol. 9, pp. 2391, 2019.

[29] I. Sturm, S. Lapuschkin, W. Samek, and K.-R. Müller, "Interpretable deep neural networks for single-trial eeg classification," *Journal of Neuroscience Methods*, vol. 274, pp. 141–145, 2016.

[30] J. Gu, Y. Yang, and V. Tresp, "Understanding individual decisions of cnns via contrastive backpropagation," in *Proc. of Asian Conference on Computer Vision (ACCV)*, 2018, pp. 119–134.

[31] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. of Advances in Neural Information Processing Systems (NIPS)*, 2012, pp. 1106–1114.

[32] C. J. Anders, T. Marinč, Neumann D., W. Samek, K.-R. Müller, and S. Lapuschkin, "Analyzing imagenet with spectral relevance analysis: Towards imagenet un-hans'ed," *CoRR*, vol. abs/1912.11425, 2019.

[33] M. Alber, S. Lapuschkin, P. Seegerer, M. Hägele, K. T. Schütt, G. Montavon, W. Samek, K.-R. Müller, S. Dähne, and P.-J. Kindermans, "innvestigate neural networks!," *Journal of Machine Learning Research (JMLR)*, vol. 20, pp. 93:1–93:8, 2019.

[34] F. Chollet et al., "Keras," https://keras.io, 2015.

[35] M. Abadi et al., "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *CoRR*, vol. abs/1603.04467, 2016.

[36] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proc. of 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 675–678.

[37] B. Zhou, A. Khosla, À. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2921–2929.

In Fig. 4 we provide further illustrative examples similar to Fig. 1, using different input images and object classes.



Fig. 4. Different attributions for the output classes "Bernese Mountain Dog" and "French Bulldog" (A), "Persian Cat" and "Siamese Cat" (B), "Zebra" and "African Elephant" (C) and "Sunglasses" and "Windsor Tie" (D and E), using the pretrained VGG-16 model. For details *cf.* Fig. 1.