

Analyzing ImageNet with Spectral Relevance Analysis: Towards ImageNet un-Hans’ed

Christopher J. Anders¹, Talmaj Marinč², David Neumann², Wojciech Samek^{2,*},
Klaus-Robert Müller^{1,3,4,*}, and Sebastian Lapuschkin^{2,*}

¹Dept. of Electrical Engineering and Computer Science, Technische Universität Berlin, Germany

²Dept. of Video Coding and Analytics, Fraunhofer Heinrich Hertz Institute, Berlin, Germany

³Dept. of Brain and Cognitive Engineering, Korea University, Seoul, Korea

⁴Max Planck Institute for Informatics, Saarbrücken, Germany

*corresponding author

Abstract

Today’s machine learning models for computer vision are typically trained on very large (benchmark) data sets with millions of samples. These may, however, contain biases, artifacts, or errors that have gone unnoticed and are exploited by the model. In the worst case, the trained model may become a ‘Clever Hans’ predictor that does not learn a valid and generalizable strategy to solve the problem it was trained for, but bases its decisions on spurious correlations in the training data. Recently developed techniques allow to explain individual model decisions and thus to gain deeper insights into the model’s prediction strategies. In this paper, we contribute by providing a comprehensive analysis framework based on a scalable statistical analysis of attributions from explanation methods for large data corpora, here ImageNet. Based on a recent technique – Spectral Relevance Analysis (SpRAy) – we propose three technical contributions and resulting findings: (a) novel similarity metrics based on Wasserstein for comparing attributions to allow for the first time scale, translational, and rotational invariant comparisons of attributions, (b) a scalable quantification of artifactual and poisoned classes where the ML models under study exhibit Clever Hans behavior, (c) a cleaning procedure that allows to relief data of artifacts and biases in a systematic manner yielding significantly reduced Clever Hans behavior, i.e. we un-Hans the ImageNet data corpus. Using this novel method set, we provide qualitative and quantitative analyses of the biases and artifacts in ImageNet and demonstrate that the usage of these insights can give rise to improved models and functionally

cleaned data corpora.

1. Introduction

Throughout the last decade, (black box) machine learning (ML) techniques have made impressive performance leaps on even the most complex tasks [1–4], especially in the form of Deep Neural Networks (DNN) [5]. These models are typically trained (or pretrained) on very large datasets, *e.g.*, ImageNet [6], with millions of samples. Recently, it was discovered that biases, spurious correlations, as well as errors in the training dataset may have a detrimental effect on the training resulting in “Clever Hans” predictors [7], which only superficially solve the task they have been trained for. Unfortunately, due to the immense size of today’s datasets, a direct manual inspection and removal of artifactual samples can be regarded hopeless. However, analyzing the biases and artifacts in the *model* instead, may provide insights about the biases and artifacts in the training data indirectly. For that purpose we would, however, need to inspect the learning models and operate them beyond black box mode.

Only recently methods of explainable AI (XAI) (*c.f.* [8] for an overview) were developed. They provide deeper insights into how an ML classifier arrives at its decisions and potentially help to unmask Clever Hans predictors. XAI methods can be roughly categorized into two groups: methods providing *local* explanations and those providing *global* explanations [9]. Here, *local* explanations increase transparency on individual predictions of the model and assess the importance of

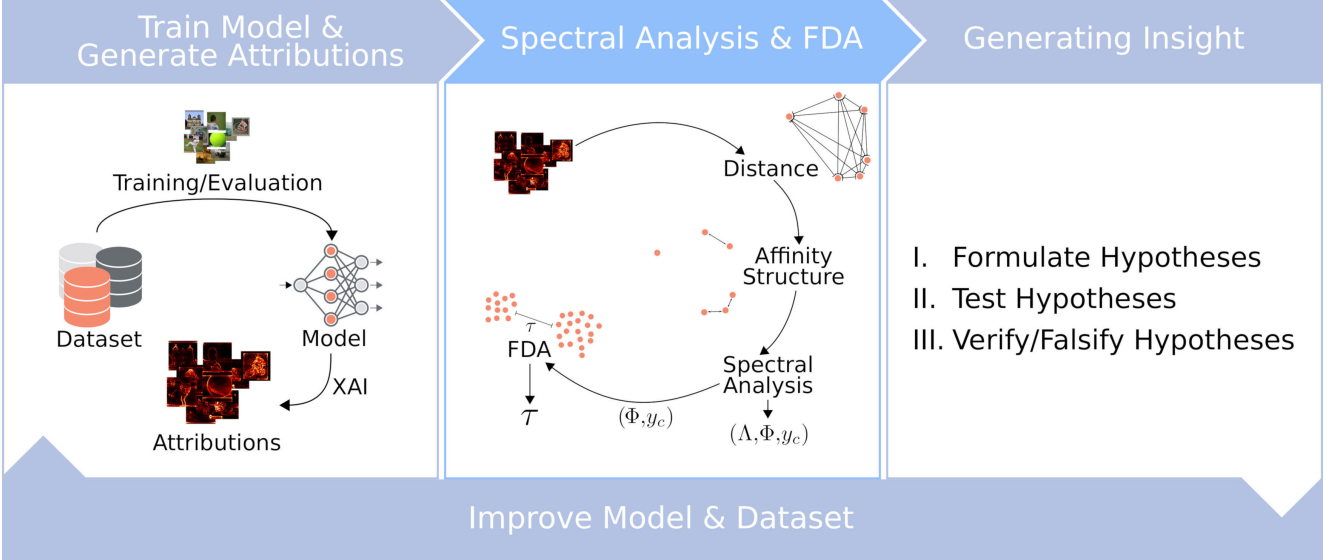


Figure 1. Overview of the SpRAY approach. *Left*: Large corpora of data can be used to train models for specific tasks. To gain insights into local model behavior, explanation methods can be employed. *Middle*: Using SpRAY, one can deduce global model behavior from a set of local explanations (see Algorithm 1). *Right*: Based on this analysis, *striking* classification strategies can be identified and further investigated. Obtained insights can be used to improve the model and/or the dataset.

input features w.r.t. specific samples. Local explanations are commonly presented in the form of attribution or heatmaps aligned to the input space, which can be computed *e.g.* from propagation-based (*e.g.* [10–13]) or surrogate-based techniques (*e.g.* [14–16]). *Global* methods on the other hand aim at obtaining model understanding by assessing the general importance of features a model relies on. These features often are high-level concepts and are either chosen/designed manually, or directly and discretely are accessible in input space [17–19].

Both the local and global approaches suffer from a (human) investigator bias during analysis, thus are of limited use when searching for biases, spurious correlations, and errors in the training dataset. Global methods can only measure the impact of pre-determined, and expected or known features (*c.f.* [18, 19]), which limits the applicability when aiming to discover *all* behavioral facets of a model (including specifically the ones unknown beforehand). Local methods, on the other hand, have the potential to provide much more detailed information per sample, but compiling information about model behavior over thousands or millions of samples and explanations is a tiring and laborious process. Furthermore, the success of such an analysis depends on the examiner’s keen perception and domain knowledge, at times limiting the potential for knowledge discovery.

Our current contribution aims at bridging the vast gap between the discussed two extremes in an *objective*

and *automated* manner making use of scalable statistical inference on millions of heatmaps, specifically employing the Spectral Relevance Analysis (SpRAY) [7] technique. SpRAY constitutes a semi-automated statistical analysis of large numbers of locally explaining attribution maps, with the intent of automated summarization and clustering of model strategies via local attribution maps. Thus, SpRAY is aligned with the assumption that an ML model might base its predictions on multiple sub-strategies (instead of only globally effective ones) for recognizing a target class and distinguishing it from others. Early applications of SpRAY [7] demonstrate its utility for knowledge discovery via strategy summarization on gameplay sequences of Atari-2600 playing DNN agents [2], and the detection of multiple flaws in dataset composition and model architecture design in the context of the formerly widely used Pascal VOC image recognition benchmark [7, 20].

In this paper, we extend SpRAY to make it better applicable for large-scale analyses on datasets with hundreds of classes and millions of samples, such as ImageNet [6]. Our technical contributions are: (a) a new Wasserstein-based similarity metric for scale, translational, and rotational invariant comparisons of attributions, (b) the identification of artifactual and biased samples in the ImageNet corpus, and a quantitative analysis of their impact on the Clever Hans’ness of the classifier, (c) a systematic cleaning procedure for the artifacts and biases to reduce Clever-Hans behavior,

i.e. we un-Hans the ImageNet dataset. These allow interesting findings that are illuminating beyond our specific technical approach.

2. Methods

We will first briefly summarize SpRAY from [7] (see Figure 1 for a procedural overview), emphasizing and motivating where and how we go beyond [7]. An algorithmic summary of the technique can be found in Algorithm 1.

2.1. Spectral Relevance Analysis brought to scale

The SpRAY technique is a meta-analysis tool for finding patterns in model behavior, given sets of instance-based explanatory attribution maps. The algorithm is based on Spectral Clustering (SC) [21, 22] and performs the following sequence of computations.

Computing attributions SpRAY analyzes the (spatial) structure described by a set of heat/attribution maps, each locally explaining single model decisions w.r.t. to a model prediction the user is interested in. Following [7], we provide attributions computed with Layer-wise Relevance Propagation (LRP) [10] according to the recommended composite (or layer-dependent) strategy [23, 24]. Specifically, since our analyses are based on a pre-trained VGG-16 [25] type DNN we follow [24] and apply LRP_ϵ to the model’s dense layers, LRP_b in the lowest convolutional layer and $LRP_{\alpha\beta}$ ($\alpha=1$) in all other convolutional layers, using the pre-configured LRP-analyzers of the iNNvestigate [26] toolbox. We sum attribution scores along the color channel axis to obtain a single attribution value per pixel coordinate.

Preprocessing of attributions The work of [7] analyzes the behavior of DNN predictors, as well as a former state-of-the-art model from the bag-of-words family, the improved Fisher Kernel SVM [27]. Other than the DNN, the latter predictor does not expect inputs of a fixed size and therefore LRP computes non-uniformly sized attribution maps over the analyzed data. The authors of [7] resort to sum-pooling attribution scores from arbitrarily sized explanation maps onto a 20×20 sized grid, resulting in a 400-dimensional representation per attribution map for further processing. The authors justify the often extreme size or dimensionality reduction with an increased (regional) stability of the analysis and a decrease in computational cost. They also point out that this step – albeit useful – is no practical necessity.

Since in this paper we only process uniformly sized attribution maps of 224×224 pixels as output of LRP and the VGG-16 model, we omit the optional preprocessing and rather preserve the complete and unaltered

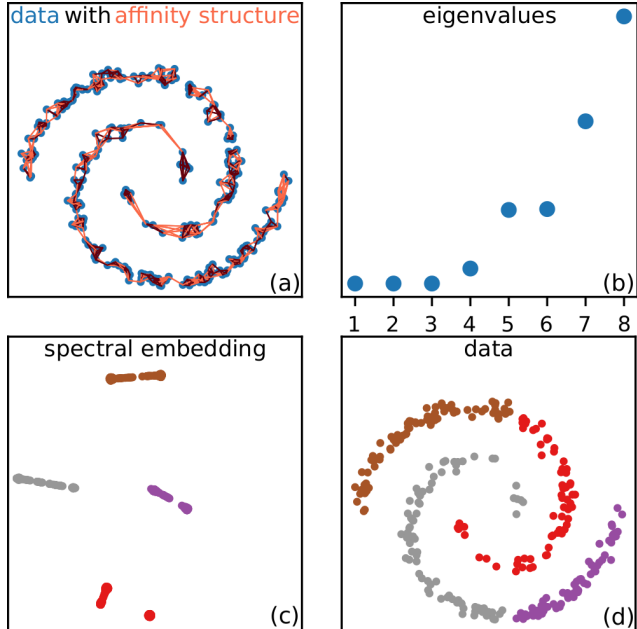


Figure 2. Toy example on spectral clustering. (a) Data with affinity structure based on KNN. (b) Eigenvalues computed from the graph laplacian. Note the large gap after the fourth eigenvalue, indicating four clusters. (c) Spectral embedding Φ of the input data, labelled using k -means clustering with $k=4$, with the choice of k depending on the eigenvalue spectrum. (d) Embedding labels assigned to the corresponding original data.

structural information within the LRP heatmap attributions. We show that albeit absent this preprocessing step, we can discover class-specific and consistent clever-hans strategies on ImageNet.

Computing distances and affinities An ingredient of SpRAY as performed in [7] is a comparison of heatmaps based on the euclidean distance, which focuses on the task of finding *regional consistencies* between attribution maps. In order to further expand the scope of SpRAY, we consider the Gromov-Wasserstein distance (GWD) [28] as an additional distance measure.

As [7], we compute pair-wise (binary) affinity scores a_{ij} between all samples i, j using the k -nearest-neighbor (KNN) algorithm [29] on the previously computed distance measures and store them in a matrix A . We visualize this step of the procedure, along with all following steps, in Figure 2, intuitively demonstrating of the spectral Clustering (SC) process.

We find that for ImageNet classes, choosing a fixed $k=10$ instead of $k=\lceil \log(n) \rceil=8$ (*c.f.* [7, 30]. $n=1,300$ is the number of samples per ImageNet object category) yields the most consistent and robust groupings throughout the spectral analysis. Values for k larger

than 10 seem to have no, or only little effect, while smaller values lead to an overly-fractured affinity structure. The affinity matrix A is then symmetrized with

$$A \leftarrow \frac{1}{2} (A + A^\top) , \quad (1)$$

which leads to $\forall i, j : a_{ij} = a_{ji}$ and $a_{ij} = 1$ if i and j already have been mutually connected prior to the symmetrization, and $a_{ij} = 0.5$ otherwise.

Computing the spectral embedding The computation of the spectral embedding is the analytic core of the SpRAY method. Given a matrix A describing the affinity structure of the data, following [7] we first compute the *symmetrized and normalized* p.s.d. graph laplacian [22, 30]

$$L_{\text{sym}} = D^{-1/2} L D^{-1/2} , \quad (2)$$

where D is the diagonal degree matrix of the connectivity graph described by A , and $L = D - A$ the *un-normalized* graph laplacian [30]. The entries d_i of the diagonal matrix D are computed as

$$d_i = \sum_{j=1}^n a_{ij} . \quad (3)$$

Then, an eigenvalue decomposition of L_{sym} is performed, which yields eigenvalues $\Lambda = \{\lambda_i\}_{i=1\dots q}$ and eigenvectors as *columns* of $\Phi \in \mathbb{R}^{n \times q}$.

The set of returned eigenvalues informs about the cluster structure discovered in the data [30], where completely separated clusters are indicated by the number of eigenvalues $\lambda_i=0$, which is rarely the case with real world data. Figure 2(b) demonstrates that the structure and number of embedded clusters can also be inferred from the *eigengap*, a sudden increase in difference of neighboring eigenvalues. The first eigengap in the toy example identifies four almost completely disjoint groups of points.

The *rows* in Φ constitute the *spectral embedding* of the data in \mathbb{R}^q , reflecting the affinity structure encoded in A and L_{sym} . Cluster labels can be assigned to the embedding by using *e.g.* k -means-clustering on the spectral embedding. In the toy example in Figure 2(c) we choose $k=4$ according to the eigenvalue spectrum Λ and show the assigned labels using color coding on the spectral embedding projected to \mathbb{R}^2 . Due to the correspondence of the input data to the rows of Φ , we can directly assign the computed labels in input space as well (Figure 2(d)).

On ImageNet, where $n=1,300$, we restrict $q=32$ for the computation of the eigen-decomposition. Specifically, we use the computationally efficient and iterative

Lanczos algorithm¹ [32], which has in general a computational complexity of $\mathcal{O}(n^3)$, which reduces to $\mathcal{O}(qn^2)$ when only the $q \leq n$ most discriminative eigenvectors are to be computed.

2.2. Alternative distance measures

SpRAY has originally used euclidean distances to compute the neighborhood graph [7]. In the application of images, this means that image similarity is identified by spatial properties, *i.e.* having the same attribution intensities at the same pixel renders high similarity. This is a reasonable approach, especially if one would like to focus on spatial properties such as watermarks or padding. However, when the domain of interest are spatially unrelated shapes or color distributions, other measures of similarity may be needed.

A recently very popular distance metric is the Optimal Transport, or *Wasserstein-Distance*. In the context of computer vision, it is also known as the *Earth-Mover’s Distance* [33]. Its benefit is that it “feels” like a very natural distance metric [34].

Wasserstein distances use distances between spatially fixed points over the same identical image grid. Gromov-Wasserstein [28] distances matches points by their pairwise distances, instead of using a fixed image grid with a fixed amount of points. This means that however points are spatially distributed, if in both sets there are points whose pairwise relations are similar, then their Gromov-Wasserstein distance will be small. A somewhat intuitive visualization of euclidean distance, Wasserstein distance, and Gromov-Wasserstein distances is shown in Figure 3. We show 4 samples of hand-written digits [35] in 4 corners, translated and rotated. All images that lie on the line between the corners are barycenters [36] of the corner images, weighted by the Chebyshev distance to all samples. The metrics used to compute the barycenters are the 3 previously mentioned metrics. Wasserstein barycenters are computed as in [34]. For the Gromov-Wasserstein distance, we need to compute pairwise distances between points in the image. Points are extracted from the images by choosing each pixel one after another, starting with the largest, until 99 percent of the total sum of all pixel values is reached. We can nicely see that the Wasserstein distance seems translation invariant, but fails with different rotations. Gromov-Wasserstein distance shows to be invariant to rotation, translation, and mirroring, since all the information is contained in only the pairwise relations. Thus, to enable invariant comparison and in this manner maximally match regionally independent shapes in our analysis, we use the Gromov-

¹via the `sparse.linalg.eigsh` function provided by the SciPy [31] package for Python

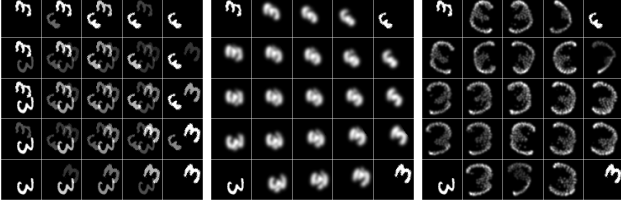


Figure 3. Barycenters of four rotated and translated MNIST digits. The original images are in the four corners. Used distance metrics are euclidean (*left*), Wasserstein distance (*middle*) and Gromov-Wasserstein distance (*right*).

Wasserstein distance measure.

2.3. Fisher Discriminant Analysis for Clever Hans identification

A critical decision in clustering approaches is the number of desired clusters. Since [7] analysed the comparatively small Pascal VOC dataset [20] with 20 classes and almost 10,000 samples, going over each class individually and exploring the eigenvalues for a possibly significant *eigengap* is a straight-forward task. For ImageNet [6] with its 1,000 classes and 1.3 million samples, looking at each and every class eigengap manually is a time-consuming and exhausting venture. We would therefore like to establish some score to measure how *interesting* a class is in terms of classification strategies that separate themselves significantly from all other possible strategies, as they are good candidates for Clever Hans behaviors. Such a score could furthermore be used to rank all classes, so that only those which produce a high cluster separability score τ , can be selected for thorough investigation. Fisher Discriminant Analysis (FDA) [37, 38] is a widely popular method for classification as well as class- (or cluster-) structure preserving dimensionality reduction. FDA finds an embedding space by maximizing between-class scatter $S^{(b)}$ and minimizing within-class scatter $S^{(w)}$. It can be understood as finding the direction(s) of maximal separability between classes. The criterion as chosen by [38] for multiple classes is given as

$$J(W) = \text{tr}((W^\top S^{(w)} W)^{-1} W^\top S^{(b)} W) \quad (4)$$

$$\hat{W} = \arg \max_W J(W) \quad (5)$$

where \hat{W} is the projection matrix that minimizes within-class scatter

$$S^{(w)} = \sum_{k=1}^K \sum_{x_i \in \mathbb{C}_k} (x_i - \mu_k)(x_i - \mu_k)^\top \quad (6)$$

while maximizing between-class scatter

$$S^{(b)} = \sum_{k=1}^K (\mu_k - \mu)(\mu_k - \mu)^\top. \quad (7)$$

Here, \mathbb{C} is the set of clusters and K its cardinality, μ_k the mean of cluster k and μ the mean over *all* samples. The projection matrix \hat{W} can be found by solving the generalized eigenvalue problem

$$S^{(b)} v = \lambda S^{(w)} v \quad (8)$$

with $\hat{W} = (v_1 | v_2 | \dots | v_d)$ where $\{v_i\}_{i=1}^d$ are the generalized eigenvectors corresponding to the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$. By reducing the projection matrix \hat{W} to only use the eigenvectors corresponding to the q largest eigenvalues and inserting into (4), we obtain a score of separability. In our application, since we already act on embedded data, we found that using only the eigenvector corresponding to the largest eigenvalue, *e.g.* the direction of largest separability, gives a reasonable score for ranking the classes. We compute k -means clusterings for each class individually with k ranging from 2 to 30, and compute the separability score for each. Then, we compute the mean separability over all clusterings of each class as τ . Classes with a high mean separability score are Clever Hans candidates since then the strategies to classify exhibit highly variable behavior. Table 1 in Section 3 lists the ImageNet classes with the highest and lowest τ values.

2.4. Cluster structure visualization

Following [7] t-distributed Stochastic Neighborhood Embedding (t-SNE) [39] is used to visualize the cluster structure described by the analyzed heatmaps. To preserve a strong connection between the cluster assignments from SC and the 2D-embeddings visualized with t-SNE, the authors re-use the affinity structure in A as input to t-SNE.

We extend this approach by using the spectral embedding Φ as input to compute visualizable point clouds in \mathbb{R}^2 . We observe significantly less overlap among differently labelled clusters in the visualized point clouds when embedding the Φ from SC instead of computing embeddings from A . See Figure 2(c) for projections of the spectral embedding into \mathbb{R}^2 , computed with t-SNE for in context of a toy example.

Alternatively to t-SNE, we consider the recent UMAP [40]. The qualitative difference between t-SNE- and UMAP-embeddings are minimal yet result in reduced intra-cluster scatter and increased inter-cluster scatter for UMAP. This is however to be expected, as we merely aim to embed Φ for visualization, and representative spectral embeddings have already been computed with SC.

Algorithm 1: Spectral Relevance Analysis	
Data:	Input samples $X = \{x\}$, a model f operating on X .
Result:	eigenvalues $\Lambda = \{\lambda\}$, spectral embeddings $\Phi \in \mathbb{R}^{n \times q}$, cluster labels $Y = \{y\}$, cluster separability score τ , visualization embeddings $Z \in \mathbb{R}^2$
	<i>/* compute attributions for $x \in X$ */</i>
1	$R = \{ \}$;
2	for $x \in X$ do
3	$R_x = \text{LRP}(f, x)$;
4	$R.append(R_x)$;
5	end
	<i>/* (optionally preprocess R) */</i>
6	for $R_x \in R$ do
7	$R_x \leftarrow \text{maybe_preprocess}(R_x)$;
8	end
	<i>/* compute affinity and laplacian */</i>
9	$A = \text{affinity}(R)$;
10	$L = \text{laplacian}(A)$;
	<i>/* compute analytic quantities and visualization */</i>
11	$\Lambda, \Phi, Y = \text{spectral}(L)$;
12	$\tau = \text{FDA}(\Phi, Y)$;
13	$Z = \text{visualization_embedding}(\Phi)$;

3. Experiments and Evaluations

We apply our extended SpRAY pipeline to the ImageNet dataset as described in Section 2 and Algorithm 1. That is, per class, we compute relevance attribution heatmaps for the training data of ImageNet using LRP w.r.t. to the true label. Note, that the proposed procedure is also readily applicable to all members of the XAI zoo (cf. [8]).

3.1. Identifying clusters and classes with FDA

From the separability scores computed with FDA on the spectral embeddings Φ we can now readily identify those classes with concentrated and isolated clusters of attribution maps. Some example scores for classes with high and low separability scores are shown in Table 1. The highest separability score is achieved by class “laptop”, of which the UMAP of its spectral embedding with a significant cluster is depicted in Figure 4 (top). Clearly the visualized cluster is extremely well separated from all the other samples. For this clustering in particular, we see that a cluster of highly similar examples results in a high separability score: The “laptop” class of the ImageNet corpus contains a set of samples showing rendered instances of laptop computers in

(on pixel-level) identical poses. In contrast, class “sliding door” achieves the lowest separability score of all classes. Figure 4 (bottom) shows an exemplary cluster for this class. Its UMAP shows a distribution of attributions that seems to be hard to separate, *i.e.* there is no decision strategy that can be clearly categorized, *e.g.* by importance of region.

top classes	τ	τ	bottom classes
laptop	4.77	0.44	fountain
stethoscope	1.28	0.44	home_theater
book_jacket	1.14	0.43	wallet
bottlecap	1.13	0.43	thresher
tennis_ball	1.13	0.43	pencil_sharpener
clumber	1.12	0.42	bannister
stole	1.06	0.41	sliding_door

Table 1. Mean attribution spectral clustering separability based on Fisher Discriminant Analysis. A high separability score τ means there are significantly different decision strategies being used, potentially of Clever Hans type.

3.2. Model understanding and hypothesis testing

By closely inspecting the clusters identified using FDA, we can formulate hypotheses about the model’s decision strategies. We can recognize groupings of complicated shapes, invariant of scale, location or translation on clusters found with Gromov-Wasserstein distance at the base of SpRAY. Examples for two distilled clusters from classes “ring-neck snake” — where the snake’s head and its brightly colored neck appear to be the relevant features — and “great grey owl” — where the patterns highlighting the face (eyes and beak) and shape of the head seem to be the common denominator — can be seen in Figure 5. However, despite the favorable invariance properties, deducting distinct hypotheses for these strategies turns out to be a nontrivial task, since clusters are semantically much harder to interpret compared to groupings found with a euclidean distance at the root. As expected, euclidean distance-based results exhibit tight groupings of attribution maps with shared regional concentrations of attribution scores, which is rather intuitive to interpret, even without much domain knowledge. We thus continue with a euclidean-based analysis.

We have discovered various interesting Clever-Hans’ moments in ImageNet. In the following we will go into detail for four examples of such strategies, which are depicted in Figure 6. Two significant clusters can be found in the class “stole”: One which contains samples where all four corners are digitally rounded, ap-

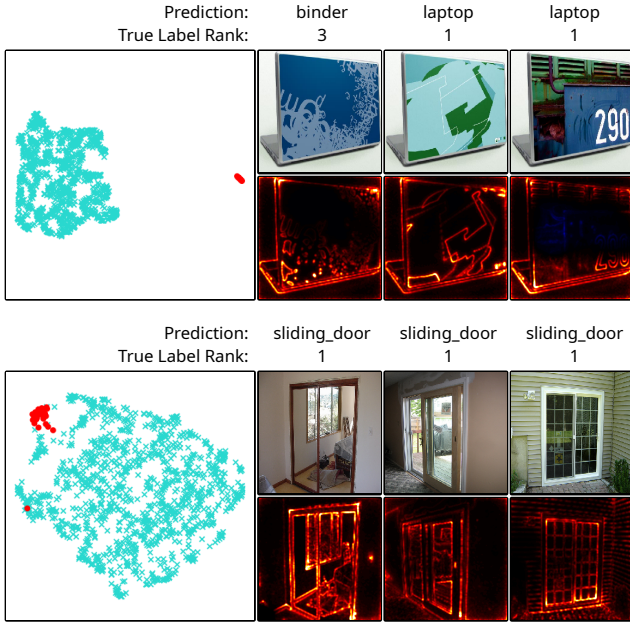


Figure 4. *Top*: The UMAP (left) shows a highly separated clusters for class “laptop”. *Bottom*: For class “sliding door”, the UMAP (left) shows no sign of separability. Red dots in the UMAPs identify the clusters the samples to the right have been grouped into. Relevancy maps to the right are color coded to identify relevant image regions *supporting* the classifier decision in hot colors (red to yellow) irrelevant regions in black color and relevant regions *contradicting* the final prediction in cold (blue to cyan) hues. The text above the sample images shows the classifier’s top-1 predicted class, and the prediction rank of the true label.

parently created by the same author (photographer). Heatmaps for those samples show very high relevance for this rounded corner property (pointed at by a red marker). It is worth to note that for the leftmost example shown in the figure, which was in contrast to all others incorrectly classified, the bottom-left corner of the image is not deemed relevant by the model (pointed at by a blue marker). The second, densely packed cluster for class “stole” shows a wooden mannequin “head” consistently used by the model for predicting the true class. Another intriguing cluster is found in the class “garbage truck”. In the identified cluster, we can see a significant number of images with a characteristic watermark placed in the bottom left corner, which is consistently picked up by the model as a relevant feature. The class “stethoscope”, which received the second highest τ score, shows padding added to the top and bottom of the image, which has not been introduced during a model-specific preprocessing step, but rather are part of the images themselves. The class “jigsaw puzzle” shows a presumably digitally pasted identical patterns on top several the source images. At

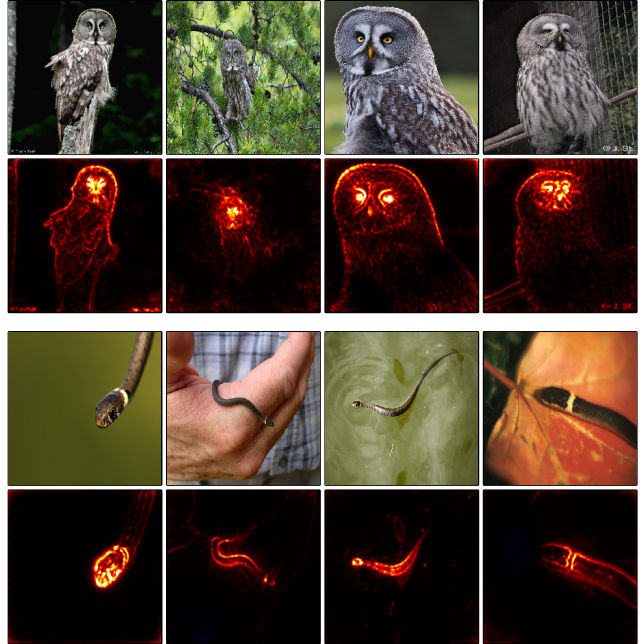


Figure 5. *Top*: Significant Gromov-Wasserstein distance based SpRAY cluster of class “great grey owl” with the corresponding attribution maps below the samples. *Bottom*: Significant clusters for class “ring-neck snake”.

least three distinct variants of this puzzle pattern have been discovered using SpRAY. The consistency of this pattern across multiple samples allows the model to overfit to this data artefact. Finally, a series of samples from class “mountain bike” shows a near-identical gray border padding picked up by the model.

For each of these observations, we can formulate the hypothesis that the model is biased on the via heatmaps highlighted properties towards their respective class. To investigate whether this is true, we construct an ablation study by isolating the artifact source and adding it to samples of other classes. That is, for class “stole” we create a digital mask of rounded corners from the affected original images and extract one of the shown wooden “mannequin heads” as a freely placable image layer. For class “mountain bike” we replicate the gray border and for “jigsaw puzzle” we extract all three discovered digital puzzle patterns. If the model then shifts its decision from the the ground truth label of the samples from the “other classes” towards the class of the artifact, we can safely deduce that the model is biased with respect to this artifact. We summarized the results for a quantitative verification of selected hypotheses in Table 3, with mean prediction rank difference $\mu(\Delta(\text{rk}))$ and mean prediction difference $\mu(\Delta f(x))$. A significant increase in prediction rank is clearly visible towards the shown

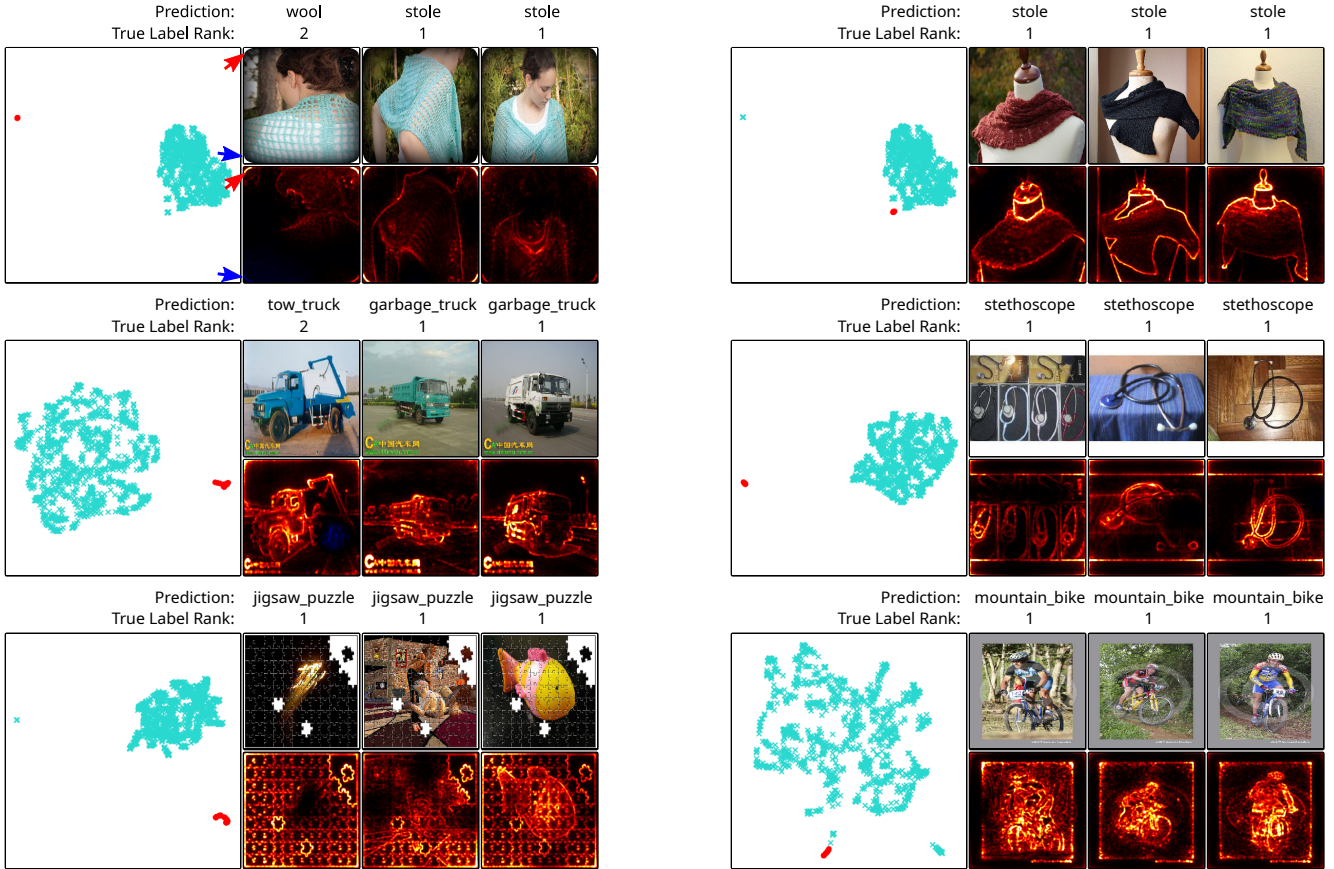


Figure 6. UMAP with samples and heatmaps of significant clusters for classes classes “stole“ (top), “garbage truck” and “stethoscope” (mid) and “jigsaw puzzle” and “mountain bike” (bottom). All significant clusters are highly separated from the rest of the samples. For each class, some images and their respective attributions from the identified cluster are shown.

artefact classes, except for class “mountain bike”. A possible explanation of this deviation of class “mountain bike” from the expected trend is revealed upon closer inspection of additional data artefacts. Several other classes, such as (e.g. classes “expresso maker” and “guillotine”), show similar border padding effects in some of the images. Furthermore, we expect to also see increased relevance on the artifact in the attribution for the modified image. Examples demonstrating isolated artifacts added onto different class samples, and the models’ reaction to the artefact when computing heatmaps for the artefact’s class of origin, are visualized in Figure 7.

By *adding* a discovered data artifact to samples of other classes we are able to quantify its importance to the detection of the labelled concept. By *removing* an artefact, we can estimate to what degree a model has learned to (solely) base its decision on the artefactual feature. If the model reacts strongly to the removal of the artefactual image feature, it has (with high probability) resorted to the artefact as a main

source of information for the respective target class. If the model does only show a weak reaction or none at all, it may have learned (several) backup strategies for detecting the concept of the target label. We measure the model’s sensitivity to the artefacts discussed in this section, by using digital inpainting techniques on the affected samples in the validation set. Table 4 compiles measurements $\mu(\Delta(\text{rk}))$ and $\mu(\Delta f(x))$ for artefact removals on classes “stole”, “jigsaw puzzle” and “mountain bike”. While the prediction for class “mountain bike” is almost completely unaffected again, and the classifier seems to have developed backup plans for predicting class “stole” in the absence of rounded image corners and wooden mannequin heads, the “jigsaw puzzle” classifier catastrophically fails in two out of three cases when a discovered digitally pasted jigsaw puzzle pattern is removed from the affected samples. The model has thus, for the class “jigsaw puzzle”, strongly overfitted to the discovered dataset bias.

As an additional interesting observation, we have also found classes with examples in the validation set

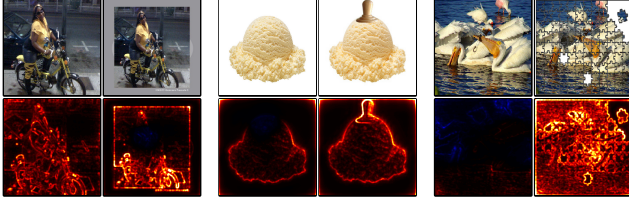


Figure 7. Addition of discovered artefacts to samples of other classes. Relevance maps are computed w.r.t. the class of origin of the artefact. *Left*: Addition of a border which transforms a “moped” into “mountain bike”. *Mid*: The addition of the “mannequin head” increases the classifier output for class “stole”. Note how the model interprets the lack of a “mannequin head” on top of the ball of ice cream in the left heatmap as contradictory feature. *Right*: Further note how the model considers the white color in the image corners as features for “stole”. Adding a digital puzzle pattern to any image forces a high probability “jigsaw puzzle” prediction.

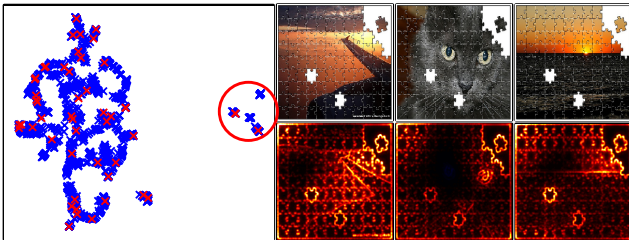


Figure 8. *Left*: UMAP of Spectral Embedding on union of training (red) and validation set (blue) for class “jigsaw puzzle”. *Right*: Images of the validation set in the previously identified “jigsaw puzzle” bias (top) with attributions (bottom).

that show the same type of artifacts as used in some of the discovered Clever Hans prediction strategies (e.g. see Figure 8), putting the model’s performance on the validation set for any of the affected classes in question.

3.3. Un-Hans’ing the model by fixing the data

A range of the prediction biases of Clever Hans type identified within our study exhibits variations of highly systematic artifactual patterns. The yellow watermark used to predict class “garbage truck”, e.g., always occurs in the bottom left corner of the affected images in Figure 6 and Figures 9 and 10. It should thus now be possible to use this understanding of the biased decision strategy for the purpose of model rectification, with the intent of *making the model forget* its association of the yellow watermark artefact to class “garbage truck”.

For this purpose, we design an experiment on a reduced set of ImageNet classes, including the affected class “garbage truck”, as well as four other randomly chosen classes. This subset (henceforth called subset

A) consists of $1,300 \cdot 5$ training samples and $50 \cdot 5$ validation samples and shall serve as a baseline. We then create a copy B of subset A , and isolate the watermark (see Figure 9 *Mid*) supposedly causing the prediction bias in class “garbage truck” from the affected training images. Within the training partition of the subset B , the isolated watermark is then added to the bottom left corner of *all* (yet unaffected) training samples in order to disable the watermark as a source of information which can be associated to one specific class, i.e. “garbage truck”. Starting from the original (pre-trained) VGG-16 model, we fine-tune two neural networks, one on the sets A and B each. We compare the prediction performance on the unaltered validation partition of subset A , as well as the model’s behavior via LRP heatmaps. First of all, it is to be expected that both fine-tuned models perform well on this reduced (and thus simpler) problem set. However, it is interesting to note that the model trained on subset B significantly outperforms model A with 99.8% vs. 98.2% top-1 accuracy on the validation set of the ImageNet subset A fine-tuning.

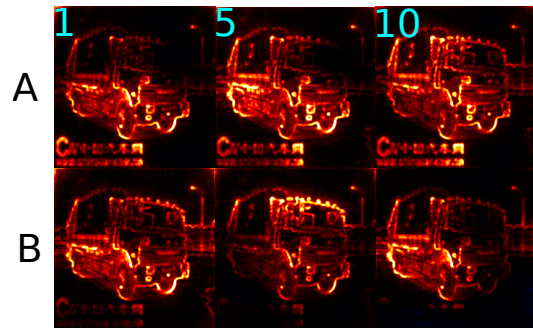


Figure 9. *Top*: Input sample with heatmap from the pre-trained model. *Mid*: The isolated artifactual watermark. *Bottom*: Heatmaps after 1, 5 and 10 epochs (numbers in cyan color) of training w.r.t. to ImageNet subsets A and B respectively.

Observing the relevance attribution maps for both

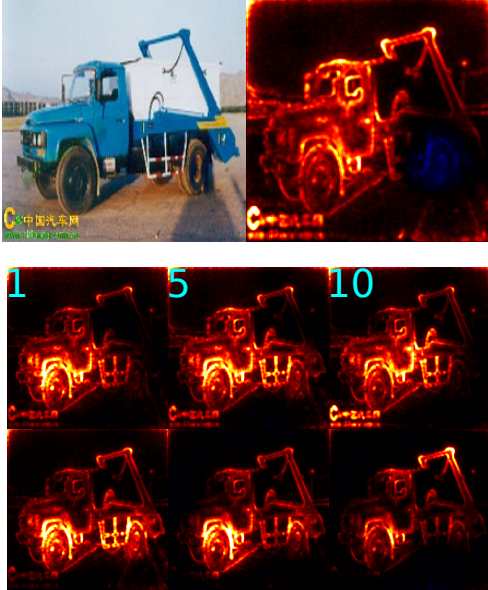


Figure 10. *Top*: Input sample with heatmap from the pre-trained model. *Bottom*: Heatmaps after 1, 5 and 10 epochs (numbers in cyan color) of training w.r.t. to ImageNet subsets *A* and *B* respectively.

fine-tuned models shown in Figures 9 and Figures 10 reveals that the model trained on subset variant *B* has consistently ceased to rely on the watermark for decision making after only 10 epochs, (already starting from the first epoch), while gradually shifting the base of its decision-making to the storage container of the vehicle. The model trained on *A*, however, retains the copyright watermark as part of its prediction strategy. The absence of the watermark within the unaltered validation set *A* explains the lower performance of model *A*, which obviously relies on the watermark to “generalize” on unseen data. We investigate this finding by exposing both fine-tuned models to a validation set which has seen the treatment of ImageNet subset *B*, *i.e.*, where all validation samples show an artificially added copyright tag. In this case, the accuracy of the model trained on *A* drops to 96.4%, while the model trained on *B* retains its original performance, further verifying that model *A* still associates the watermark to a higher degree to class “garbage truck”, while model *B* has largely disassociated class and artefactual feature. The reported accuracy ratings are summarized in Table 2. We have thus – albeit in a toy scenario – for the first time successfully *unlearned* a Clever-Hans strategy discovered with attributions and SpRAY, resulting in an improved model *B* above its baseline *A*. This marks an excellent starting point to un-Hans data corpora.

	validation on <i>A</i>	validation on <i>B</i>
training on <i>A</i>	98.2%	96.4%
training on <i>B</i>	99.8%	99.8%

Table 2. Prediction performance of the VGG-16 models fine-tuned (and evaluated) on the ImageNet subsets *A* and *B* for the “garbage truck” un-Hans’ing experiment.

class	bias	samples	$\mu(\Delta(\text{rk}))$	$\mu(\Delta f(x))$
stole	rounded corners	2000	58.14	0.0004
stole	mannequin “head”	10	106.10	0.0081
jigsaw puzzle	jigsaw pattern 1	2000	220.98	0.0160
jigsaw puzzle	jigsaw pattern 2	2000	355.60	0.8415
jigsaw puzzle	jigsaw pattern 3	2000	356.42	0.9540
mountain bike	watermark	2000	-101.02	0.0001

Table 3. The effect of *adding* a class-related artefact to samples of other classes, towards the prediction of the artefact’s class of origin. For all artefacts except the freely placable “mannequin head”, we randomly selected 2000 samples from other classes and measured the effect of the artefact addition. The $\mu(\Delta(\text{rk}))$ measures the mean change in prediction *rank* due to the artefact addition and $\mu(\Delta f(x))$ measure the mean change in the artefact’s *class probability*. High(er) values mean that the model is (strongly) affected by the artefact in its decision for the artefact’s class of origin.

class	bias	samples	$\mu(\Delta(\text{rk}))$	$\mu(\Delta f(x))$
stole	rounded corners	10	-0.70	-0.1756
stole	mannequin “head”	13	-0.62	-0.3713
jigsaw puzzle	jigsaw pattern 1	44	-0.11	-0.0146
jigsaw puzzle	jigsaw pattern 2	44	-112.52	-0.9160
jigsaw puzzle	jigsaw pattern 3	44	-208.41	-0.9305
mountain bike	watermark	17	0.00	0.0206

Table 4. The effect of *removing* a class-related artefact from image samples, towards the prediction of the artefact’s class of origin. The $\mu(\Delta(\text{rk}))$ measures the mean change in prediction *rank* due to the artefact addition and $\mu(\Delta f(x))$ measure the mean change in the artefact’s *class probability*. Low(er) values mean that the model is (strongly) affected by the artefact removal in its decision for the artefact’s class of origin.

4. Conclusion

Deep Learning models have gained high practical usability by pre-training on large corpora and then reusing the learned representation for transferring to novel related data. A prerequisite for this practice is the availability of large corpora of rather standardized and, most importantly, representative data. If artifacts or biases are present in data corpora, then the representations formed are prone to inherit these flaws. This is clearly to be avoided, however, it requires either clean data or detection and subsequent removal of artifacts,

biases *etc.* of data bases that would cause dysfunctional representation learning.

In this paper we have used explanation methods (LRP attributions [10], for an overview see [8, 41]) and specifically extended SpRAY, a technique that has been successfully used to unmask Clever Hans behavior, for automatically and scalably detecting subtle and less subtle flaws in the ImageNet corpus. One strand of our analysis was devoted to the question of how to properly reflect scaling, translations and rotations when comparing attributions in the clustering step of SpRAY. Here we found the Wasserstein distance to be a versatile candidate for achieving a metric encompassing the mentioned crucial invariances. Furthermore, our comprehensive qualitative and quantitative analysis based on the scalable technique proposed above reveals for different classes in ImageNet rather unexpected Clever Hans-type strategies [7] of the popular VGG-16 deep learning model (to which also other architectures are sensitive to; see Appendix). These are caused by a zoo of artifacts and biases isolated by our framework in the corpus; these encompass: copyright tags, unusual image formatting, specific co-occurrences of unrelated objects, cropping artifacts, just to name a few. Detecting this zoo gives not only insight but also the possibility for relieving ImageNet from its Clever Hans moments, *i.e.* we are now able to un-Hans the ImageNet corpus and provide a more consistent basis for pretrained models. We demonstrated this in an *unlearning experiment* for class “garbage truck” (see above and Figures 9 and 10), and further ImageNet classes in the Appendix. Note that without removing such data artifacts, learning models are prone to adopt Clever Hans strategies [7], thus, giving the correct prediction for an artifactual/wrong reason. This makes them especially vulnerable to adversarial attacks that can harvest all such artifactual issues in a data corpus [42]. Future work will therefore focus on the important intersection between security and functional cleaning of data corpora, *e.g.*, to lower the attack risk when building on top of pretrained models.

Acknowledgement

This work was supported by the Brain Korea 21 Plus Program through the National Research Foundation of Korea; the Institute for Information & Communications Technology Promotion (IITP) grant funded by the Korea government [No. 2017-0-00451]; the Deutsche Forschungsgemeinschaft (DFG) [grant Math+, EXC 2046/1, Project ID 390685689]; and the German Ministry for Education and Research (BMBF) as Berlin Big Data Center (BBDC) [01IS14013A], Berlin Center for Machine Learning

(BZML) [01IS18037I] and TraMeExCo [01IS18056A].

References

- [1] Y. LeCun, Y. Bengio, and G. E. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [3] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, *et al.*, “Mastering the game of go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [4] K. T. Schütt, F. Arbabzadah, S. Chmiela, K.-R. Müller, and A. Tkatchenko, “Quantum-chemical insights from deep tensor neural networks,” *Nature Communications*, vol. 8, p. 13890, 2017.
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems (NIPS)*, pp. 1097–1105, 2012.
- [6] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, *et al.*, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [7] 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, p. 1096, 2019.
- [8] W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K.-R. Müller (Eds.), “Explainable AI: Interpreting, explaining and visualizing deep learning,” *Springer LNCS 11700*, 2019.
- [9] S. M. Lundberg, G. G. Erion, H. Chen, A. De-Grave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S. Lee, “Explainable AI for trees: From local explanations to global understanding,” *CoRR*, vol. abs/1905.04610, 2019.
- [10] 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, p. e0130140, 2015.
- [11] 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 *Proc. of IEEE International Conference on Computer Vision (ICCV)*, pp. 618–626, 2017.
- [12] M. Sundararajan, A. Taly, and Q. Yan, “Axiomatic attribution for deep networks,” in *Proc. International Conference on Machine Learning (ICML)*, pp. 3319–3328, JMLR.org, 2017.
- [13] A. Shrikumar, P. Greenside, and A. Kundaje, “Learning important features through propagating activation differences,” in *Proc. of International Conference on Machine Learning (ICML)*, pp. 3145–3153, 2017.
- [14] M. T. Ribeiro, S. Singh, and C. Guestrin, “why should I trust you?: Explaining the predictions of any classifier,” in *Proc. of ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, pp. 1135–1144, 2016.
- [15] L. M. Zintgraf, T. S. Cohen, T. Adel, and M. Welling, “Visualizing deep neural network decisions: Prediction difference analysis,” in *Proc. of International Conference on Learning Representations (ICLR)*, 2017.
- [16] R. C. Fong and A. Vedaldi, “Interpretable explanations of black boxes by meaningful perturbation,” in *Proc. of IEEE International Conference on Computer Vision (ICCV)*, pp. 3449–3457, 2017.
- [17] I. Guyon and A. Elisseeff, “An introduction to variable and feature selection,” *Journal of machine learning research*, vol. 3, no. Mar, pp. 1157–1182, 2003.
- [18] B. Kim, M. Wattenberg, J. Gilmer, C. J. Cai, J. Wexler, F. B. Viégas, and R. Sayres, “Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV),” in *Proc. of International Conference on Machine Learning (ICML)*, pp. 2673–2682, 2018.
- [19] R. Rajalingham, E. B. Issa, P. Bashivan, K. Kar, K. Schmidt, and J. J. DiCarlo, “Large-scale, high-resolution comparison of the core visual object recognition behavior of humans, monkeys, and state-of-the-art deep artificial neural networks,” *Journal of Neuroscience*, vol. 38, no. 33, pp. 7255–7269, 2018.
- [20] M. Everingham, L. Gool, C. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes challenge results,” URL: http://host.robots.ox.ac.uk/pascal/VOC/voc2007/workshop/everingham_cls.pdf, 2007.
- [21] M. Meila and J. Shi, “A random walks view of spectral segmentation,” in *Proceedings of the Eighth International Workshop on Artificial Intelligence and Statistics, AISTATS 2001, Key West, Florida, US, January 4-7, 2001*, 2001.
- [22] A. Y. Ng, M. I. Jordan, and Y. Weiss, “On spectral clustering: Analysis and an algorithm,” in *Advances in Neural Information Processing Systems*, pp. 849–856, 2002.
- [23] 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 LNCS 11700, 2019.
- [24] M. Kohlbrenner, A. Bauer, S. Nakajima, A. Binder, W. Samek, and S. Lapuschkin, “Towards best practice in explaining neural network decisions with LRP,” *CoRR*, vol. abs/1910.09840, 2019.
- [25] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *CoRR*, vol. abs/1409.1556, 2014.
- [26] M. Alber, S. Lapuschkin, P. Seegerer, M. Hägele, K. T. Schütt, G. Montavon, W. Samek, K.-R. Müller, S. Dähne, and P.-J. Kindermans, “investigate neural networks!,” *Journal of Machine Learning Research*, vol. 20, pp. 93:1–93:8, 2019.
- [27] F. Perronin, J. Sánchez, and T. Mensink, “Improving the fisher kernel for large-scale image classification,” in *Proc. of European Conference on Computer (ECCV)*, pp. 143–156, 2010.
- [28] G. Peyré, M. Cuturi, and J. Solomon, “Gromov-wasserstein averaging of kernel and distance matrices,” in *Proc. of International Conference on Machine Learning (ICML)*, pp. 2664–2672, 2016.
- [29] N. S. Altman, “An introduction to kernel and nearest-neighbor nonparametric regression,” *The American Statistician*, vol. 46, no. 3, pp. 175–185, 1992.

- [30] U. von Luxburg, “A tutorial on spectral clustering,” *Statistics and Computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [31] E. Jones, T. Oliphant, P. Peterson, *et al.*, “SciPy: Open source scientific tools for Python,” 2001–. [Online; accessed].
- [32] C. Lanczos, *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*. United States Governm. Press Office Los Angeles, CA, 1950.
- [33] Y. Rubner, C. Tomasi, and L. J. Guibas, “A metric for distributions with applications to image databases,” in *Proceedings of the Sixth International Conference on Computer Vision (ICCV-98), Bombay, India, January 4-7, 1998*, pp. 59–66, 1998.
- [34] J. Solomon, F. de Goes, G. Peyré, M. Cuturi, A. Butscher, A. Nguyen, T. Du, and L. J. Guibas, “Convolutional wasserstein distances: efficient optimal transportation on geometric domains,” *ACM Trans. Graph.*, vol. 34, no. 4, pp. 66:1–66:11, 2015.
- [35] Y. LeCun, “The mnist database of handwritten digits.” <http://yann.lecun.com/exdb/mnist/>, 1998.
- [36] M. Cuturi and A. Doucet, “Fast computation of Wasserstein barycenters,” in *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 685–693, 2014.
- [37] R. A. Fisher, “The use of multiple measurements in taxonomic problems,” *Annals of eugenics*, vol. 7, no. 2, pp. 179–188, 1936.
- [38] K. Fukunaga, “Chapter 1 - introduction,” in *Introduction to statistical pattern recognition*, Boston: Academic Press Professional, Inc., 1990.
- [39] L. v. d. Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of Machine Learning Research*, vol. 9, no. Nov, pp. 2579–2605, 2008.
- [40] L. McInnes and J. Healy, “UMAP: uniform manifold approximation and projection for dimension reduction,” *CoRR*, vol. abs/1802.03426, 2018.
- [41] G. Montavon, W. Samek, and K.-R. Müller, “Methods for interpreting and understanding deep neural networks,” *Digital Signal Processing*, vol. 73, pp. 1–15, 2018.
- [42] N. Carlini and D. Wagner, “Towards evaluating the robustness of neural networks,” in *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57, IEEE, 2017.
- [43] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 2261–2269, 2017.

Appendix

We complement our findings from our manuscript with additional data within this appendix, by (1) demonstrating how the prediction artefacts discovered for the VGG-16 classifier also reappear on other DNN models and (2) showing the effects of the un-Hans’ing experiment, which has been described in detail in Section 3.3, when performed on other classes and artifacts.

4.1. ImageNet Artefacts Affecting Additional DNNs

In Section 3.2 we describe a series of systematic prediction biases discovered using the SpRAY technique for several affected classes. In all these cases, the downloaded VGG-16 model has overfit on input features which are characteristic for certain object classes in context of the ImageNet dataset. We thus assume that other neural network architectures sharing the same data source for training may also share certain Clever Hans strategies with the investigated VGG-16 classifier.

Figure 11 exemplarily shows LRP heatmaps computed for the VGG-19 [25] and the DenseNet-121 [43] model — which have also been downloaded as pre-trained predictors optimized on the ImageNet data corpus — for samples which exhibit data artefacts as discovered for the VGG-16 model. We notice that both the architecturally very similar VGG-16 and VGG-19 architectures heatmaps are very concentrated on shape features such as edges and color-gradient rich image areas. The heatmaps computed for the DenseNet-121 model on the other hand are much more focused on class- and object-specific textures and colors. For all investigated samples, we however notice that all three models tend to use the same w.r.t. to the true class semantically unrelated yet correlated features for prediction. That is, for class “carton”, all three models support their predictions with a set of barely visible and centered watermark consisting of asian characters for prediction, as well as a second orange and small watermark appearing in the bottom right corner of “carton” images with high frequency. Similarly for classes “garbage truck”, “jigsaw puzzle” and “stole” shown in Figure 11 all three models support their prediction based on the discovered yellow watermark, the cut-outs of the digitally added puzzle pattern, the rounded image corners and the wooden mannequin head.

Considering the systematicity of use of these data artefacts by all three models, we strongly recommend a thorough categorization of Clever Hans behavior of machine learning models and their data sources essential components of future dataset creation efforts.

4.2. Additional Cases of Un-Hans’ing

In Section 3.3, we described the un-Hans’ing procedure with the intent to force the neural network model to *forget* learned yet unintended feature-to-class associations in detail for class “garbage truck” affected a yellow watermark artefact. Here, we repeat the experiment of Section 3.3 for artefacts discovered for classes “stole” and “jigsaw puzzle”.

Figure 12 demonstrates a setting highly similar to the one for class “garbage truck” discussed in Section 3.3: The discovered data artefact — here a digitally rounded image corners with white background — exhibits extremely high regional consistency and only covers very limited parts of the image area. Once the isolated corner feature has been added to *all* samples during our experiment, the model quickly has disassociated the artefact from the label “stole”. After continued re-training, LRP begins to attribute negative relevance to rounded image corners, indicating that the process of un-Hansing went beyond mere forgetting by creating a negative association between corner artefact and class label.

The second data artefact discovered for class “scole” is a frequently shown wooden “mannequin head” co-appearing with the woven stoles themselves. Since here, the expression of the artefact was much more diverse in pose and position and has shown almost no regional consistency, we manually isolated a (very) limited amount of prototypical “mannequin heads” from the data and randomly (within reason) added wooden stump as an image element to each sample of each batch during re-training. Figure 13 shows the progression of un-Hansing at hand of two different input sample. While for the sample shown at the top of the figure the model has not disassociated between this particular expression of the “mannequin head” feature (at times, the feature’s accumulated positive relevance even increased), the model has ceased to support its prediction for class “stole” with the artefactual feature for the bottom image.

Lastly, we investigate the “digital jigsaw puzzle pattern” artefact discovered for class “jigsaw puzzle”, which appears in multiple variants. Each variant, however, is expressed with almost complete and pixel-identical consistency. We therefore select one variant of the artefact and add it as a mask to *all* training samples of the un-Hansing training subset B extracted from the ImageNet corpus. Here again, we can observe that the model *forgets* the association between this particular pattern and the class label “jigsaw pattern”: In Figure 14, positive relevance completely disappears from the digital jigsaw pattern during un-Hansing, such that the feature is not used anymore for predicting “jigsaw”.

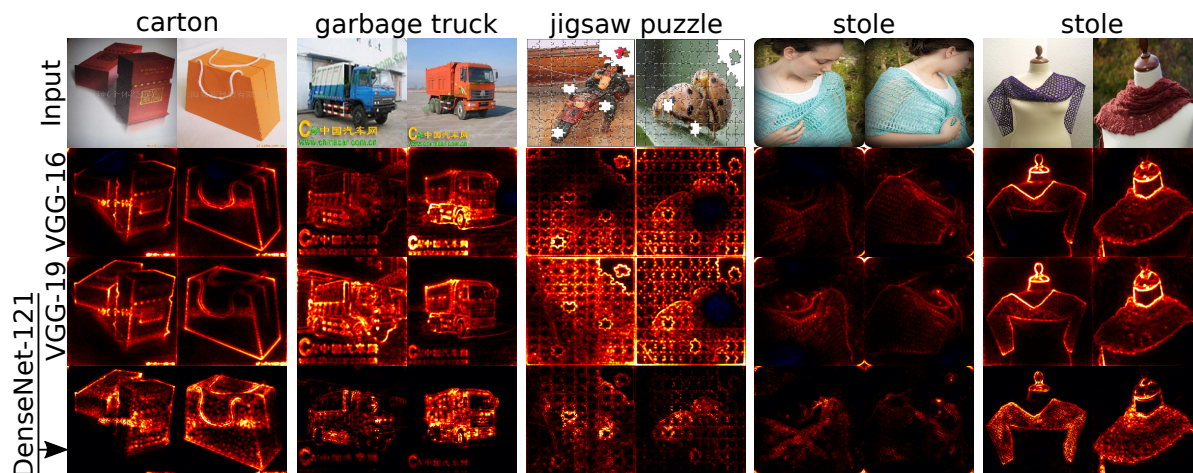


Figure 11. Heatmaps for classes and samples with on ImageNet and VGG-16 discovered data and prediction artefacts for the VGG-19 and DenseNet-121 models. *Top to bottom*: Input samples, heatmaps for VGG-16, VGG-19 and DenseNet-121. *Left to right*: Columns show (in pairs) artefacts for classes “carton”, “garbage truck”, “jigsaw puzzle”, “stole” (rounded corners) and “stole” (prediction supported by mannequin head), which have been discovered from a VGG-16 classifier, but apply to all three models.

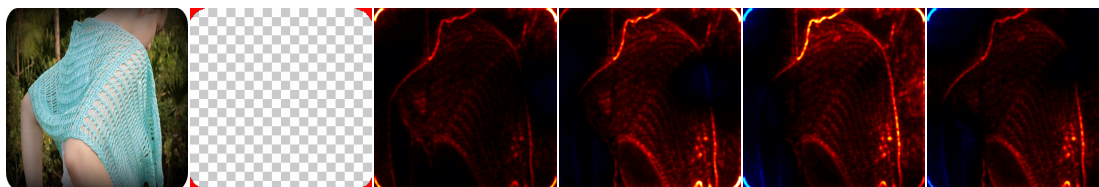


Figure 12. Un-Hans’ing Experiment for class “stole” and the “rounded corners” artefact. *Left to right*: Example input, the artefact (with transparent background, and the white corner pattern here shown in read for visibility reasons), heatmap expressions computed during the un-Hans’ing process.

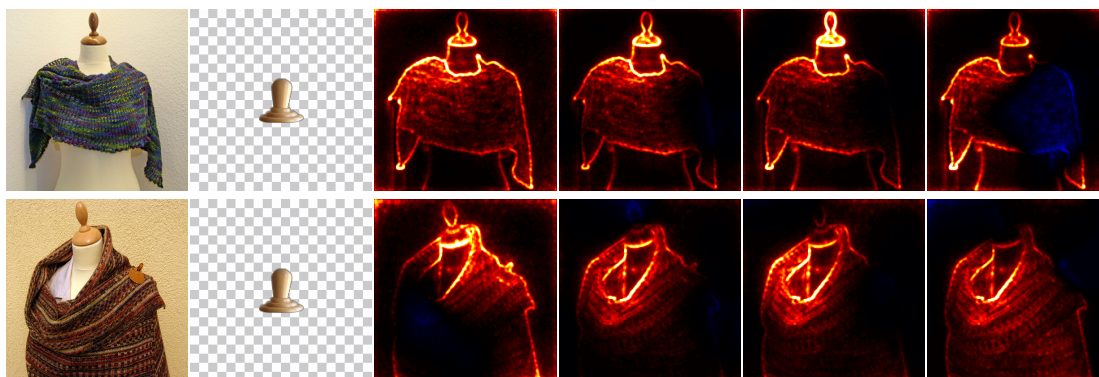


Figure 13. Un-Hans’ing Experiment for class “stole” and the “mannequin head” artefact. *Left to right*: Example inputs, the artefact (a manually isolated wooden mannequin head), heatmap expressions computed during the un-Hans’ing process.

What prevails, however, is a strongly negative relevance map on the fornicating ladybug pair of ladybugs, indicating the model’s reasoning that the insects’ presence speaks *against* class “jigsaw” (and rather for a competing network output). The effect of forgetting this consistently expressed yet very large artefactual feature has an understandably catastrophic effect to the model’s capability to predict the original ImageNet

label for affected samples (*c.f.* Table 4).

We conclude that while according to our experiments *precisely applying brain damage* to a pre-trained neural network model is definitely possible, its execution may in some cases — at least while through manipulation of the training data in pixel space — be non-trivial.

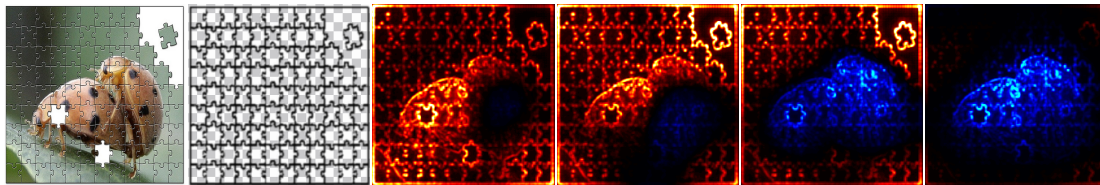


Figure 14. Un-Hans'ing Experiment for class “jigsaw puzzle” and the “digital puzzle pattern” artefact. *Left to right:* Example input, the artefact (a manually isolated wooden mannequin head), heatmap expressions computed during the un-Hans'ing process.