



This Looks Like That, Because ... Explaining Prototypes for Interpretable Image Recognition

Meike Nauta^{1,2}(✉), Annemarie Jutte¹, Jesper Provoost¹, and Christin Seifert^{1,2}

¹ University of Twente, Enschede, The Netherlands
m.nauta@utwente.nl, {a.m.p.jutte,j.c.provoost}@student.utwente.nl
² University of Duisburg-Essen, Essen, Germany
christin.seifert@uni-due.de

Abstract. Image recognition with prototypes is considered an interpretable alternative for black box deep learning models. Classification depends on the extent to which a test image “looks like” a prototype. However, perceptual similarity for humans can be different from the similarity learned by the classification model. Hence, only visualising prototypes can be insufficient for a user to understand what a prototype exactly represents, and why the model considers a prototype and an image to be similar. We address this ambiguity and argue that prototypes should be explained. We improve interpretability by automatically enhancing visual prototypes with quantitative information about visual characteristics deemed important by the classification model. Specifically, our method clarifies the meaning of a prototype by quantifying the influence of colour hue, shape, texture, contrast and saturation and can generate both global and local explanations. Because of the generality of our approach, it can improve the interpretability of any similarity-based method for prototypical image recognition. In our experiments, we apply our method to the existing Prototypical Part Network (ProtoPNet). Our analysis confirms that the global explanations are generalisable, and often correspond to the visually perceptible properties of a prototype. Our explanations are especially relevant for prototypes which might have been interpreted incorrectly otherwise. By explaining such ‘misleading’ prototypes, we improve the interpretability and simulatability of a prototype-based classification model. We also use our method to check whether visually similar prototypes have similar explanations, and are able to discover redundancy. Code is available at https://github.com/M-Nauta/Explaining_Prototypes.

A. Jutte and J. Provoost—Both authors contributed equally to this work.

© Springer Nature Switzerland AG 2021
M. Kamp et al. (Eds.): ECML PKDD 2021 Workshops, CCIS 1524, pp. 441–456, 2021.
https://doi.org/10.1007/978-3-030-93736-2_34

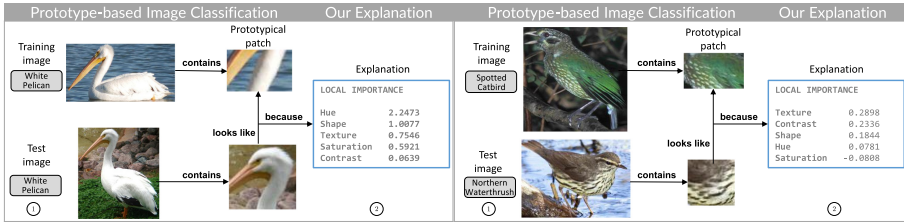


Fig. 1. ① Prototype-based image classification (e.g. ProtoPNet [7]). ② Our contribution: Quantifying the importance of visual characteristics to explain why the classification model deemed an image patch and prototype similar. Left: Logical explanation for the clear similarity between the patches. Right: a ‘misleading’ prototype: humans might expect these patches to be dissimilar, but our method explains that the classification model considers these patches alike because of similar texture.

1 Introduction

Convolutional Neural Networks (CNNs) are the de-facto standard for object detection due to their impressive performance in numerous automated image classification tasks. However, the black box nature of neural networks prevents a human to assess the model’s decision making process, which is especially problematic in domains with high stakes decisions [21]. Following this demand on understanding automated decision making, explainable Artificial Intelligence (XAI) has been actively researched [10]. *Post-hoc* explanation methods learn a second, transparent model to approximate the first black box model [10], but these reverse-engineering approaches are not guaranteed to show the *actual* reasoning of the black box model [21]. *Intrinsically* interpretable models on the other hand, are faithful by design and allow *simulatability*: a user should be able to reproduce the model’s decision making process based on the input data together with the explanations of the interpretable model and come to the same prediction [15]. One type of such models is prototypical learning, which has a transparent, built-in case-based decision making process. We focus on the problem of supervised image recognition where a machine learning model should label an image. Prototypes in this context are usually ‘nearest neighbours’, i.e., images from the training set that look similar to the image being classified [1, 4]. The similarity between a prototype and an image is often measured in latent space, learned by the neural network, where images from the same class are close and dissimilar images are far apart with respect to a certain distance or similarity metric. Recently, the Prototypical Part Network (ProtoPNet) [7] and ProtoTree [17] were introduced which use prototypical *parts* and identify similar patches in an image. The classification depends on the extent to which *this* part of the image “looks like” *that* prototypical part, measured by a similarity score. An example of this reasoning is shown in Fig. 1.

Prototype Ambiguity. In this paper, we address the ambiguity that prototypes can have and present a method to *explain prototypes*. Consider the left part in Fig. 1, showing a prototypical patch (‘prototype’) of a white pelican. Although the similarity between this prototype and the patch in the test image

is not surprising, it is unclear what this prototype exactly represents. Is the prototype looking for a white neck, an orange-coloured beak or is the shape of the beak specifically important? The similarity score between the two patches assigned by the model depends on its classification strategy, and hence its learned latent space. Explanations are especially needed when similarity is not so obvious. When seeing the two patches in the right part of Fig. 1, a human might argue that these patches are dissimilar because of the colour differences. The classification model however assigns these patches a high similarity score, and thus considers them alike, even though the test image is from a different class than the prototype. This shows that a human and the CNN might have different reasoning processes, despite using the same prototypes. The classification strategy of a neural network, dependent on the learned latent space, determines the reason for considering two patches as being similar or different. It has been shown that CNNs trained on ImageNet are strongly biased towards recognizing texture [8], although other work shows that CNNs can be biased towards shape [18] or colour [11]. Perceptual similarity for humans however is biased towards shape [3], but also based on e.g. colour, size, semantic similarity and complexity [13, 20]. It is also questionable whether humans and CNNs will ever follow the exact same similarity reasoning, since Rosenfeld et al. found that neural networks fall short on predicting human similarity perception [19]. Since a user is not aware of the underlying classification strategy of the trained CNN and might also be unaware of personal biases, only visualising prototypes is insufficient for understanding what a prototype exactly represents, and why a prototype and an image are considered similar. This issue may also arise with other explainability methods that show or highlight image parts, such as attention mechanisms [6], components [22] and other part-based models e.g. [26, 27]. Including our explanations can help users to increase the simulatability [15] and general understanding of the model.

Contribution. We improve the interpretability of a prototype-based CNN by automatically enhancing prototypes with extra quantitative information about visual characteristics used by the model. Specifically, we present a methodology to quantify the influence of colour hue, saturation, shape, texture, and contrast in a prototype. This clarifies what the model pays attention to and why a model considers two images to be similar. Hence, we disentangle localisation and explanation. Our method can extend any prototype-based model for image recognition, such as ProtoPNet [7] and ProtoTree [17]. In this paper, we show its applicability for the prototypical parts of ProtoPNet. For example, again considering the left part of Fig. 1, our explanation shows that ProtoPNet considers the prototype and patch from the test image to be similar because of the similar colour hue and shape of the beak in the test image. Our method is especially useful when similarity is not so obvious. It can explain potentially misleading prototypes such as the right prototype in Fig. 1. Whereas a human might look for something green, our explanation reveals that ProtoPNet considers these two patches similar because of texture, contrast and shape. The similarity is thus because of the dotted pattern and colour hue was not important. This explanation seems reasonable given that the prototype is from the class “Spotted Catbird”.

Our method automatically modifies images to change their hue, shape, texture, contrast or saturation. We forward both the original image and the modified image through the prototype-based network and analyse the resulting similarity scores. Specifically, the similarity score between the prototype and the original image is compared with the similarity score of the prototype and a modified image. The intuition is that a visual characteristic is considered *unimportant* by the classification model when the difference between these two similarity scores is small (and will therefore get a low importance score), and is deemed *important* when the similarity scores differ sufficiently. For example, a blue bird is changed to a purple bird by changing the hue of the image. If hue would have been important for the specific prototype, it would be expected that the model assigns a low similarity between the prototype and the purple bird, whereas the similarity with the blue bird was high. As shown in Fig. 1, the prototypes can subsequently be explained by quantifying the importance of visual characteristics.

2 Prototypical Part Network

We apply the methodology presented in this paper to ProtoPNet, the Prototypical Part Network from Chen et al. [7] that follows the “*this looks like that*” reasoning. Prototypical parts learned by ProtoPNet are subsequently explained by our method. Key for presenting our explanation methodology is having a global understanding of the workings of ProtoPNet.

The ProtoPNet architecture consists of a standard CNN (e.g. DenseNet), followed by a prototype layer and a fully-connected layer. The prototype layer consists of a pre-determined number of class-specific prototypes, usually 10 prototypes per class [7]. The fully-connected layer learns a weight for each prototype. During training, prototypes are vectors in latent space that should learn discriminative, prototypical parts of a class. An input image is forwarded through the CNN, after which the prototype layer compares the resulting latent embedding with the prototype. A kernel slides over the latent image and at each location, the distance between the latent prototype and a patch in the latent image is calculated. This creates an activation map, containing the distance to the prototype at each location in the latent image. To ensure that the prototype can be visualised, the training procedure of ProtoPNet requires that each prototype is *identical* to some latent training patch such that it can be upsampled to the size of the original image and visualised as an image patch (Fig. 2).

After training, ProtoPNet classifies a test image k by calculating the similarity between a prototype and image k . The distance $d_{j,k}$ between the nearest patch in latent image k to the j -th prototype is converted to a similarity score:

$$g_{j,k} = \log \left(\frac{d_{j,k} + 1}{d_{j,k} + \epsilon} \right), \quad (1)$$

where ϵ is an arbitrarily small positive quantity to prevent zero division. To classify this image, the similarity scores of the image and each prototype are weighted by the fully-connected layer and summed per class, resulting in a final score for an image belonging to each class. The left part of Fig. 4 illustrates this reasoning process.

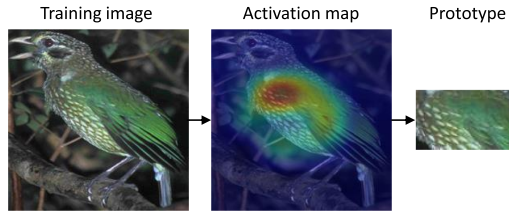


Fig. 2. A prototype from ProtoPNet is the nearest patch of a training image.

3 Methodology

In order to obtain importance scores for visual characteristics of prototypes, images are automatically modified. The characteristics in focus are contrast, colour hue and saturation, shape, and texture (cf. Sect. 3.1). Our approach for automatically modifying these characteristics is described in Sect. 3.2. Section 3.3 presents a methodology to explain prototypes by quantifying the importance of a visual characteristic.

3.1 Important Visual Characteristics

The perceptual and cognitive processing in the human visual system is influenced by various features. To determine which image modifications we need to effectively explain prototypes, we discuss important visual characteristics for the human perceptual system.

The data visualisation domain has a ranking of channels to control the appearance of so-called *marks* [16]. A ‘mark’ is a basic graphical element in an image, such as a black triangle or moving red dot. Important visual channels for marks are position, size, angle, spatial region, colour hue, colour luminance, colour saturation, curvature, motion and shape [16]. For static 2-dimensional natural images, motion is not applicable and we consider curvature related to shape. Furthermore, it is not necessary to modify the size, position, angle or spatial region of objects in images, since CNNs with pooling, possibly combined with suitable data augmentation, are invariant to these characteristics [9, 23]. Moreover, research in neuroscience shows that the human eye can recognise objects independent of ambient light level during the day [24], whereas contrast (spatial variation in luminance) is needed for edge detection and delineation of objects [16]. The human visual system is thus more sensitive to contrast than absolute luminance [24]. We therefore will not modify the absolute luminance, but the contrast in an image. Thus, the visual characteristics from the data visualisation domain that we deem important for explaining a prototype are **hue**, **contrast**, **saturation** and **shape**.

The channels for marks mentioned in the previous paragraph are however too simplistic, because they do not include the texture or material of an object. Research in neuroscience also emphasises the importance of texture for classifying objects in the natural world [5, 14]. Related to this, Bau et al. [2] disentangled

visual representations by layers in a CNN and found that self-supervised models, especially in the earlier layers of the network, learn many texture detectors. We therefore also include **texture** as an important visual characteristic.

3.2 Image Modifications

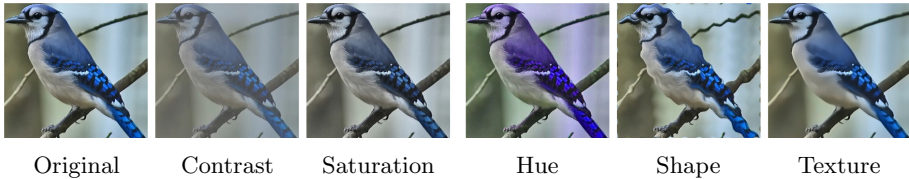


Fig. 3. Image modifications for corresponding visual characteristics.

For each of the visual characteristics, an image set is created. Each of these sets contains modified images, which are designed to be harder to classify based on the respective characteristic. For example, we generate a set of low-contrast images, such that contrast information can not (or hardly) be used by the model. Using these modified images, the importance of a characteristic for a specific prototype can be determined by comparing the differences between the prototype-image similarity of the original and modified image.

To create the modified images, we apply image transformations to reduce the intensity of each characteristic, i.e., we create images with reduced contrast, saturation, hue, shape and texture. Figure 3 shows an example image and its modified versions. We opt for automated image modifications instead of manual modifications used for experiments in psychology research (e.g. [18]), to be able to create a large number of modified images efficiently.

To create low **contrast** images, the original image is blended with the mean of its grayscale version. More concretely, we first create a grayscale version of the image and calculate its mean value. We then generate the modified image by pixel-wise averaging each channel (RGB) of the original image with the mean grayscale value. Similarly, the low-**saturation** image is created by averaging the original image with its grayscale counterpart. To generate an image with different **colour hues**, the RGB image is converted to the HSV colour space after which the H-dimension is modified for each pixel. In order to modify **shapes** in an image, we add a linear displacement by warping the image. Specifically, we shift pixels according to a sine wave in both the horizontal and vertical direction. To modify **texture**, we apply a non-local means denoising technique which removes small quantities of noise, and can therefore be used to blur the sophisticated texture of a bird while preserving its overall shape. Implementation details are presented in Sect. 4.1.

3.3 Importance Scores for Image Characteristics

We evaluate the importance of visual characteristics by calculating local and global importance scores. The **local** score measures the importance of the visual characteristics for a single image, and is therefore applied on previously unseen images, i.e., any image in the test dataset S_{test} . The **global** score measures the importance of visual characteristics for one prototype in general, and is independent of a specific input image. The global score is obtained from the training dataset S_{train} .

Let $i \in \{\text{contrast, saturation, hue, shape, texture}\}$ denote the type of modification, $j \in \{1, 2, \dots, n\}$ the prototype index and k the image. Furthermore, as introduced in Sect. 2, let the similarity of the original image and the prototype be denoted as g , and the similarity of the modified image and the prototype be denoted as \hat{g} . Then the local importance score $\phi_{\text{local}}^{i,j,k}$ of characteristic i for test image $k \in S_{\text{test}}$ on the j -th prototype is the difference in similarity scores:

$$\phi_{\text{local}}^{i,j,k} = g_{j,k} - \hat{g}_{i,j,k}. \quad (2)$$

For the calculation, we fix the patch location such that the part of image k compared with prototype j is the same for both the original and modified image.

Whereas the local importance score indicates to what extent a visual characteristic influences the similarity score given by the prototype-based model *for a single image*, the *global* importance score gives a general impression of the importance of visual characteristics in a prototype. The global score of characteristic i on the j th prototype can be calculated by taking all training images into account. A naive approach would be to average over the local scores of all training images. However, prototype j might not be present in all images and modifying those images will therefore not (or barely) influence the resulting similarity score. For example, if prototype j represents a specific beak which is not present in original image k , ProtoPNet will give a low similarity score $g_{j,k}$. Since this prototype will also be absent in the modified image, the difference between the similarity scores, Eq. 2, is near zero. This result could indicate that a certain characteristic is not important, although the result is actually indicating that the prototype was simply not present. Therefore, we create a more informative global importance score by calculating a weighted arithmetic mean by weighing the local scores of all images in S_{train} by their similarity score with the prototype:

$$\phi_{\text{global}}^{i,j} = \frac{\sum_{k=1}^{|S_{\text{train}}|} \phi_{\text{local}}^{i,j,k} \cdot g_{j,k}}{\sum_{k=1}^{|S_{\text{train}}|} g_{j,k}}. \quad (3)$$

Hence, if unmodified image k gets a low similarity score with prototype j , it will get a low weight for the global importance calculation. In contrast, if prototype j is clearly present in image k , ProtoPNet will assign a high similarity score and hence k gets a high weight.

These importance scores can be used to create *global* explanations that explain a prototype, and *local* explanations that explain the similarity score between a given image and a prototype. The global explanation for the j -th

prototype lists for each visual characteristic i its importance by showing the importance scores $\phi_{\text{global}}^{i,j}$. This explanation is thus input independent and can be created before applying the prototype model to unseen images. The local explanation is of use during testing and explains a single prediction.

4 Experimental Setup

For our experiments, we use the Caltech-UCSD Birds dataset [25], a dataset for bird species identification also used by Chen et al. [7] for training their ProtoPNet. It contains 200 different classes with approximately 60 images per class. The dataset provides a train-test split, leading to S_{train} with 5994 images and S_{test} with 5794 images. To evaluate our method for explaining prototypes, we first train a ProtoPNet [7] that results in an interpretable predictive model with prototypical parts for fine-grained image recognition. We apply the same data processing techniques as the original work [7]. We cropped the images according to the bounding boxes provided with the dataset and apply data augmentation on S_{train} as described in ProtoPNet’s supplementary material [7]. All images are resized to 224×224 . We opted for DenseNet-121 [12] as backbone of ProtoPNet, as this was reported to be the best-performing network on the Caltech-UCSD dataset [7]. The DenseNet-121 network has been pre-trained on ImageNet.¹ When forwarding the resized images through DenseNet, the input image dimensions, $H_{\text{in}} = W_{\text{in}} = 224$ and $D_{\text{in}} = 3$, are transformed to the output dimensions $H = 7$, $W = 7$ and $D = 128$. Depth D is a hyperparameter in ProtoPNet determining the number of channels for the network output and the prototypes, and is set to 128 as in ProtoPNet [7]. As in the original paper [7] we use 10 prototypes per class, leading to 2000 prototypes in total. All other training parameters are also replicated from the implementation by Chen et al. [7].

We apply our method to the resulting prototypes for generating global and local explanations. Section 4.1 presents the implementation details for our image modifications. The design of our experiments to evaluate our explanations is presented in Sect. 4.2.

4.1 Modification Implementation

When implementing the image modifications as described in Sect. 3.2, we aim for a similar modification ‘strength’ for all characteristics in order to compare importance scores. Furthermore, the image modifications should be modest, since too extreme modifications can lead to out-of-distribution images that result in erratic behaviour of the underlying neural network of ProtoPNet. A similar modification degree depends on how ProtoPNet perceives the images. Therefore, we find a suitable modification degree by forwarding both the unmodified image and

¹ We use the same methodology as ProtoPNet [7] in order to reproduce results, although it is known that there is some overlap between Caltech-UCSD and ImageNet.

Classification					Explanation	
Test image (most activated area)	Prototype from training image	Similarity score	Weight last layer	Points from this prototype	Activation map	Local Importance Scores
		4.784	x 1.145	= 5.478		Shape 1.9613 Contrast 1.2462 Hue 0.9496 Texture 0.7774 Saturation 0.0298
		4.028	x 1.173	= 4.725		Hue 1.9574 Shape 0.8927 Contrast 0.2185 Saturation 0.1970 Texture 0.0365
		3.668	x 1.178	= 4.321		Texture 0.7038 Hue 0.6413 Shape 0.5852 Contrast 0.3011 Saturation 0.0674
				⋮		
Total points White Pelican			=	29.448		

Fig. 4. Left: ProtoPNet’s reasoning with a subset of all prototypes of the White Pelican class. To classify a test image, ProtoPNet compares the class-specific prototypes of each class with the test image to calculate the total number of points for this class. An image is classified as the class with the most points. Right: The activation maps produced by ProtoPNet and our corresponding local explanations that explain which characteristics were important for a similarity score.

the modified image through the underlying CNN and compare their L1-norm distance in latent space. We tune the modifications parameters such that the mean distance between the unmodified latent training images and the latent modified images is exactly 0.0002 for all characteristics. This value is experimentally chosen such that it results in modifications that are clearly distinguishable for the human eye, while still being perceived by ProtoPNet as being close to the original images. For the colour modifications (contrast, saturation and hue), we use PyTorch’s image transformations². More specifically, we use the `ColorJitter` function where we set the contrast value to 0.45, saturation to 0.7 and hue to 0.1 for the respective modifications. The shape modification is manually implemented in Python. The texture modification is implemented with the Non-local Means Denoising algorithm for coloured images in OpenCV³. The filter strength of the denoising algorithm is set to 4 to get the correct mean latent distance.

4.2 Considerations for Evaluation

We would like to emphasise that we do not want to explain human perception, but the perception of the prototype-based model. Hence, we cannot ask users what they deem important, since we aim to explain the model’s reasoning. Also, we cannot construct a ground-truth since we are opening up a “black-box” for

² <https://pytorch.org/docs/stable/torchvision/>, accessed June 2020.

³ https://docs.opencv.org/3.4/d1/d79/group__photo__denoise.html#ga03aa4189fc3e31dafd638d90de335617, accessed June 2020.

which no ground-truth is available. However, we can still do a quantitative analysis to evaluate the generalizability and robustness of the explanations. If generalised well, one would expect that global importance scores are similar when computed for different image sets. We also evaluate the distribution of global importance scores to get more insight in the general model reasoning. Additionally, we qualitatively analyse a varied selection of local explanations (Sect. 5.1) and global explanations (Sect. 5.2). In Sect. 5.3 we use our approach to analyse potential redundancy of prototypes.

5 Results and Discussion





ProtoPNet is trained for 30 epochs, reaching a test accuracy of 78.3%. Having applied the same data and training process as the original work [7], we do not know why our accuracy is lower than the accuracy reported in the original work (80.2%). However, the aim of this paper is not to train the best ProtoPNet, but to find a reasonable well-performing model in order to explain its prototypes.

Figure 4 (left) shows a selection of prototypical patches (‘prototypes’) learned by ProtoPNet. ProtoPNet measures the similarity between a prototype and patches in a given test image. The resulting similarity scores are multiplied with learned weights resulting in a final score per class. An image is classified as the class with the highest score (i.e. most points).





5.1 Analysing Our Local Explanations

Figure 4 (right) shows how our local explanations complement the prototypical reasoning by explaining which visual characteristics were important for ProtoPNet’s similarity score between a prototype and a specific image. We show the activation map as implemented by Chen et al. [7] and list the local importances for each visual characteristic. The importances identified by our local explanations for the test image shown in Fig. 4 seems reasonable given the typical white colour of the pelican and its long neck. Furthermore, our explanations enable a user to understand why ProtoPNet gave a high similarity score to a prototype and a patch in the test image. Whereas the prototypes of the white pelican look similar to the human eye, our local importance scores can differentiate between prototypes and estimate the prototype’s purpose. The topmost prototype in Fig. 4 mostly focuses on shape and contrast, the second prototype deems hue important and the third prototype focuses on texture.

Our local explanations can also be useful to confirm the user’s expectations, such as the importance of the yellow colour for the prototype in Fig. 5a. The explanations are however especially insightful when the given similarity score is in contrast with human perceptual similarity and an explanation is needed. Figure 5b explains why different test images received a high similarity score by ProtoPNet. Our local importance scores therefore serves as an extension to a trained prototype-based model, to explain a single prediction.

Prototype	Test Image Similarity score 2.7421	Test Image Similarity score 2.5124	Test Image Similarity score 2.3628
			
Global Importance:	Local Importance:	Local Importance:	Local Importance:
Hue 0.1638	Hue 2.1849	Hue 2.0172	Hue 1.9387
Saturation 0.0397	Saturation 0.8953	Contrast 0.6743	Saturation 0.6470
Shape 0.0358	Shape 0.7201	Saturation 0.4116	Shape 0.1638
Contrast 0.0089	Texture 0.6541	Shape 0.2912	Contrast 0.0956
Texture -0.0025	Contrast -0.1759	Texture -0.2301	Texture -0.4691

(a) As expected, the yellow hue is dominant, and the local explanations correspond with the global importance.

Prototype	Test Image Similarity score 2.7772	Test Image Similarity score 2.6599	Test Image Similarity score 2.2997
			
Global Importance:	Local Importance:	Local Importance:	Local Importance:
Shape 0.1192	Shape 1.9766	Contrast 0.8402	Hue 0.5265
Texture 0.0133	Hue 0.6577	Hue 0.7093	Texture 0.4193
Contrast 0.0090	Saturation 0.3983	Shape 0.1600	Contrast 0.3327
Hue 0.0002	Contrast 0.3424	Texture 0.1248	Saturation -0.0097
Saturation -0.0027	Texture -0.0909	Saturation -0.1429	Shape -0.0148

(b) Test images can get high similarity scores for different reasons.

Fig. 5. Similarity between a class-specific prototype and test images from different classes explained by our local importance scores.

5.2 Analysing Our Global Explanations

Local explanations are useful to explain an unexpected result, but do not give a coherent, overall explanation of the prototype-based model. Our methodology therefore also produces *global* explanations that give an average view regarding the importance scores for each prototype. Specifically, the global importance scores are computed for each prototype by taking the weighted mean of all training images, as introduced in Sect. 3.3. Hence, these explanations are independent of test input.

Quantitative Evaluation. To quantitatively evaluate our global explanations, we compute global importance scores not only for S_{train} , but for evaluation purposes also for S_{test} . We confirmed with a Shapiro-Wilk test that the importance scores for each characteristic are normally distributed, such that we could apply the Welch t-test to confirm that there is no significant difference between the global importance scores of all prototypes calculated from S_{train} and from

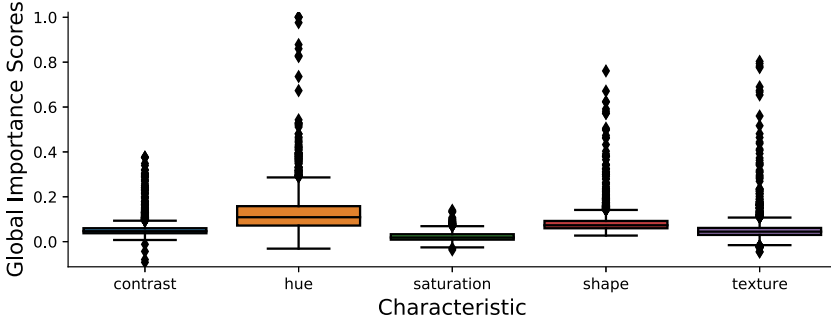


Fig. 6. Box plot of global importance scores across the training set.

S_{test} for each characteristic (p-values $< 3.5^{-11}$). This verifies that our global importance scores are generalisable and robust, since the explanations do not significantly change when computed from a different image set.

To get more insight as to how characteristics are used, Fig. 6 plots the distributions of global importance scores of all prototypes across the training set. It shows that the global scores are predominantly positive, which confirms our intuition that decreasing any of the visual characteristics usually leads to a lower similarity score. It also shows that the mean and variability of importance scores is small for saturation and contrast, meaning that those characteristics only have a moderate influence on prototype similarity. The high variability for shape, texture and especially hue means that these characteristics can be substantially important for some prototypes. This corresponds with the fact that hue and shape are considered more important and effective for humans in the data visualisation domain than saturation or contrast [16].

Qualitative Evaluation. Figure 7 shows a varied selection of prototypes with their global explanation. For the upper two rows, the importance of characteristics corresponds to the visually identifiable properties of the prototypes and hence, the explanations seem reasonable. However, the explanations in the bottom row might come as a surprise. A human might think that shape is important for the bottom left prototype and that the prototype resembles fin-footed birds. Our global explanation indicates that shape is of little importance and that colour hue is the dominant characteristic. Since a ground-truth is not available, we verify the correctness of the global explanation by analysing test images that had a high similarity with the prototype. Although a prototype is trained to be class-specific, Fig. 8a shows that images from a different class can still get assigned a high similarity score. These images confirm that the prototype deems hue important and therefore resembles *red* feet, instead of webbed feet. The reverse is true for the bottom right prototype of Fig. 7. Humans could think that the prototype resembles a red eye, and would be surprised by a high similarity score with a black-eyed bird, and hence might lower their trust in the model. Our global importance scores indicate that the importance for hue is rather low,










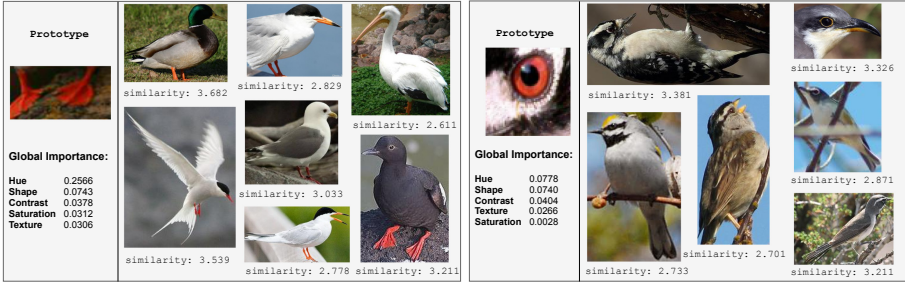
	Hue 0.2154 Texture 0.0388 Shape 0.0356 Contrast 0.0245 Saturation 0.0175		Shape 0.1200 Contrast 0.0156 Saturation 0.0082 Hue 0.0037 Texture -0.0014		Texture 0.1013 Hue 0.0751 Shape 0.0626 Contrast 0.0460 Saturation 0.0083
	Hue 0.1940 Texture 0.1484 Shape 0.1282 Contrast 0.0695 Saturation 0.0116		Hue 0.1343 Shape 0.0745 Texture 0.0325 Contrast 0.0192 Saturation 0.0101		Hue 0.2105 Shape 0.1221 Texture 0.0906 Contrast 0.0724 Saturation 0.0425
	Hue 0.2566 Shape 0.0743 Contrast 0.0378 Saturation 0.0312 Texture 0.0306		Shape 0.0894 Texture 0.0468 Contrast 0.0376 Hue 0.0115 Saturation 0.0048		Hue 0.0778 Shape 0.0740 Contrast 0.0404 Texture 0.0266 Saturation 0.0028

Fig. 7. Selection of prototypes explained with their global importance scores. Top row: Prototypes with predominantly a single important characteristic. Center row: Prototypes with intuitive explanations. Bottom row: ambiguous and potentially misleading prototypes.

and hue and shape are of similar importance. When analysing birds that get assigned a high similarity with this prototype as shown in Fig. 8b, it is easily verified that the prototype does indeed not represent a red eye. These examples show that our global explanations can clarify visual prototypes. Without our explanations, a user would not be aware of the meaning of a given prototype and correct *simulatability* [15] would not be guaranteed.

5.3 Redundant Prototypes

An interesting question is whether prototypes that are slightly different, deem the same visual characteristics important. Prototypes with different global importance scores complement each other, whereas similar explanations indicate prototype redundancy. We consider prototypes to be visually similar when they are close to each other in the latent space learned by ProtoPNet. We measure the Euclidean distance between the latent representation of two prototypes of the same class (a ‘pair’). This gives $\binom{10}{2} = 45$ unique pairs per class, and $45 \cdot 200 = 9000$ pairs in total. Let P be the set of unique pairs of two prototypes from the same class, such that $|P| = 9000$, and $V \subset P$ the set of pairs with two visually similar prototypes. We consider a pair of two prototypes i and j *identical* when the Euclidean distance in latent space $d_{i,j} = 0$ and *visually similar but not identical* when $d_{i,j} < \tau$ and $d_{i,j} > 0$, where $\tau = 0.15$ is found to be a suitable threshold for perceptual similarity. This gives 63 pairs of identical prototypes and $|V| = 93$ unique pairs of 164 visually similar prototypes. To evaluate whether these visually similar prototypes also have similar global explanations, we consider the global importance scores of a prototype as a vector of length 5 and calculate the Euclidean distance between the global explanations of two prototypes. The orange plot in Fig. 9 shows that most pairs with visually similar



(a) The prototype indeed deems hue more important than shape. (b) The global score explains that the red hue from the eye is not that important, which is validated by near test images.

Fig. 8. Test images from a different class than the prototype-class which have the highest similarity scores with the prototype.

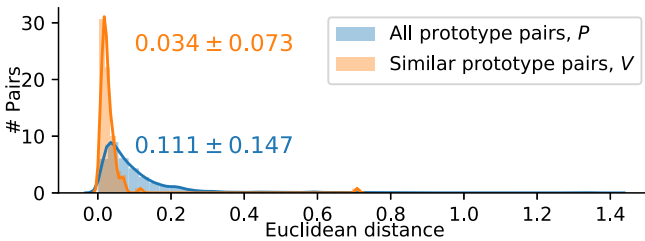


Fig. 9. Histogram with the distribution of Euclidean distances between the global explanations of two prototypes of the same class (a ‘pair’).

prototypes have a small distance between their global importance scores, which is not the case in general (blue). Therefore, these prototypes might be redundant and unnecessarily increase explanation size. Additionally, a few pairs in the orange plot have dissimilar explanations (distance of roughly 0.7) and therefore complement each other.

6 Conclusion and Future Work

A prototype-based image recognition model learns visual prototypes and localises a patch in a test image that looks alike a prototype to assign it a similarity score. We argue that these prototypes should be explained with respect to the model’s reasoning and extend localisation with explanation. We presented an automated approach to explain visual prototypes learned by any prototypical image recognition model. Our method automatically modifies the hue, texture, shape, contrast or saturation of an image, to identify which visual characteristics of a prototype the model deems important. We applied our method to the prototypes learned by ProtoPNet [7]. The importance of visual characteristics identified by our explanations often corresponded to the visually perceptible

properties of the prototypes, showing that our explanations are reasonable. We also showed that perceptual similarity for humans can be different from the similarity learned by the model, indicating the need for explaining the model’s reasoning. Such ‘misleading’ prototypes will hinder correct simulatability and only visualising prototypes can be insufficient for understanding why the model considered a prototype and an image highly similar. To the best of our knowledge, we are the first to address such ambiguity of visual prototypes and the elegant simplicity of our approach makes it a suitable stand-alone solution. We think the extra computational complexity required is justifiable given the extra insights our method provides. Furthermore, because of the stand-alone nature of our method, it can be applied to any prototypical image recognition method, including ProtoPNet [7] and ProtoTree [17]. Our approach can also easily be extended with more visual characteristics or other image modifications a user is interested in.

Future work concerns the potential interactions between characteristics. Our importance scores assume that characteristics from image modifications are mutually exclusive. However, denoising the image to lower its texture could also slightly influence shape. We implemented the image modifications in such a way to limit interactions between characteristics as much as possible, but future analysis could determine to what extent visual characteristics are correlated.

References

1. Arik, S.Ö., Pfister, T.: Attention-based prototypical learning towards interpretable, confident and robust deep neural networks. CoRR abs/1902.06292 (2019)
2. Bau, D., Zhou, B., Khosla, A., Oliva, A., Torralba, A.: Network dissection: quantifying interpretability of deep visual representations. In: CVPR (2017)
3. Op de Beeck, H.P., Torfs, K., Wagemans, J.: Perceived shape similarity among unfamiliar objects and the organization of the human object vision pathway. *J. Neurosci.* **28**(40), 10111–10123 (2008)
4. Biehl, M., Hammer, B., Villmann, T.: Prototype-based models in machine learning. *Wiley Interdisc. Rev. Cognitive Sci.* **7**(2), 92–111 (2016)
5. Cavina-Pratesi, C., Kentridge, R., Heywood, C., Milner, A.: Separate channels for processing form, texture, and color: evidence from fMRI adaptation and visual object agnosia. *Cereb. Cortex* **20**(10), 2319–2332 (2010)
6. Chaudhari, S., Polatkan, G., Ramanath, R., Mithal, V.: An attentive survey of attention models. CoRR abs/1904.02874 (2019)
7. Chen, C., Li, O., Tao, D., Barnett, A., Rudin, C., Su, J.K.: This looks like that: deep learning for interpretable image recognition. In: NeurIPS, pp. 8928–8939 (2019)
8. Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F.A., Brendel, W.: Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In: ICLR (2019)
9. Goodfellow, I., Bengio, Y., Courville, A.: *Deep Learning*. MIT Press, Cambridge (2016)
10. Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., Pedreschi, D.: A survey of methods for explaining black box models. *ACM Comput. Surv. (CSUR)* **51**(5), 1–42 (2018)

11. Hosseini, H., Xiao, B., Jaiswal, M., Poovendran, R.: Assessing shape bias property of convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops (June 2018)
12. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: CVPR, pp. 4700–4708 (2017)
13. King, M.L., Groen, I.I., Steel, A., Kravitz, D.J., Baker, C.I.: Similarity judgments and cortical visual responses reflect different properties of object and scene categories in naturalistic images. *NeuroImage* **197**, 368–382 (2019)
14. Kourtzi, Z., Kanwisher, N.: Cortical regions involved in perceiving object shape. *J. Neurosci.* **20**(9), 3310–3318 (2000)
15. Lipton, Z.C.: The mythos of model interpretability. *Queue* **16**(3), 31–57 (2018)
16. Munzner, T., Maguire, E.: *Visualization Analysis & Design*. CRC Press, Boca Raton (2015)
17. Nauta, M., van Bree, R., Seifert, C.: Neural prototype trees for interpretable fine-grained image recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 14933–14943 (June 2021)
18. Ritter, S., Barrett, D.G.T., Santoro, A., Botvinick, M.M.: Cognitive psychology for deep neural networks: a shape bias case study. In: ICML Proceedings of Machine Learning Research, vol. 70, pp. 2940–2949. PMLR (2017)
19. Rosenfeld, A., Solbach, M.D., Tsotsos, J.K.: Totally looks like - how humans compare, compared to machines. In: CVPR Workshops (2018)
20. Rossion, B., Pourtois, G.: Revisiting snodgrass and vanderwart’s object pictorial set: The role of surface detail in basic-level object recognition. *Perception* **33**(2), 217–236 (2004). pMID: 15109163
21. Rudin, C.: Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat. Mach. Intell.* **1**(5), 206–215 (2019)
22. Saralajew, S., Holdijk, L., Rees, M., Asan, E., Villmann, T.: Classification-by-components: probabilistic modeling of reasoning over a set of components. In: NeurIPS, pp. 2792–2803 (2019)
23. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: Bengio, Y., LeCun, Y. (eds.) ICLR (2015)
24. Striedter, G.F.: *Neurobiology: a Functional Approach*. Oxford University Press, Oxford (2016)
25. Welinder, P., et al.: Caltech-UCSD Birds 200. Technical Report. CNS-TR-2010-001 (2010)
26. Zhang, N., Donahue, J., Girshick, R., Darrell, T.: Part-based R-CNNs for fine-grained category detection. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) ECCV 2014. LNCS, vol. 8689, pp. 834–849. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-10590-1_54
27. Zheng, H., Fu, J., Zha, Z., Luo, J., Mei, T.: Learning rich part hierarchies with progressive attention networks for fine-grained image recognition. *IEEE Trans. Image Process.* **29**, 476–488 (2020)