# Evaluation of Explainable Artificial Intelligence techniques

The latest generation of artificial intelligence has proven to be extremely powerful and valuable, with applications in various fields like medicine and computer science. Its learning techniques can be difficult for humans to understand. But can we trust a program that can learn and decide independently without understanding how it works?

To address this issue, a new field of research called XAI (eXplainable Artificial Intelligence) has emerged. Its goal is to make the decisions made by artificial intelligence understandable to humans. This is important not only for understanding how AI works, but also for increasing its transparency and people's trust in it.

Your task is to answer a series of questions to evaluate different XAI techniques, considering the outputs of these techniques, which are provided in this document.

The AI system that we want to explain with this experiment is a Neural Network capable of classifying images by predicting the correct class. We present four correctly classified examples, which will be used throughout the whole document:

**Cassette Player**
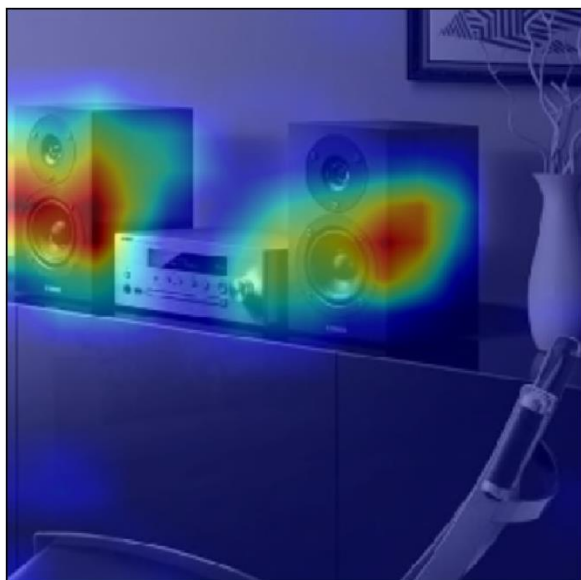


**English Springer**



**Golf Ball**



**Tench**

# Grad-CAM

This technique produces an explanation in the form of a heatmap where the "hot regions" (the part coloured in red) represent the most important parts towards the prediction.
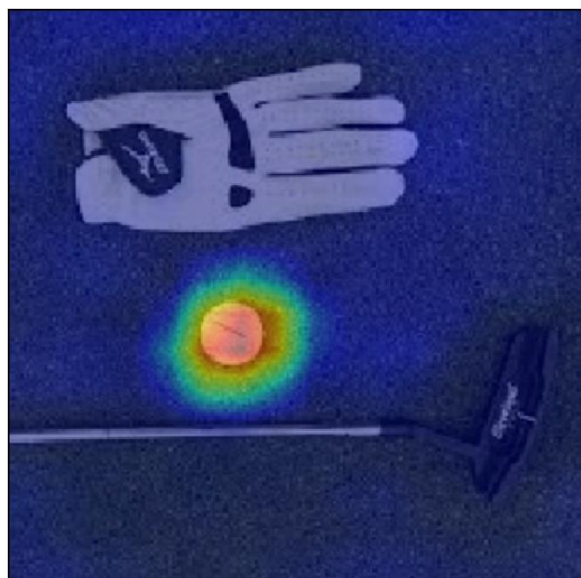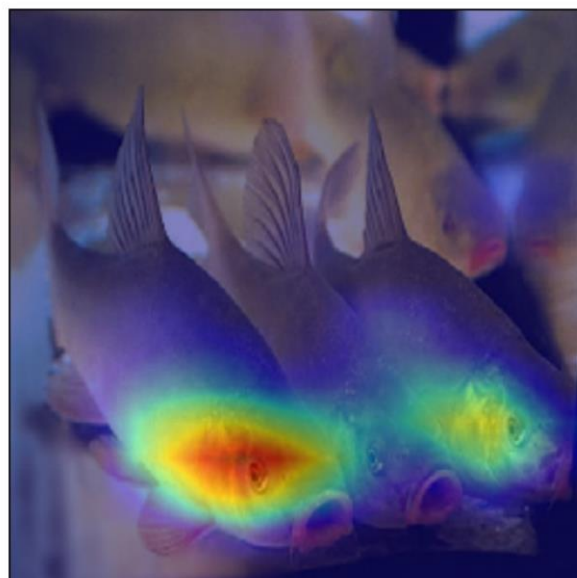
**Cassette Player**
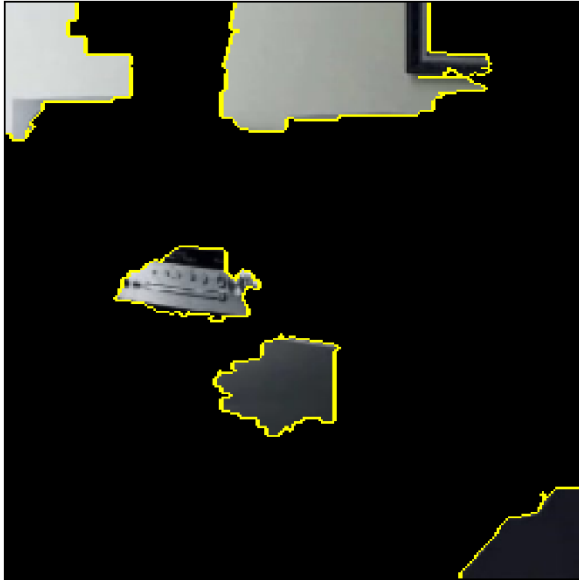


**English Springer**



**Golf Ball**



**Tench**

# LIME

This method divides an image into different regions, called superpixels, and shows the ones that, if removed, change the network's output the most.
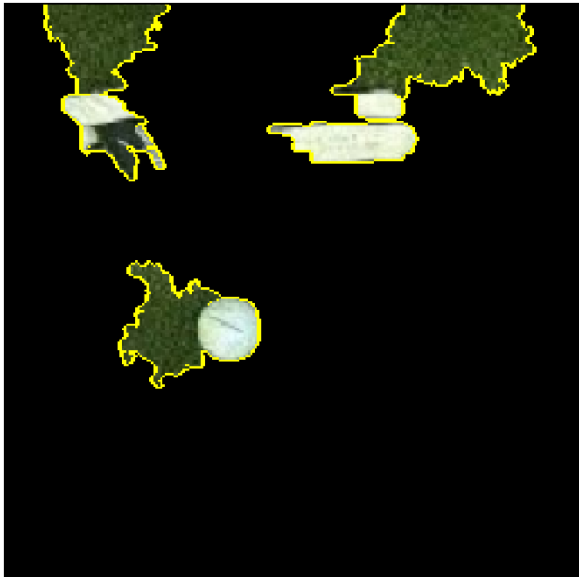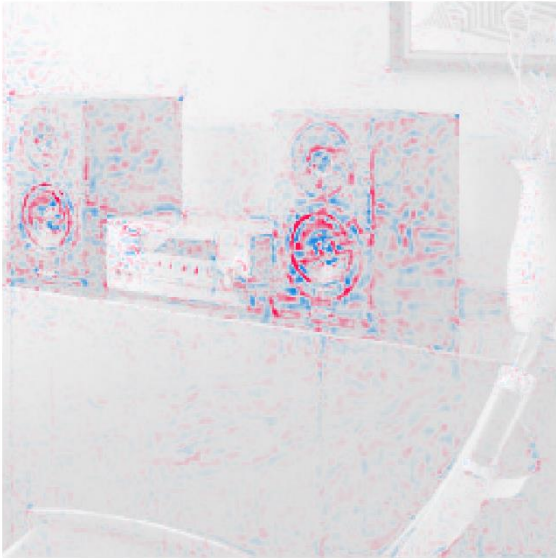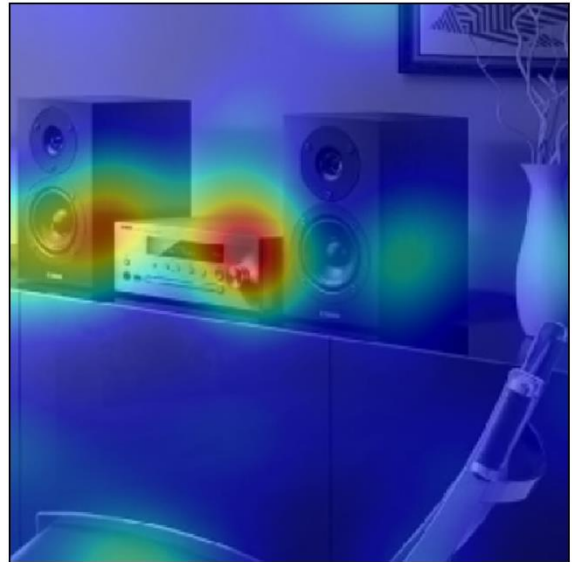
**Cassette Player**



**English Springer**



**Golf Ball**



**Tench**

# SHAP

It computes the contribution of each pixel of the image, using a game theory approach. More specifically, red pixels represent a positive contribution, while blue pixels represent a negative contribution towards the prediction.

**Cassette Player**

**English Springer**





**Golf Ball**

**Tench**

# INV

This method extracts features that the network used to make the prediction, with their corresponding contribution towards the prediction (w). The features are represented as heatmaps, like Grad-CAM, and are associated with a name.

## Cassette Player

### w=73.9% - speakers



### w=19.0% - volume



## English Springer

### w=32.4% - paws



### w=24.2% - muzzle

**Golf Ball**

# w=93.6% - ball



**Tench**

# w=73.1% - eye



# w=16.0% - fins

# Extended INV

The building blocks of a neural network are called layers. They represent groups of neurons which extract features from the input image. The initial layer extract simple features like shapes, colours, etc., which are sent to the subsequent layer to extract more complex features that are more helpful towards making a correct prediction. This process is repeated for N layers, and finally an output is produced.

This explanation extends the previous technique, which considered only the final layer, by showing the features extracted also in earlier layers, with their corresponding weight and the labels that better describe the portion of image shown. Each column represents a layer, identified by its name (for instance, block5_conv3)

## Cassette Player



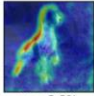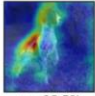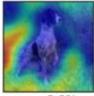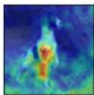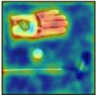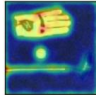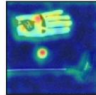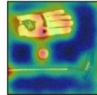| block3_conv1 | block3_conv2 | block3_conv3 | block4_conv1 | block4_conv2 | block4_conv3 | block5_conv1 | block5_conv2 | block5_conv3 |
|---|---|---|---|---|---|---|---|---|
| w = 36.2% | w = 41.9% | w = 73.9% | w = 58.8% | w = 53.0% | w = 57.5% | w = 44.7% | w = 28.3% | w = 73.9% |
| interior | display | speakers | speakers | speakers | speakers | speakers | speakers | speakers |
| parallelepiped | input | buttons | subwoofer | parts | buttons | crate | black | case |
| speakers | location | knobs | tweeter | amplifier | case | colors | case | rectangular |
| w = 15.5% | w = 16.6% | | w = 9.2% | w = 21.4% | w = 11.6% | w = 16.0% | w = 18.4% | w = 19.0% |
| speakers | volume | | keys | display | mobile | cassette | case | knobs |
| rectangular | case | | aux | exclusion | colors | form | display | volume |
| buttons | form | | | jack | notebook | input | speakers | case |
| | w = 14.9% | | | | w = 7.7% | w = 10.5% | w = 14.2% | |
| | buttons | | | | cables | buttons | buttons | |
| | display | | | | case | cassette | cd | |
| | colors | | | | decorations | picture | display | |
| | | | | | | w = 5.4% | w = 8.8% | |
| | | | | | | flowerpot | interior | |
| | | | | | | home | design | |
| | | | | | | rectangle | picture | |
| | | | | | | | w = 6.2% | |
| | | | | | | | cabinet | |
| | | | | | | | buttons | |
| | | | | | | | lines | |

# English Springer



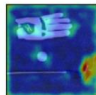| block3_conv1 | block3_conv2 | block3_conv3 | block4_conv1 | block4_conv2 | block4_conv3 | block5_conv1 | block5_conv2 | block5_conv3 |
|---|---|---|---|---|---|---|---|---|
| w = 37.7% | w = 57.8% | w = 47.1% | w = 34.2% | w = 77.4% | w = 69.1% | w = 63.0% | w = 31.4% | w = 32.4% |
| paws | muzzle | muzzle | muzzle | muzzle | muzzle | muzzle | paws | paws |
| hair | nose | nose | color | paws | hair | fur | hair | dalmatian |
| snout | spots | hair | face | fur | fur | nose | body | spots |
| w = 37.6% | | w = 13.1% | w = 21.3% | | | w = 18.1% | w = 24.6% | w = 24.2% |
| muzzle | | color | grass | | | legs | muzzle | muzzle |
| color | | dots | coat | | | paws | | face |
| nose | | height | | | | spots | | |
| | | w = 8.8% | w = 19.5% | | | | w = 9.0% | |
| | | grass | hair | | | | meadow | |
| | | legs | fur | | | | grass | |
| | | dots | ears | | | | | |
| | | w = 6.0% | | | | | | |
| | | spots | | | | | | |
| | | hair | | | | | | |
| | | color | | | | | | |

# Golf Ball

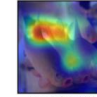| block3_conv1 | block3_conv2 | block3_conv3 | block4_conv1 | block4_conv2 | block4_conv3 | block5_conv1 | block5_conv2 | block5_conv3 |
|---|---|---|---|---|---|---|---|---|
| w = 56.3% | w = 39.0% | w = 27.4% | w = 38.7% | w = 60.8% | w = 33.9% | w = 46.7% | w = 35.6% | w = 93.6% |
| glove | glove | ball | ball | ball | lawn | ball | ball | ball |
| grass | ball | glove | see | lawn | exclusion | grass | white | lawn |
| ball | club | see | lawn | golfball | golfball | visible | grass | white |
| w = 22.6% | w = 20.9% | w = 15.2% | w = 34.2% | w = 14.1% | w = 26.9% | w = 19.1% | w = 33.6% | |
| ball | white | lawn | white | lawn | ball | lawn | lawn | |
| glove | field | texture | glove | | glove | grass | grass | |
| stand | lawn | white | grass | | logo | | glove | |
| | w = 13.7% | w = 15.1% | | | w = 13.1% | | w = 6.9% | |
| | ball | nike | | | grass | | club | |
| | exclusion | field | | | lawn | | lawn | |
| | part | white | | | ball | | glove | |
| | | w = 12.7% | | | | | | |
| | | glove | | | | | | |
| | | ball | | | | | | |
| | | grass | | | | | | |

# Tench

| block3_conv1 | block3_conv2 | block3_conv3 | block4_conv1 | block4_conv2 | block4_conv3 | block5_conv1 | block5_conv2 | block5_conv3 |
|---|---|---|---|---|---|---|---|---|
|  |  |  |  |  |  |  |  |  |
| w = 38.3% | w = 74.5% | w = 37.4% | w = 57.5% | w = 27.2% | w = 43.7% | w = 66.8% | w = 61.3% | w = 73.1% |
| fins | mouth | fins | eye | eye | fins | eye | eye | eye |
| eye | fins | eye | mouth | fins | mouth | mouth | fisheye | mouth |
| background | eye | color | scales | head | randomly | head | snout | dull |
|  | |  |  |  |  | |  |  |
| w = 22.6% | | w = 27.7% | w = 16.7% | w = 23.9% | w = 23.7% | | w = 12.0% | w = 16.0% |
| mouth | | mouth | fins | fins | eye | | mouth | fins |
| silouette | | sequins | wings | eye | mouth | | eye | conformation |
| eye | | eye | eye | fisheye | face | | aquarium | |
|  | | | |  | | | | |
| w = 7.7% | | | | w = 16.5% | | | | |
| tail | | | | face | | | | |
| brass | | | | eye | | | | |
| | | | | mouth | | | | |
| | | | |  | | | | |
| | | | | w = 5.7% | | | | |
| | | | | mouth | | | | |
| | | | | color | | | | |
| | | | | fins | | | | |