



DEEL XPLIQUE

Explainability Toolbox for Neural Networks

Presentation given by Lucas Hervier and Antonin Poché
first_name dot last_name at irt-saintexupery dot com



Table of content

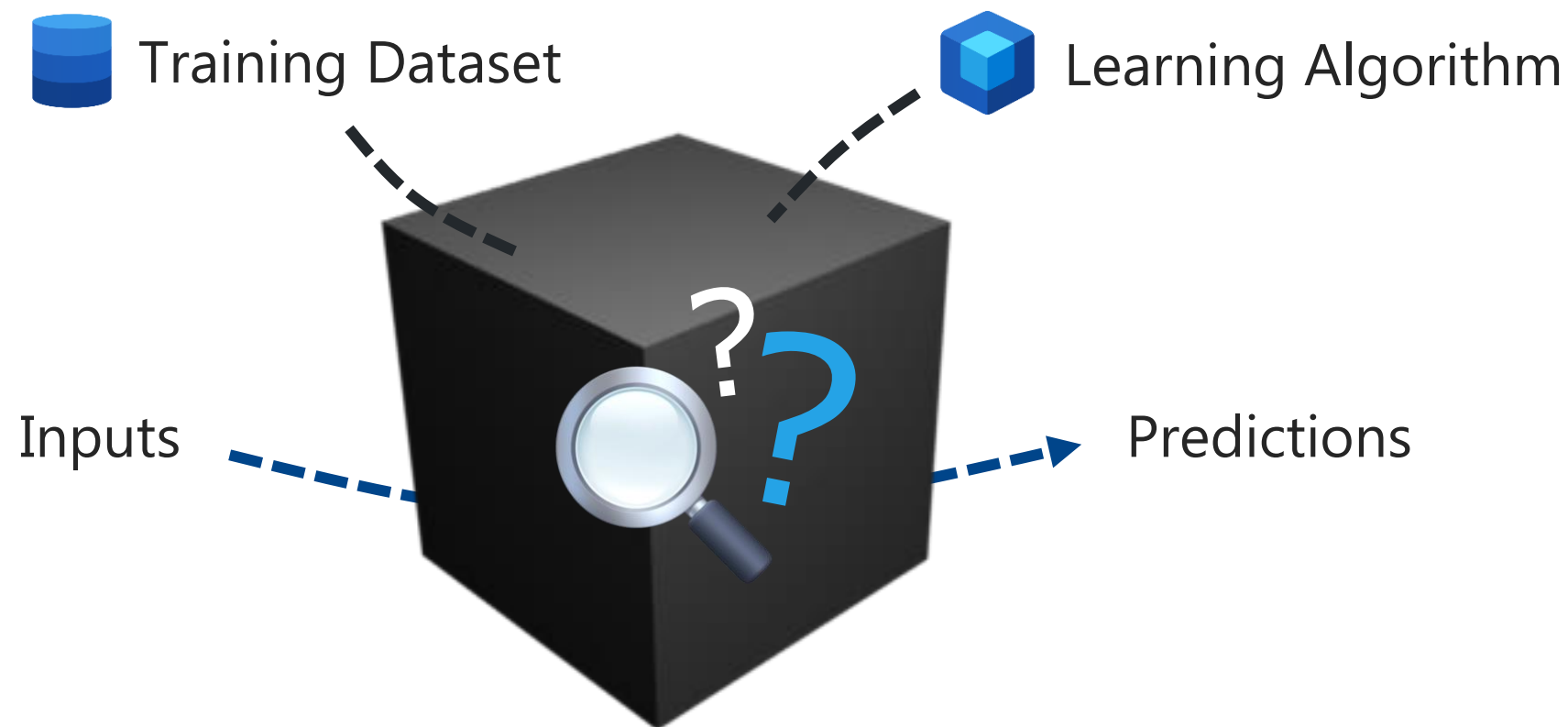
- Introduction to explainability and Xplique
- What's new?
- Two demonstrations
- Feed-back from an user



Explainability and Xplique



The black-box problem



Why do we need explainability?

- Build **trust** in the model predictions

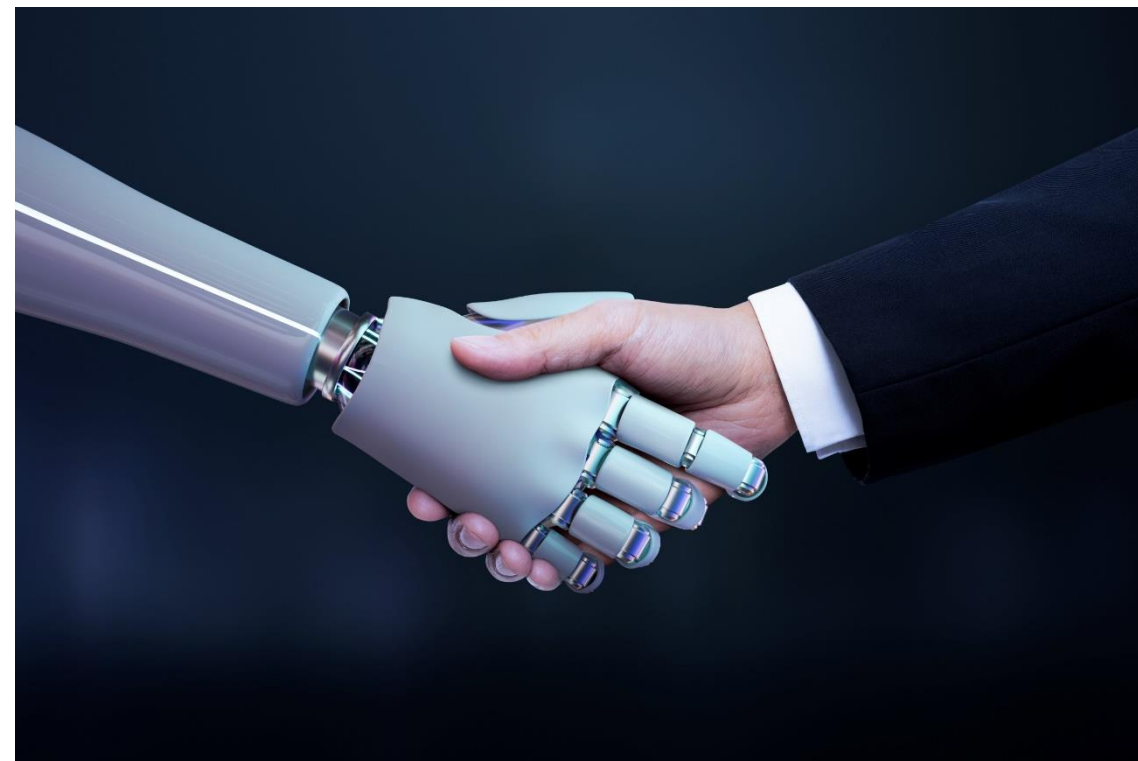


Image by rawpixel.com on Freepik

Why do we need explainability?

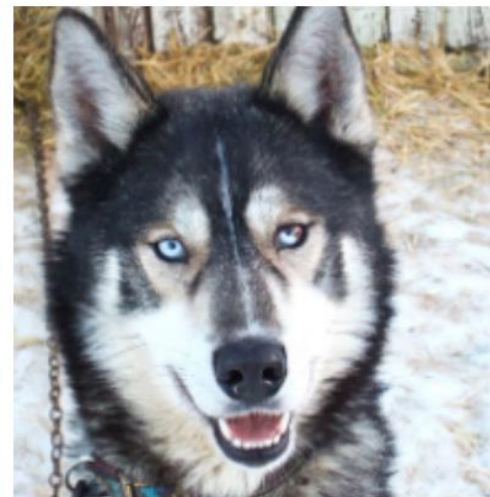
- Build **trust** in the model predictions
- **Satisfy regulatory requirements** and Certification process



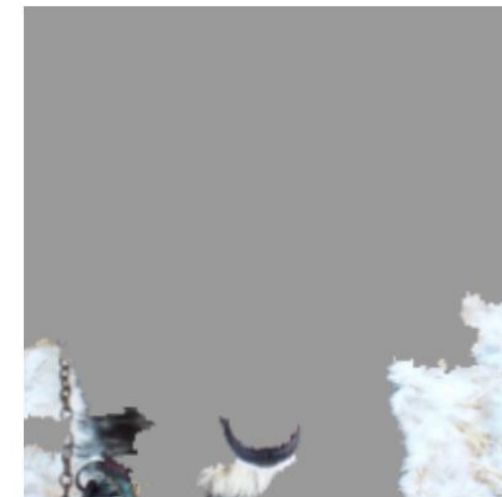
Image by rawpixel.com on Freepik

Why do we need explainability?

- Build **trust** in the model predictions
- **Satisfy regulatory requirements** and Certification process
- **Reveal bias** or other unintended effects learned by a model
- Understand to **intervene on models**
- ...



(a) Husky classified as wolf



(b) Explanation

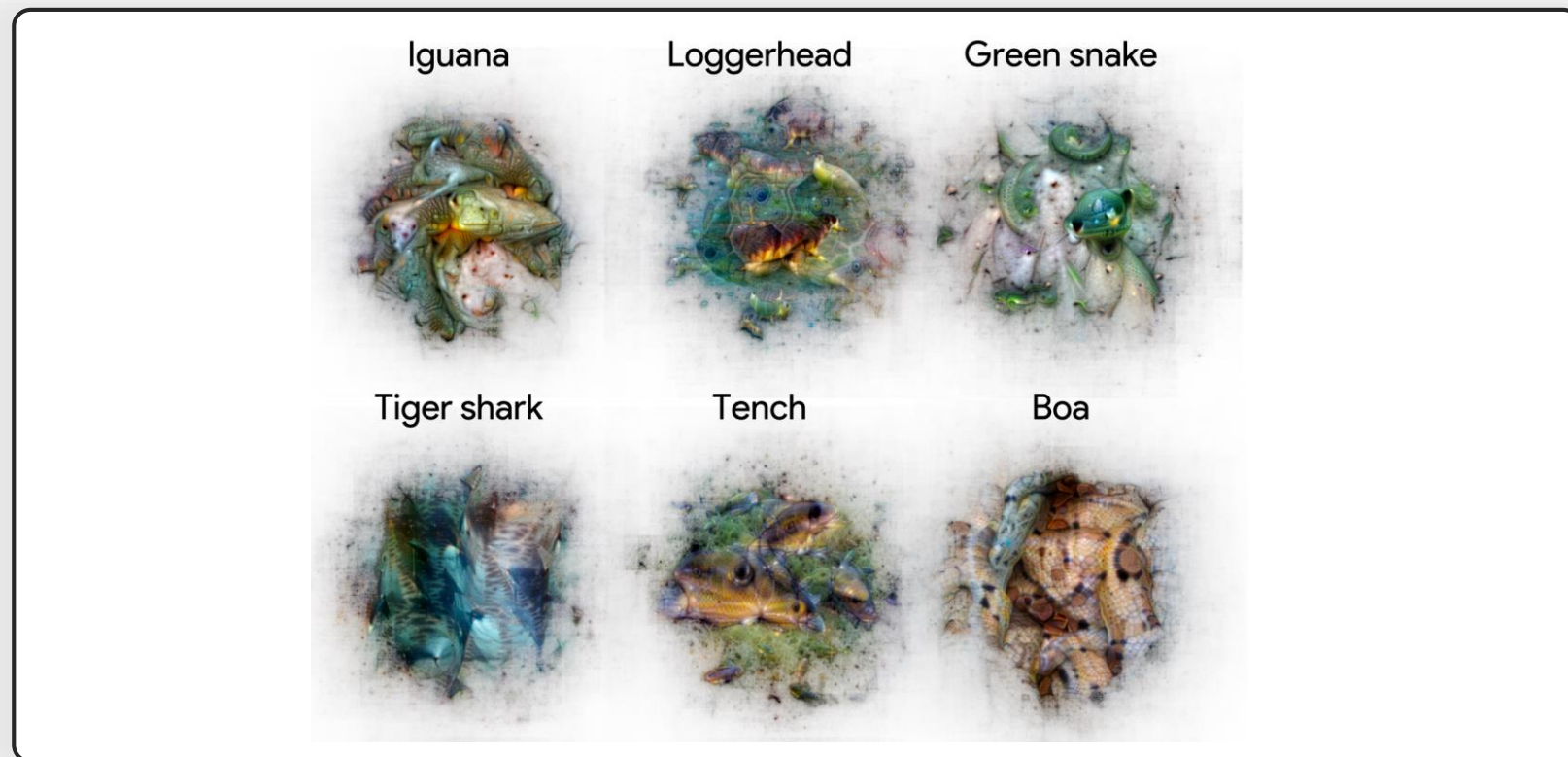
From Ribeiro et al.: "Why Should I Trust You?"

A technical challenge

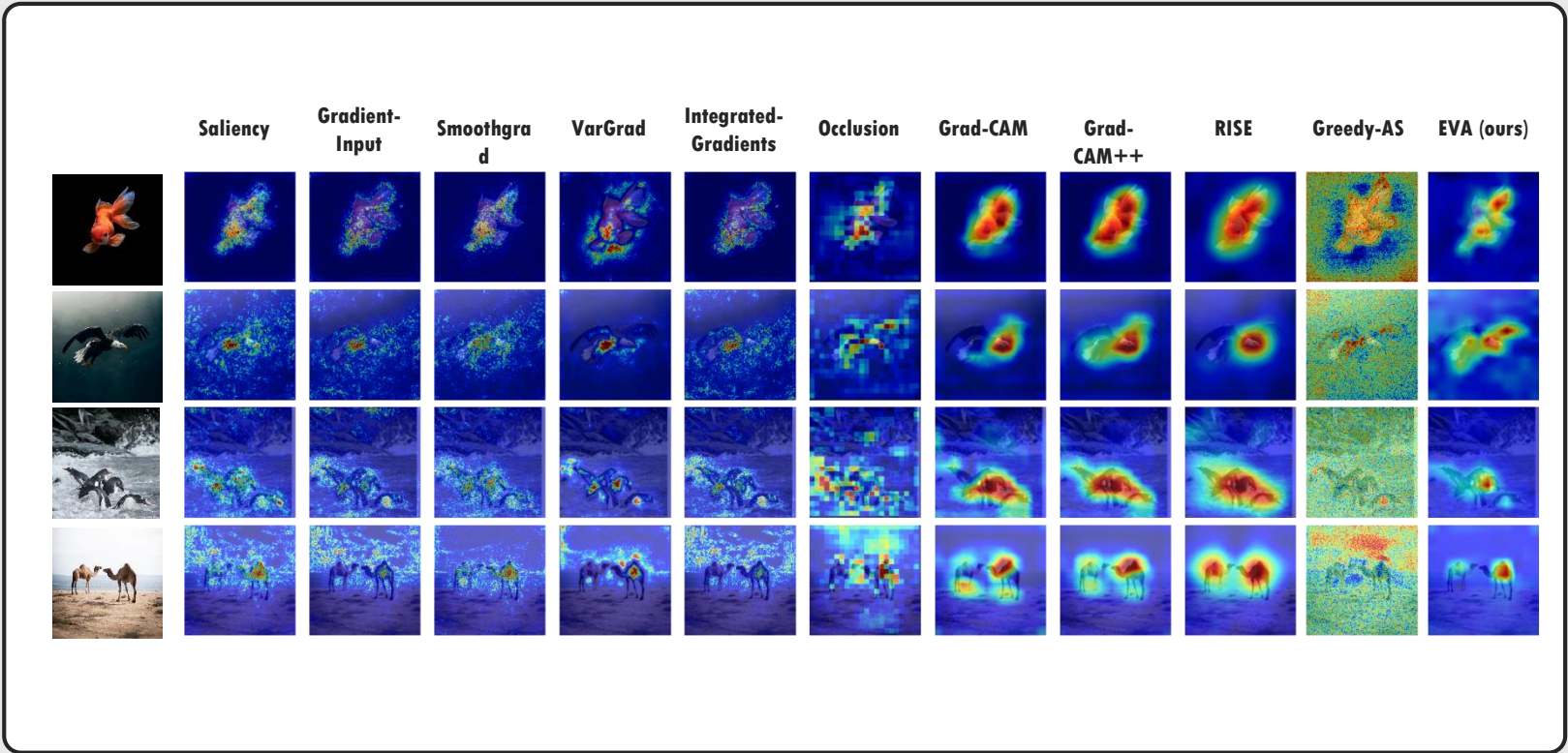
© DEEL - All rights reserved to IVADO, IIR, IRT Saint Exupéry, CRIAQ and ANITI. Confidential and proprietary document

Model

Feature Viz,
Concept Activation Vector
Explanation 'by design'
...

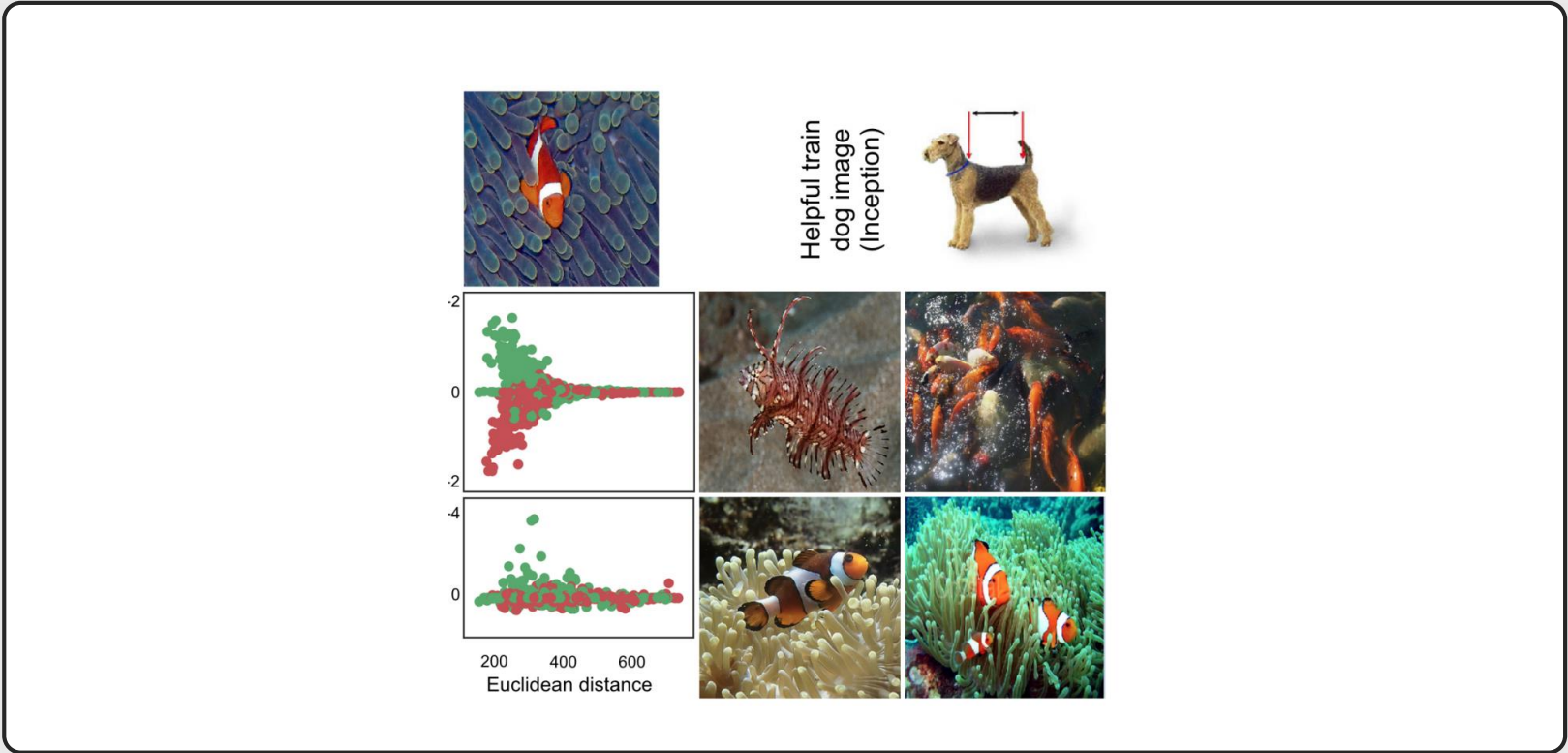


A technical challenge



Predictions
Feature Attribution
Feature Inversion
...

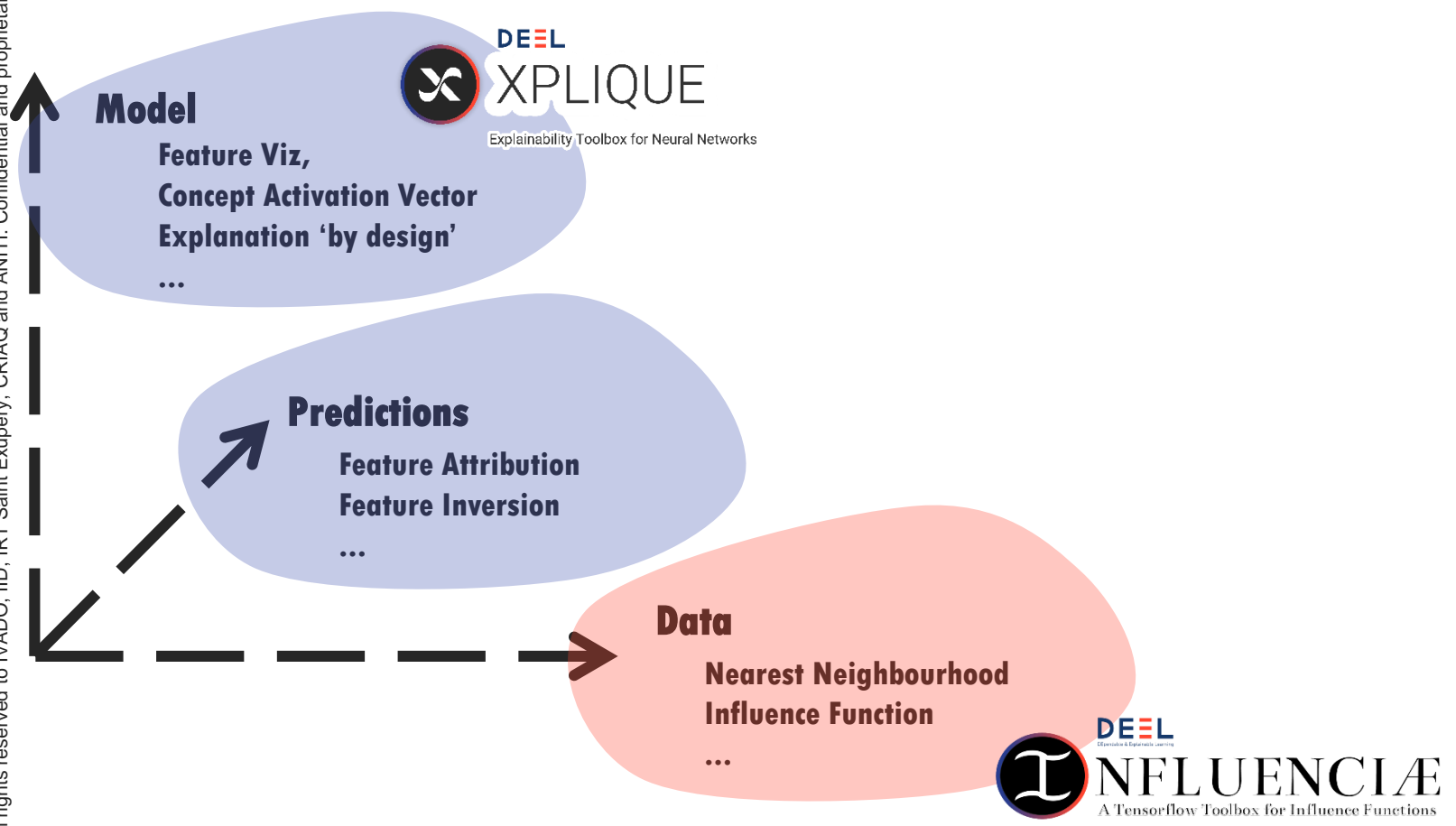
A technical challenge



Data
**Nearest Neighbourhood
Influence Function**
...

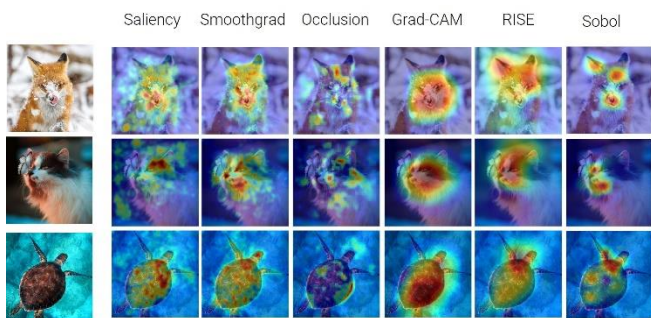
A technical challenge

© DEEL - All rights reserved to IVADO, IIR, IRT Saint Exupéry, CRIAQ and ANITI. Confidential and proprietary document



Xplique: A DL Explainability Toolbox

Attribution Methods more than 14 black-box / white-box methods*



```

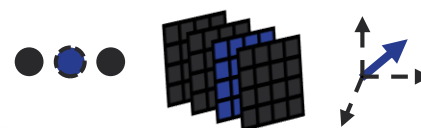
from xplique.attributions import GradCAM

explainer = GradCAM(model)
explanations = explainer(x, y)
    
```

*SOTA

Feature Visualization

• Neurons • Channels • Directions



Visualize Neurons, Channels, Vectors in activation space (e.g. CAV) or a mix of them !

```

from xplique.feature_visualization import Objective,
optimize
obj = Objective.neuron(model, 'Logits', 10)
images, obj_name = optimize(obj)
    
```

'Ladybug'



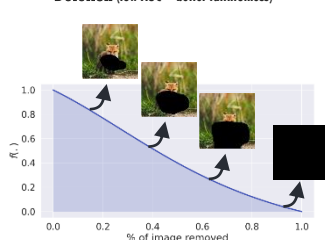
'Goldfish'



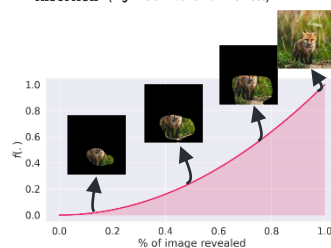
Metrics more than 6 attributions metrics each supporting multiple baselines

multiple

Deletion (low AUC = better faithfulness)



Insertion* (high AUC = better faithfulness)



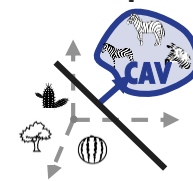
```

from xplique.metrics import Deletion
from xplique.attributions import GradCAM

metric = Deletion(model, x, y)
explanations = GradCAM(model)(x, y)
score = metric(explanations)
    
```

Concept based concept activation vector, CRAFT (new!)

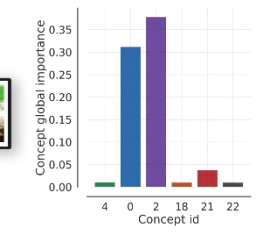
Easily extract and test CAVs:



```

from xplique.concepts import Cav

extractor = Cav(model, 'mixed3')
concept_vector = extractor(striped_samples,
                           random_samples)
    
```





What's new?



deel-ai/xplique

v1.3.3 ☆ 551 🔗 47



List of new features

- PyTorch wrapper
- Operators, they allow to treat new tasks:
 - Object detection
 - Semantic segmentation
- New state-of-the-art attribution methods:
 - Sobol, HSIC, and FORGrad
- Automatic concepts extraction
 - CRAFT
- Last feature visualization methods:
 - MaCO

PyTorch wrapper and operators API

```

1  from xplique.attributions import Saliency
2  from xplique.metrics import Deletion
3
4  # load images, targets, and model
5  # ...
6
7  # initialize the explainer with the model and method parameters
8  explainer = Saliency(model)
9
10 # call the explainer on the sample to explain
11 explanations = explainer(inputs, targets)
12
13 # compute explanation metrics
14 metric = Deletion(model, inputs, targets)
15 deletion_score = metric(explanation)

```

Initial API

```

1  import torch
2
3  from xplique.attributions import Saliency
4  from xplique.metrics import Deletion
5  from xplique.wrappers import TorchWrapper
6
7  # load images, targets, and torch_model
8  # ...
9
10 device = 'cuda' if torch.cuda.is_available() else 'cpu'
11 model = TorchWrapper(torch_model, device)
12
13 # initialize the explainer with the model and method parameters
14 explainer = Saliency(model)
15
16 # call the explainer on the sample to explain
17 explanations = explainer(inputs, targets)
18
19 # compute explanation metrics
20 metric = Deletion(model, inputs, targets)
21 deletion_score = metric(explanation)

```

PyTorch Wrapper

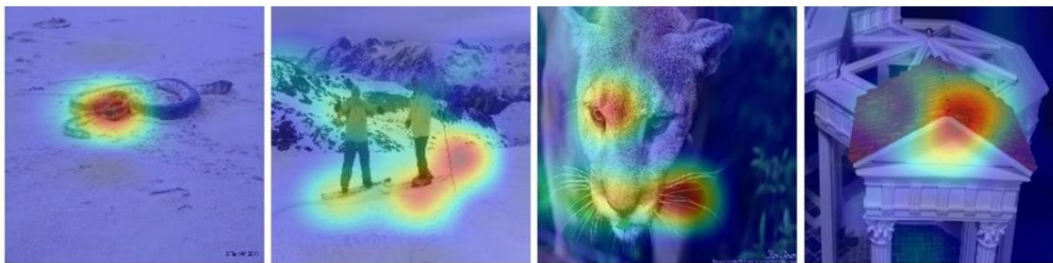
```

1  from xplique.attributions import Saliency
2  from xplique.metrics import Deletion
3
4  # load images, targets, and model
5  # ...
6
7  # initialize the explainer with the model and method parameters
8  explainer = Saliency(model, operator=xplique.Tasks.SEMANTIC_SEGMENTATION)
9
10 # call the explainer on the sample to explain
11 explanations = explainer(inputs, targets)
12
13 # compute explanation metrics
14 metric = Deletion(model, inputs, targets, operator=xplique.Tasks.SEMANTIC_SEGMENTATION)
15 deletion_score = metric(explanation)

```

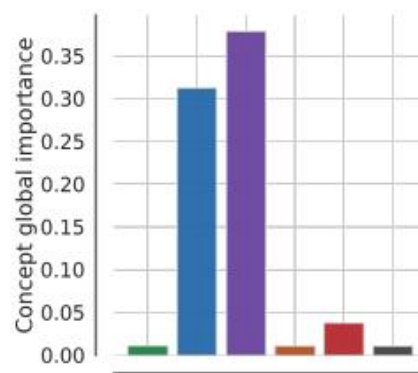
Operators

New methods



HSIC

Novello, Paul, Thomas Fel, and David Vigouroux. "Making sense of dependence: Efficient black-box explanations using dependence measure." *Advances in Neural Information Processing Systems* 35 (2022): 4344-4357.



CRAFT

Fel, Thomas, et al. "Craft: Concept recursive activation factorization for explainability." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.

Bottles



Pineapple



Espresso



Green snake



Bottles



Guitar



Beagle



Banana





Demonstrations



Demonstrations with Xplique tutorials

HSIC on a PyTorch model
for semantic segmentation

CRAFT and MaCO

- PyTorch wrapper
- Operators, they allow to treat new tasks:
 - Object detection
 - Semantic segmentation
- New state-of-the-art attribution methods:
 - Sobol and HSIC
- Automatic concepts extraction
 - CRAFT
- Last feature visualization methods:
 - MaCO

Future Works



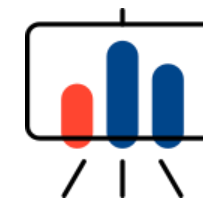
Research

- Development of new methods (example-based)
- Extend to new Use Cases
- Development of new metrics
-



Tools

- Enhance our library efficiency
- Extend to new Use Cases
- CI/CD, include the newest methods
-



Benchmarking

- Provide a thorough User Guide
- Multi-criteria evaluation
- Challenge on more Use Cases
-



Feed-back from an user



Want to know more ?



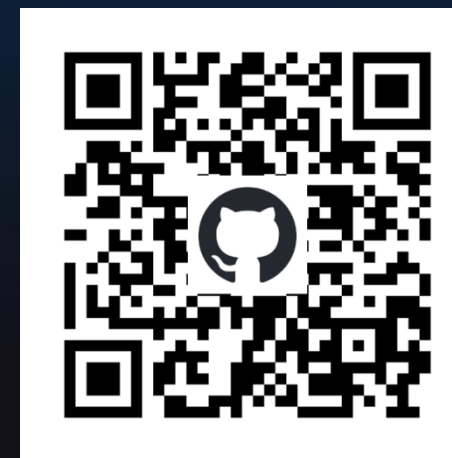
www.deel.ai



[github/deel-ai](https://github.com/deel-ai)



[linkedin/showcase/deel-ai](https://www.linkedin.com/showcase/deel-ai)



Do not forget
to leave a star

Want to join DEEL ?



[gregory dot flandin at irt-saintexupery dot com](mailto:gregory.dot.flandin@irt-saintexupery.com)



Bâtiment B612
3 Rue Tarfaya, 31400 Toulouse