# INTERPRETABLE MACHINE LEARNING

Anna Dawid

# Outline

WHY DO WE NEED INTERPRETABILITY AND WHAT IT IS

OVERVIEW OF INTERPRETABILITY METHODS

TUTORIAL EXAMPLE: CLASS ACTIVATION MAP (CAM)

# WHY DO WE NEED INTERPRETABILITY?

And what it is?

'AI IS THE NEW ELECTRICITY'

"Just as electricity transformed almost everything 100 years ago, today I actually have a hard time thinking of an industry that I don't think AI will transform in the next several years."
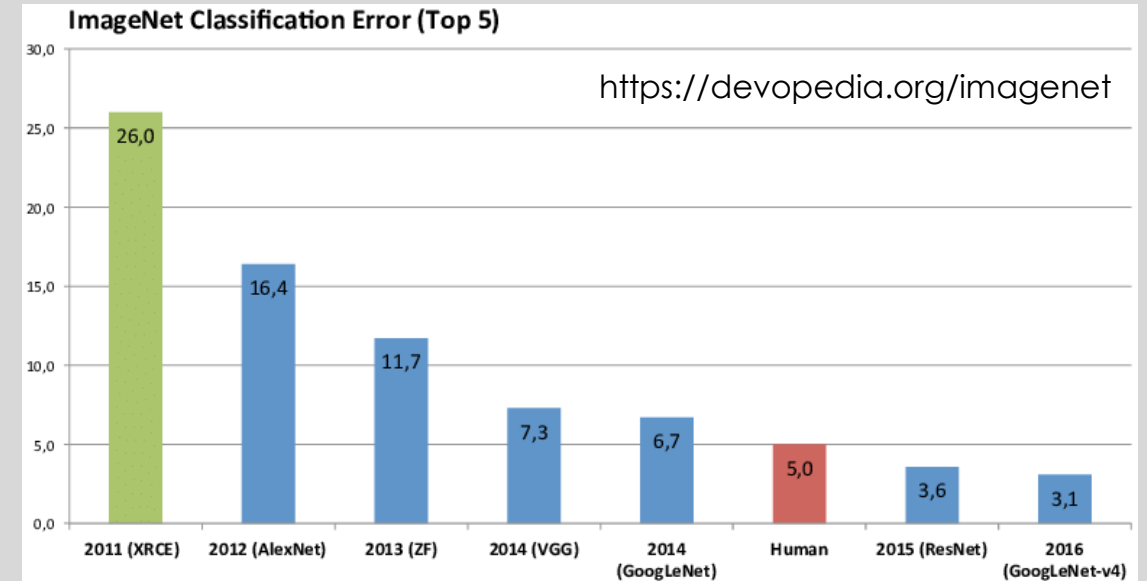
**Andrew Ng**

Former chief scientist at Baidu, Co-founder at Coursera

https://github.com/karolmajek/Mask_RCNN

Computer vision for self-driving cars

Super-human performance on image classification

**ImageNet Classification Error (Top 5)**

https://devopedia.org/imagenet

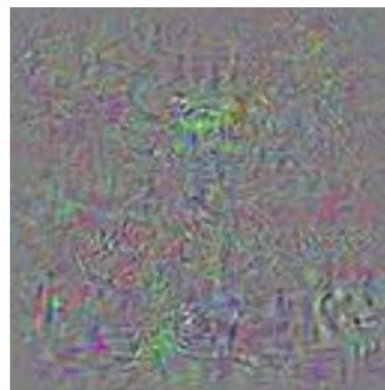| Year (Model) | Error |
|---|---|
| 2011 (XRCE) | 26,0 |
| 2012 (AlexNet) | 16,4 |
| 2013 (ZF) | 11,7 |
| 2014 (VGG) | 7,3 |
| 2014 (GoogLeNet) | 6,7 |
| Human | 5,0 |
| 2015 (ResNet) | 3,6 |
| 2016 (GoogLeNet-v4) | 3,1 |

# All is well?
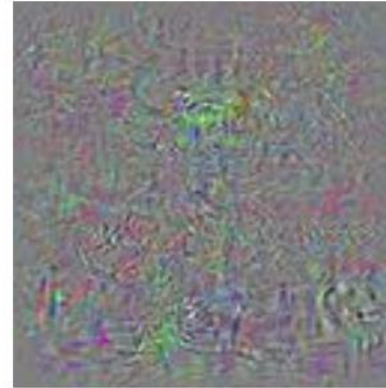
What humans see:

Schoolbus + Perturbation (rescaled for visualization) = Schoolbus

What (some) machines see:

Schoolbus + Perturbation (rescaled for visualization) = Ostrich

modified from Goodfellow (2018), arXiv:1806.04169v1

What humans see:

What (some) machines see:



Schoolbus

Schoolbus                    Perturbation                    Ostrich
                        (rescaled for visualization)

open to
adversarial attacks!

What humans see:

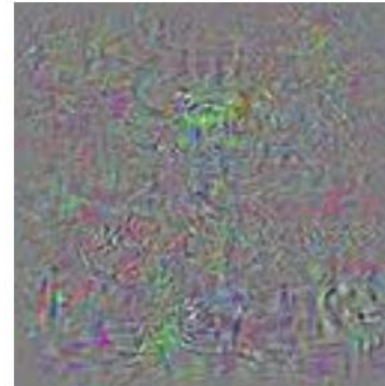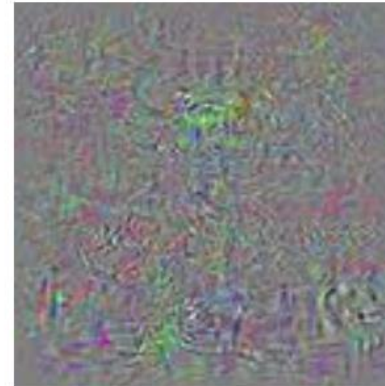Schoolbus + Perturbation (rescaled for visualization) = Schoolbus

What (some) machines see:

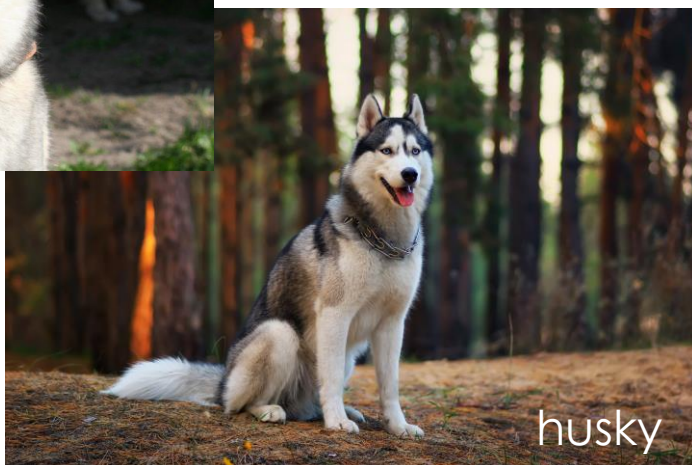Schoolbus + Perturbation (rescaled for visualization) = Ostrich

modified from Goodfellow (2018), arXiv:1806.04169v1

open to
adversarial attacks!

What humans see:

Schoolbus + Perturbation (rescaled for visualization) = Schoolbus

UNRELIABLE!

What (some) machines see:

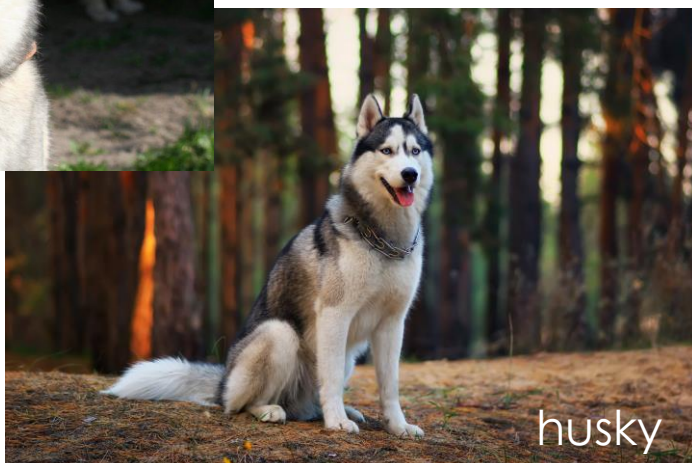Schoolbus + Perturbation (rescaled for visualization) = Ostrich

modified from Goodfellow (2018), arXiv:1806.04169v1

Husky vs. wolf classifier

husky

wolf

wolf

Classifier learnt the background...

husky

husky

wolf

wolf

Ribeiro (2016), arXiv:1602.04938v3

husky

wolf

wolf

husky

Classifier learnt the background...

Bad generalization!
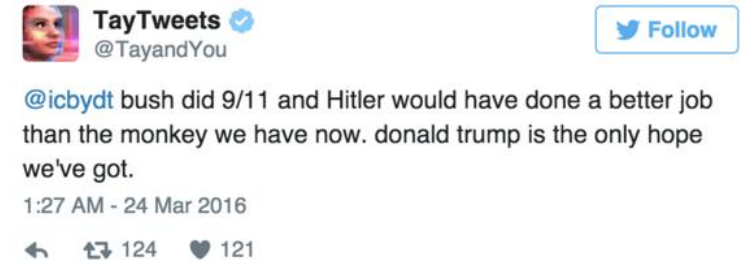
husky

wolf

wolf

Ribeiro (2016), arXiv:1602.04938v3

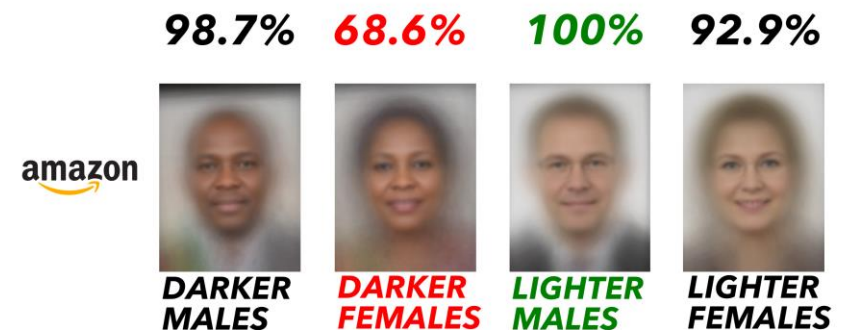AI trained by Amazon inherited biases from historical hiring decisions



https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G

TayTweets, a chat bot trained by Microsoft on Twitter data started spouting racist tweets…



https://towardsdatascience.com/biases-in-machine-learning-61186da78591

Amazon Rekognition failures in facial recognistion correlate with skin color



| 98.7% | 68.6% | 100% | 92.9% |
| DARKER MALES | DARKER FEMALES | LIGHTER MALES | LIGHTER FEMALES |

https://medium.com/@Joy.Buolamwini/response-racial-and-gender-bias-in-amazon-rekognition-commercial-ai-system-for-analyzing-faces-a289222eeced

AI trained by Amazon inherited biases from historical hiring decisions



https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G
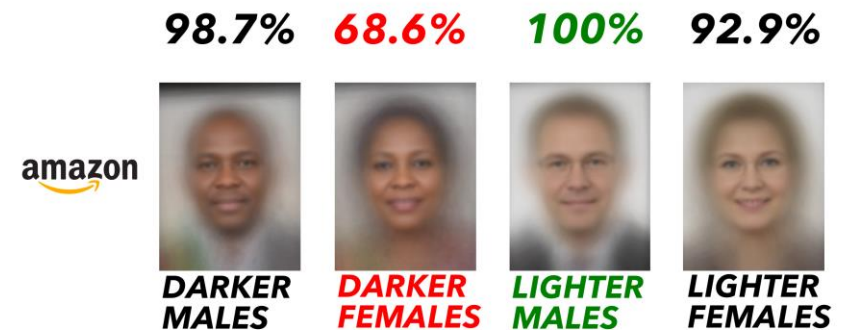
TayTweets, a chat bot trained by Microsoft on Twitter data started spouting racist tweets...



TayTweets ✔
@TayandYou

@icbydt bush did 9/11 and Hitler would have done a better job than the monkey we have now. donald trump is the only hope we've got.

1:27 AM - 24 Mar 2016

124    121

https://towardsdatascience.com/biases-in-machine-learning-61186da78591

UNFAIR!

Amazon Rekognition failures in facial recognistion correlate with skin color



**98.7%**  **68.6%**  **100%**  **92.9%**

DARKER MALES    DARKER FEMALES    LIGHTER MALES    LIGHTER FEMALES

https://medium.com/@Joy.Buolamwini/response-racial-and-gender-bias-in-amazon-rekognition-commercial-ai-system-for-analyzing-faces-a289222eeced

*People worry that the computers will get too smart and take over the world, but the real problem is that they're too stupid and they've already taken over the world.*

Pedro Domingos "The Master Algorithm"

All is well?

*No.*

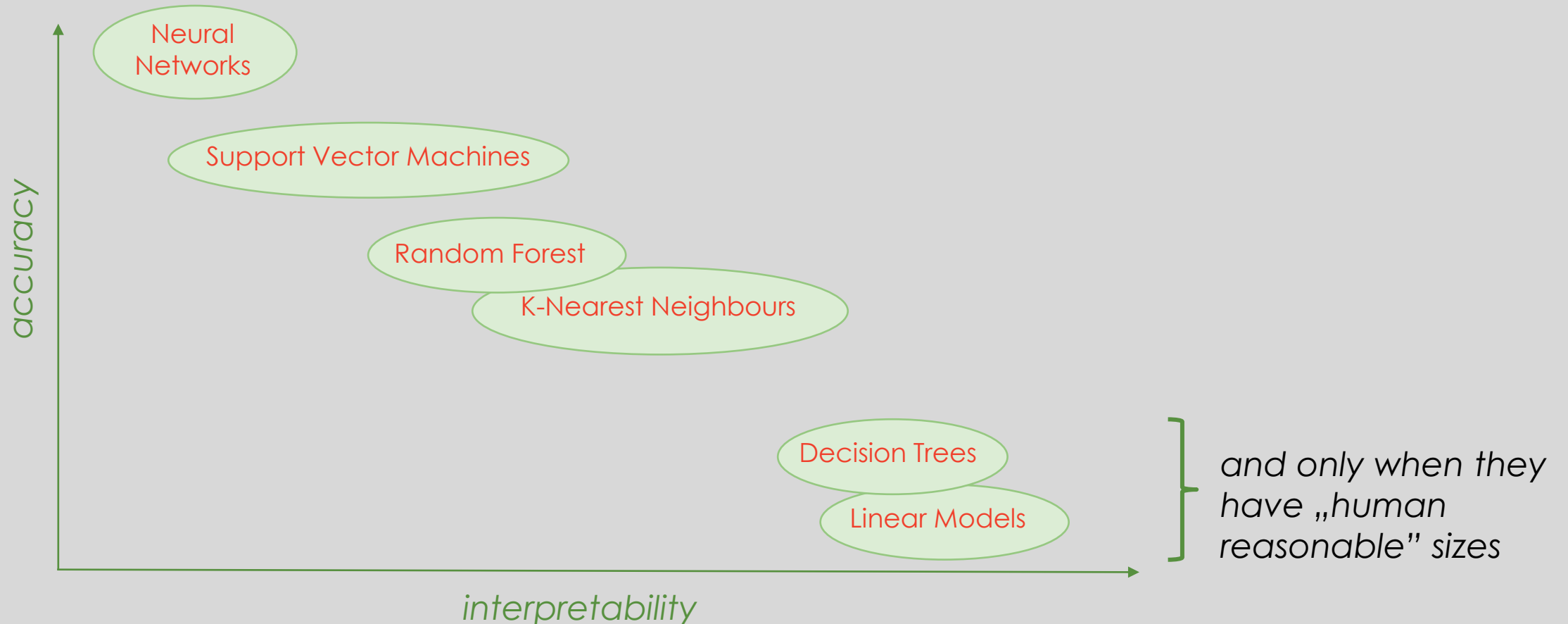We need **interpretability**.

# Some definitions

**Black boxes**
Systems that hide their internal logic to the user (either the internals are unknown or uninterpretable to humans)
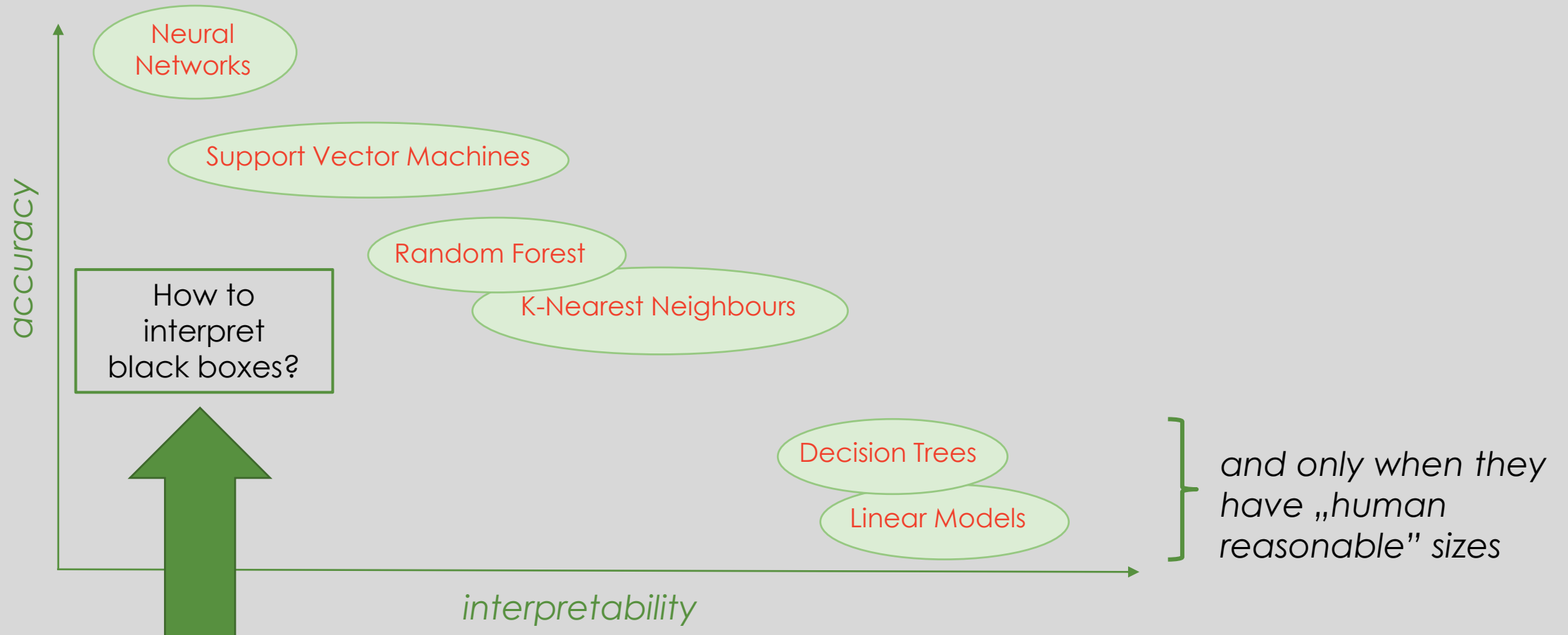
**Interpretability (also: comprehensibility)**
Ability to explain or to present in understandable terms to a human

Guidotti et al. 2018. *A Survey of Methods for Explaining Black Boxes Models*. ACM Comput. Surv. 51, 5, art. 93

# Trade-off between complexity and interpretability



*and only when they have „human reasonable" sizes*

# Trade-off between complexity and interpretability

# OVERVIEW OF INTERPRETABILITY METHODS

## Possible approaches

methods assigning meaning to individual **model components**

methods analyzing model predictions when **data is perturbed**

surrogate approach where the model is approximated by a simpler, more interpretable **surrogate model**

Molnar, Casalicchio, Bischl (2020) arXiv:2010.09337

# Possible approaches

methods assigning meaning to individual **model components** *(always model-specific)*

methods analyzing model predictions when **data is perturbed** *(mostly model-agnostic)*

surrogate approach where the model is approximated by a simpler, more interpretable **surrogate model**

methods

model-specific

model-agnostic

Molnar, Casalicchio, Bischl (2020) arXiv:2010.09337

# EXAMPLES

subjective choice of mine!

# Feature visualisation

**Dataset Examples** show us what neurons respond to in practice

**Optimization** isolates the causes of behavior from mere correlations. A neuron may not be detecting what you initially thought.

Baseball—or stripes?
*mixed4a, Unit 6*

Animal faces—or snouts?
*mixed4a, Unit 240*
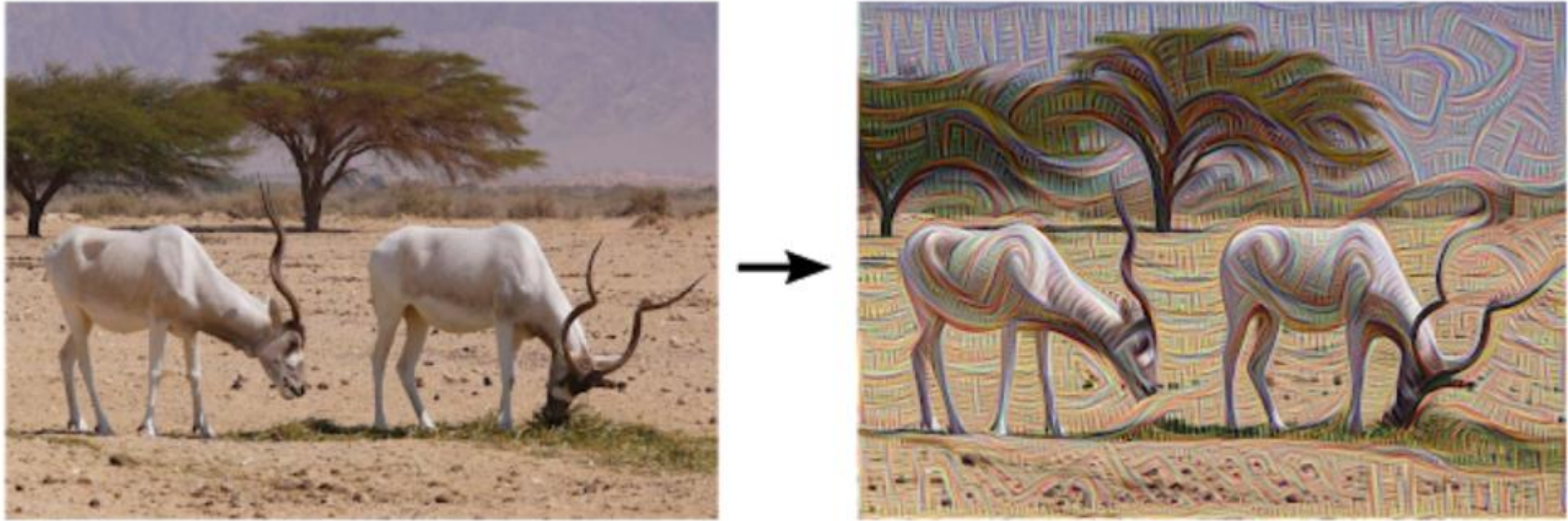
Clouds—or fluffiness?
*mixed4a, Unit 453*

Buildings—or sky?
*mixed4a, Unit 492*

https://distill.pub/2017/feature-visualization/
https://distill.pub/2017/feature-visualization/appendix/

# We can look at whole layers!
Deep Dream

enhance what was detected by a chosen layer and visualise



Left: Original photo by Zachi Evenor. Right: processed by Günther Noack, Software Engineer

# We can look at whole layers!
## Deep Dream

enhance what was
detected by a chosen layer
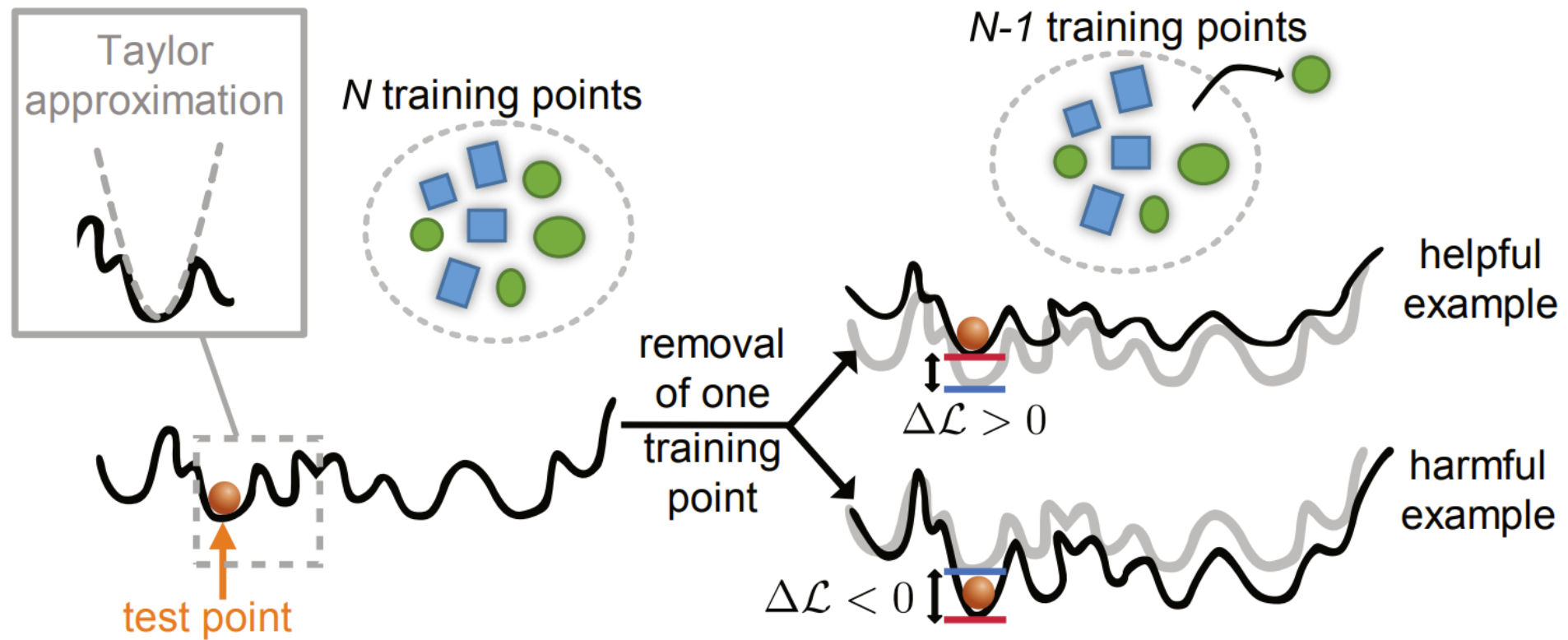and visualise
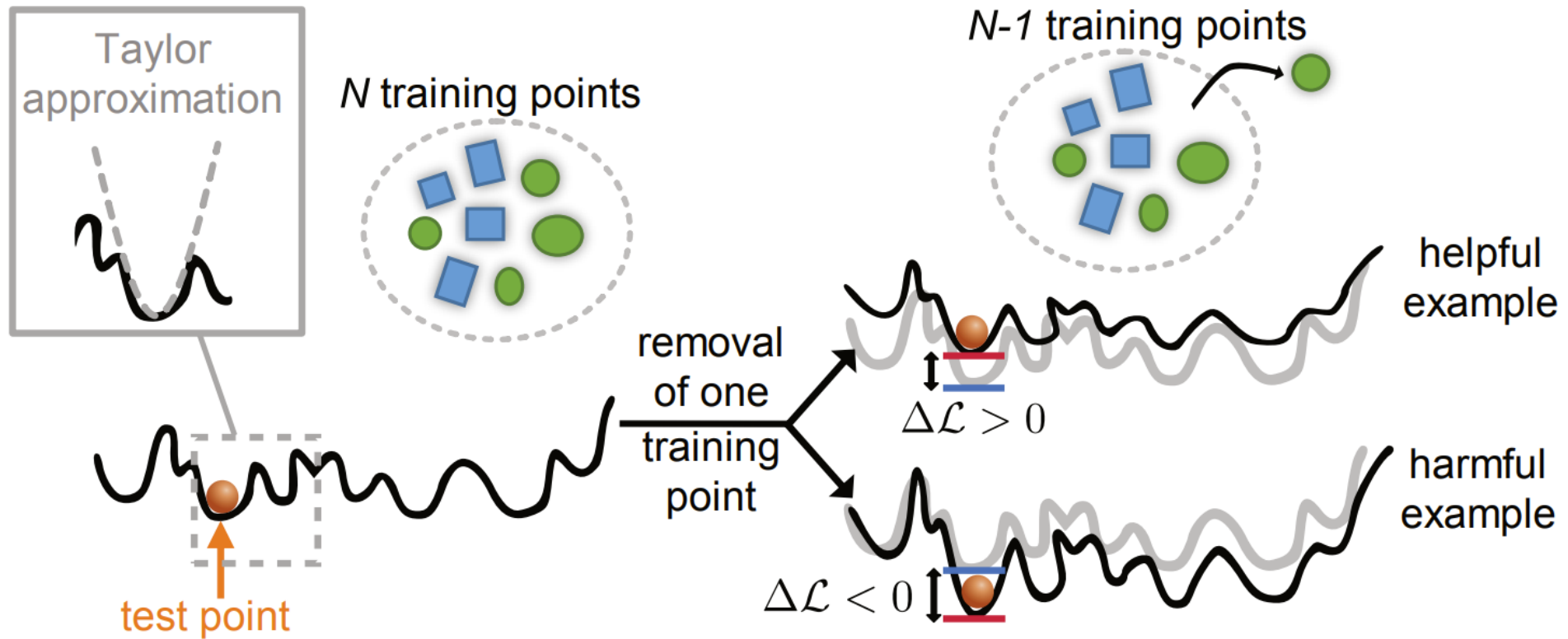


"Admiral Dog!"    "The Pig-Snail"    "The Camel-Bird"    "The Dog-Fish"

https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html
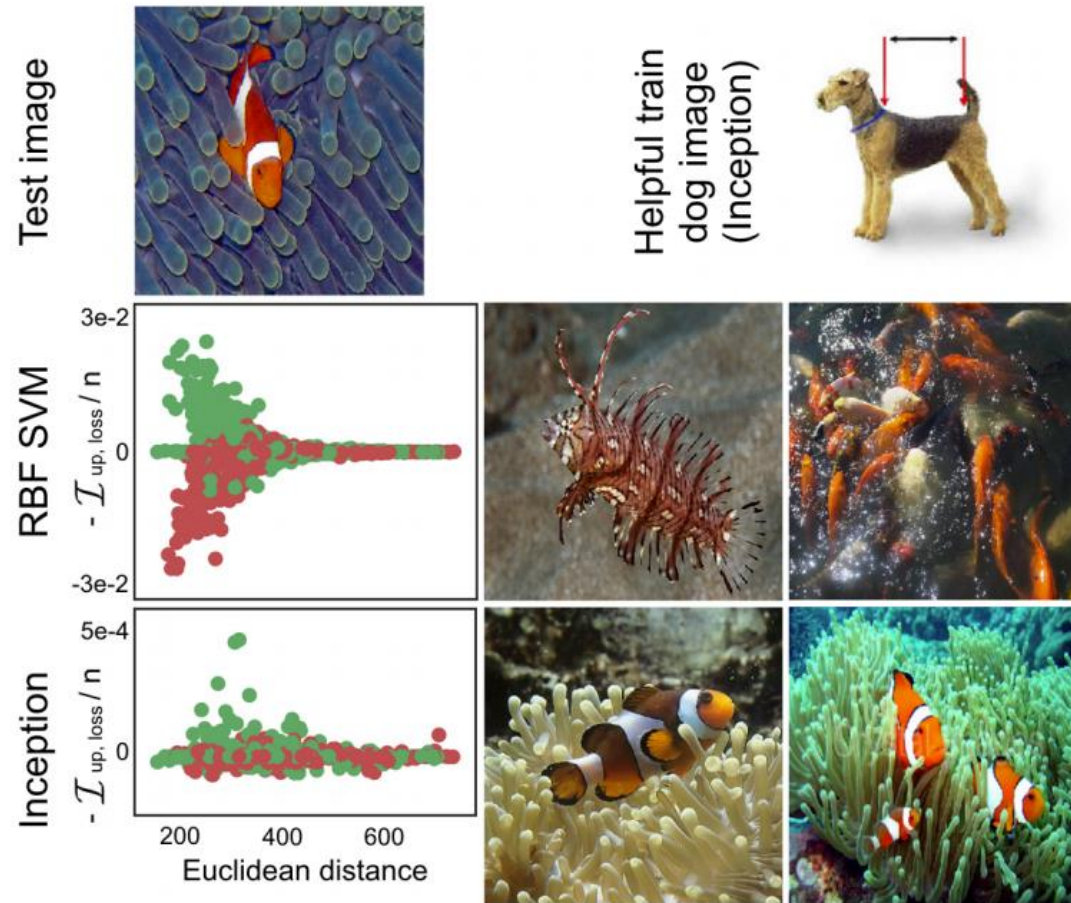
# Leave-one-out training

# Leave-one-out training

Figure 4. **Inception vs. RBF SVM. Bottom left:** $-\mathcal{I}_{up,loss}(z, z_{test})$ vs. $\|z - z_{test}\|_2^2$. Green dots are fish and red dots are dogs. **Bottom right:** The two most helpful training images, for each model, on the test. **Top right:** An image of a dog in the training set that helped the Inception model correctly classify the test image as a fish.

Its approximation are:
**influence functions**

*(They can be used to NNs after generalizing to non-convex problems)*

Koh & Liang (2017), arXiv:1703.04730v1

# LIME (Local Interpretable Model-Agnostic Explanations)



Model — Data and Prediction — Explainer (LIME) — Explanation — Human makes decision

# LIME (Local Interpretable Model-Agnostic Explanations)



Simplified LIME algorithm
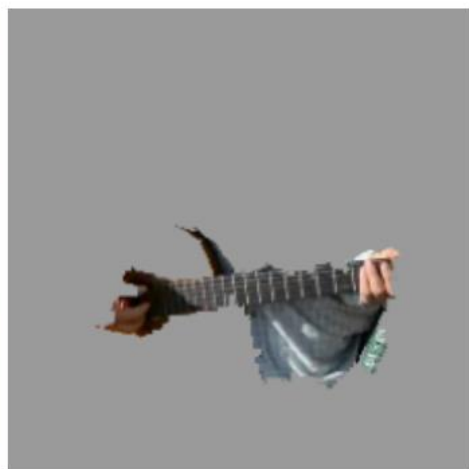
(a) Husky classified as wolf    (b) Explanation

Figure 11: Raw data and explanation of a bad model's prediction in the "Husky vs Wolf" task.
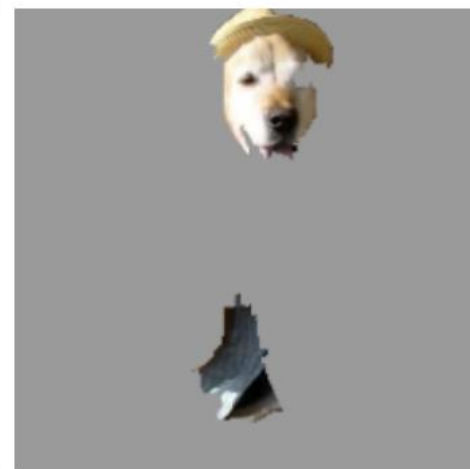


(a) Original Image    (b) Explaining *Electric guitar*    (c) Explaining *Acoustic guitar*    (d) Explaining *Labrador*

Figure 4: Explaining an image classification prediction made by Google's Inception neural network. The top 3 classes predicted are "Electric Guitar" ($p = 0.32$), "Acoustic guitar" ($p = 0.24$) and "Labrador" ($p = 0.21$)

# TUTORIAL EXAMPLE

Class Activation Map (CAM)

**Class Activation Mapping**

$$w_1 * \quad + \quad w_2 * \quad + \ldots + \quad w_n * \quad = \quad$$

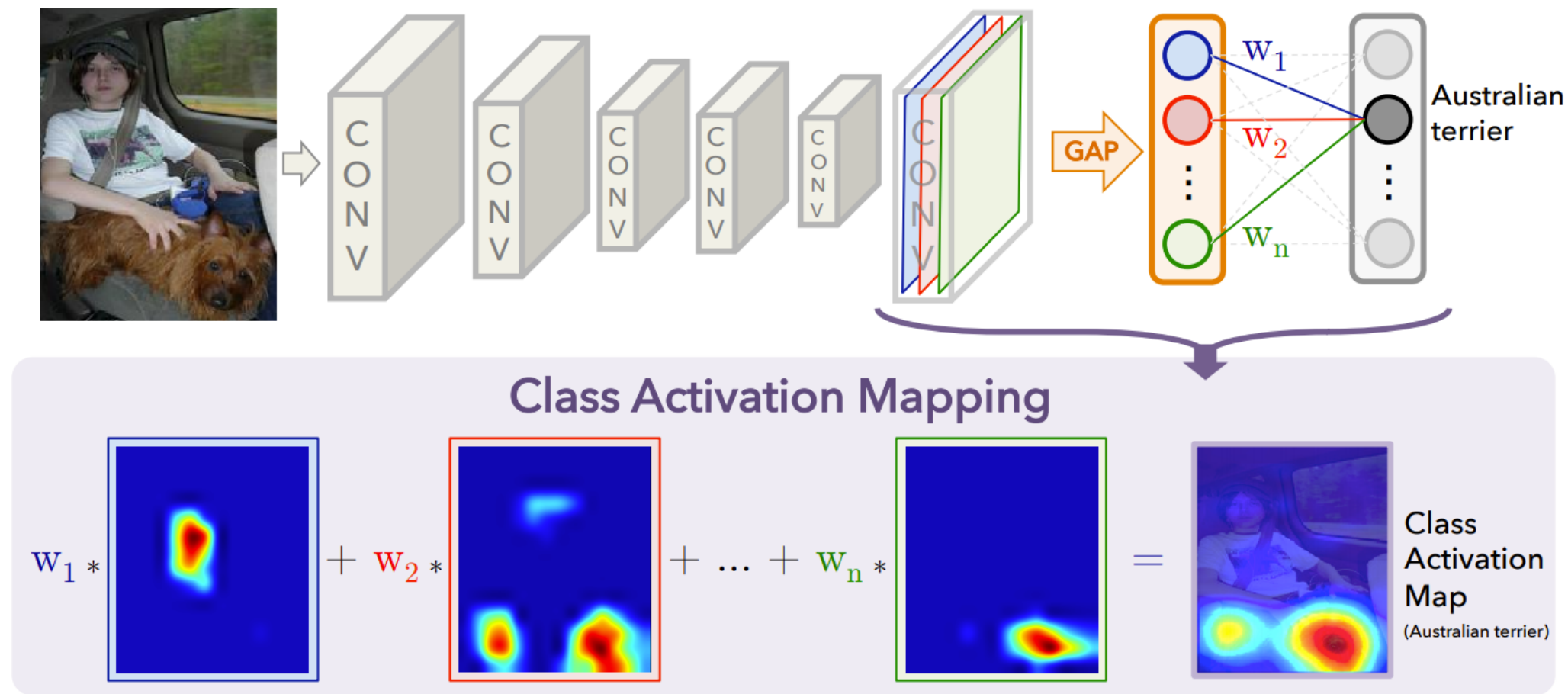Class Activation Map (Australian terrier)

Figure 2. Class Activation Mapping: the predicted class score is mapped back to the previous convolutional layer to generate the class activation maps (CAMs). The CAM highlights the class-specific discriminative regions.

Zhou et al. (2015), arXiv:1512.04150v1

Brushing teeth    Cutting trees

Figure 1. A simple modification of the global average pooling layer combined with our class activation mapping (CAM) technique allows the classification-trained CNN to both classify the image and localize class-specific image regions in a single forward-pass e.g., the toothbrush for *brushing teeth* and the chainsaw for *cutting trees*.

# Take-home messages

**ML needs to be interpretable if it is to be trusted**

Is it reliable? Is it fair? Does it generalizes well? Can we extract the learnt knowledge?

**Interpretability methods**

Some look at the components of the model (CAM, feature visualisation), some are based on perturbing the data and looking at how it changed the model (LOO training), there are also surrogate approaches (LIME).

**Class Activation Map (CAM)**

Shows which parts of the image were decisive for the prediction of the given class.

**Many other methods exist!**