

Challenges in explaining deep DIRECT learning models for data with biological variations



Lenka Tětková¹, Erik Schou Dreier², Robin Malm² and Lars Kai Hansen¹

¹ DTU Compute, Technical University of Denmark

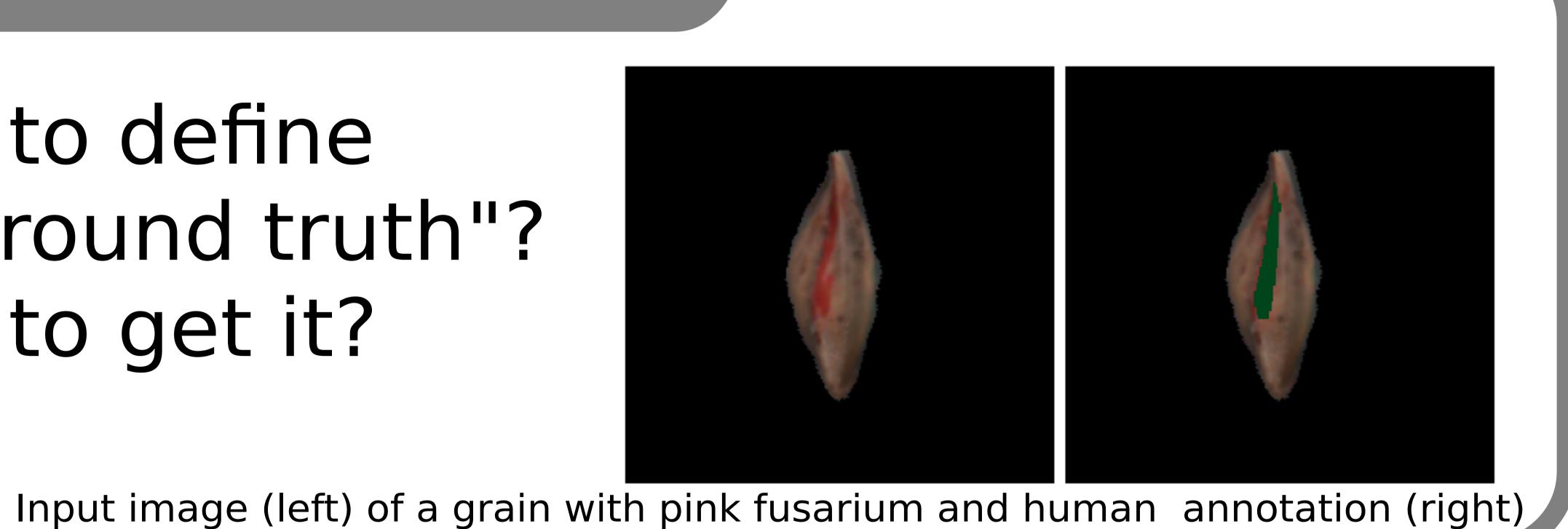


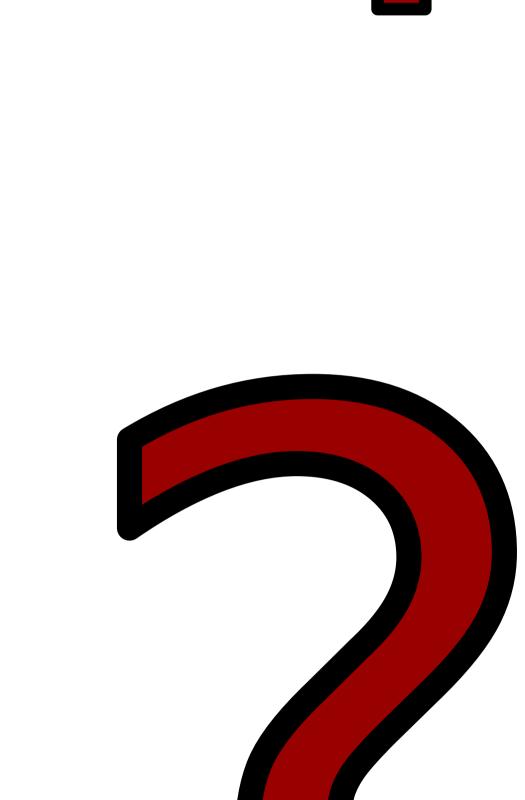
Motivation and research qustions

What explainability method is the best?

Can post-hoc explainability methods be directly applied to other of data?

How to define "ground truth"? How to get it?







- Methods

Data: images of grains Binary classification of presence of a disease (pink fusarium) or damage (broken).

Small convolutional network

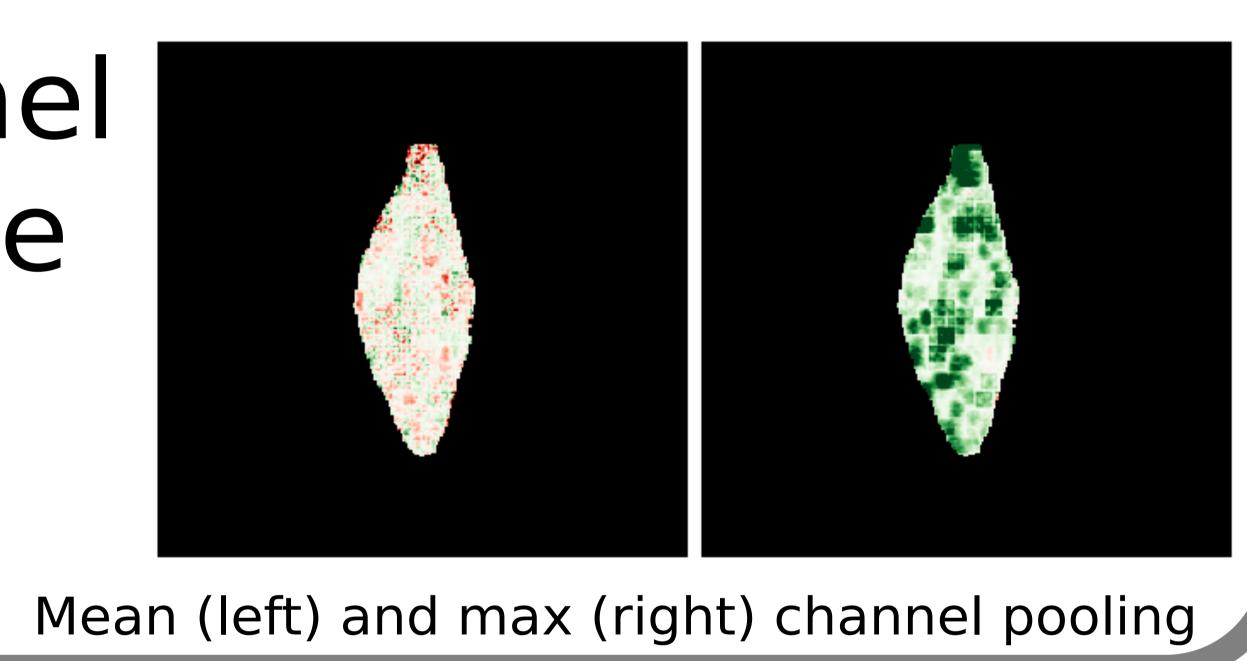
Evaluation

What is a good explanation? Given two explanations, which one is

better?

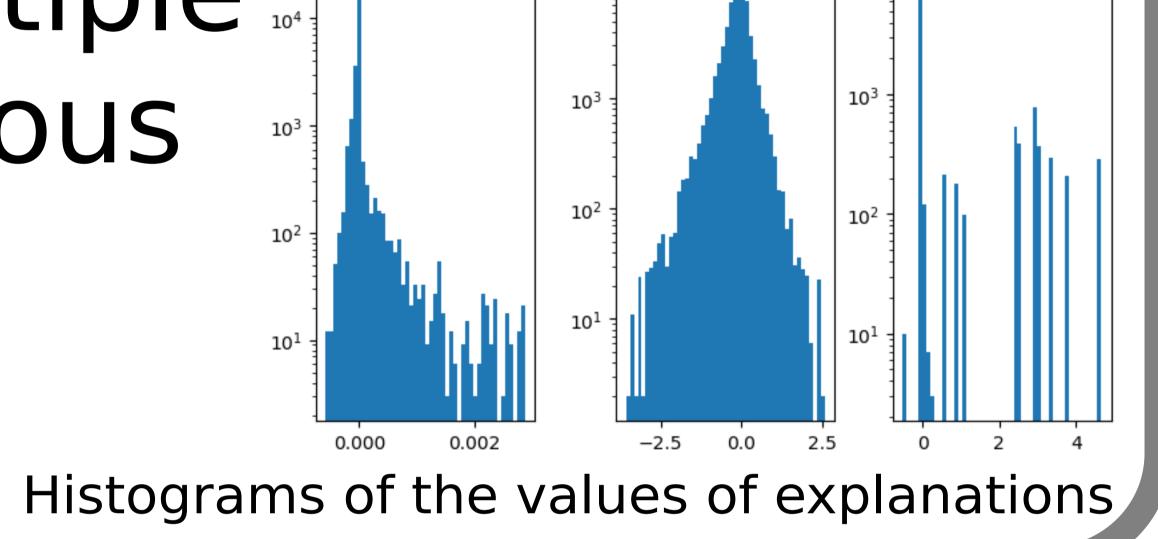
Channel pooling :

How to pool 3-channel explanations into one channel? (for visualization)



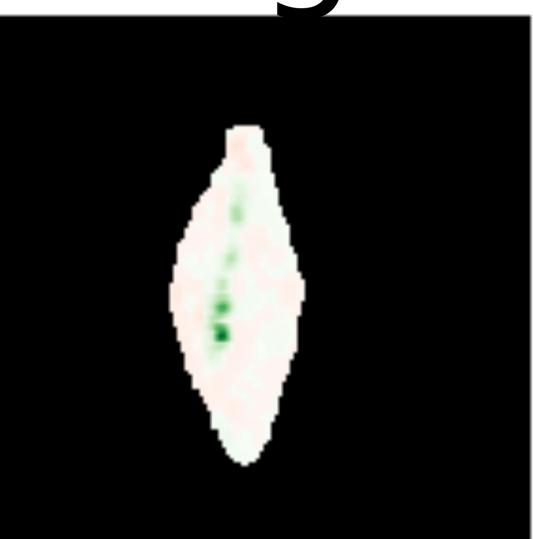
Aggregation

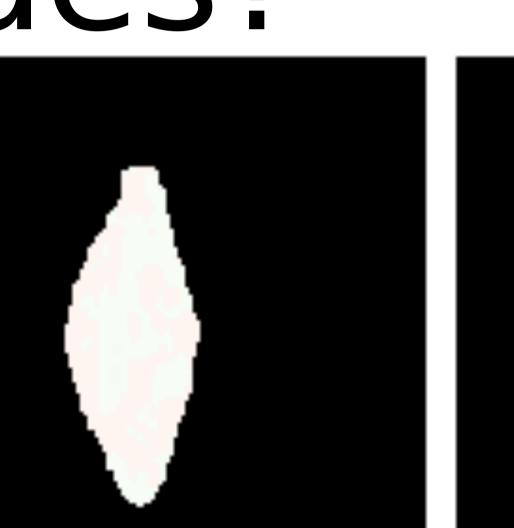
How to aggregate multiple 104 explanations with various magnitudes?

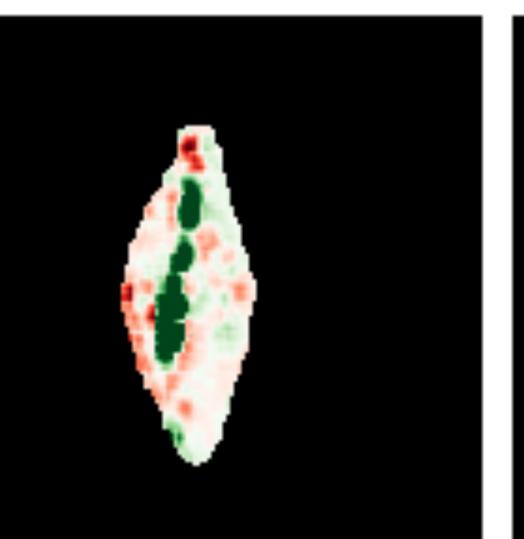


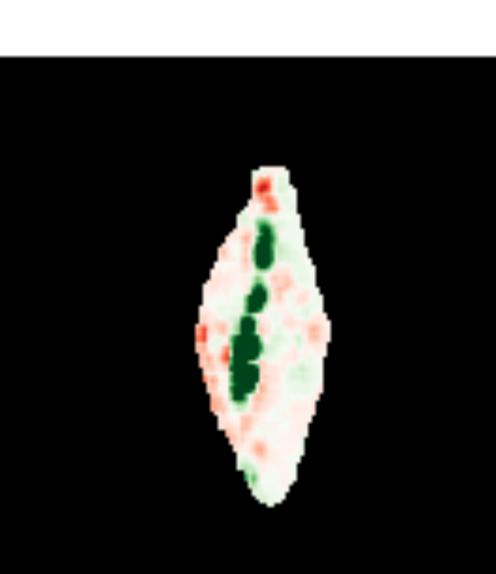
Visualization

How to visualize and compare explanations with various magnitudes?





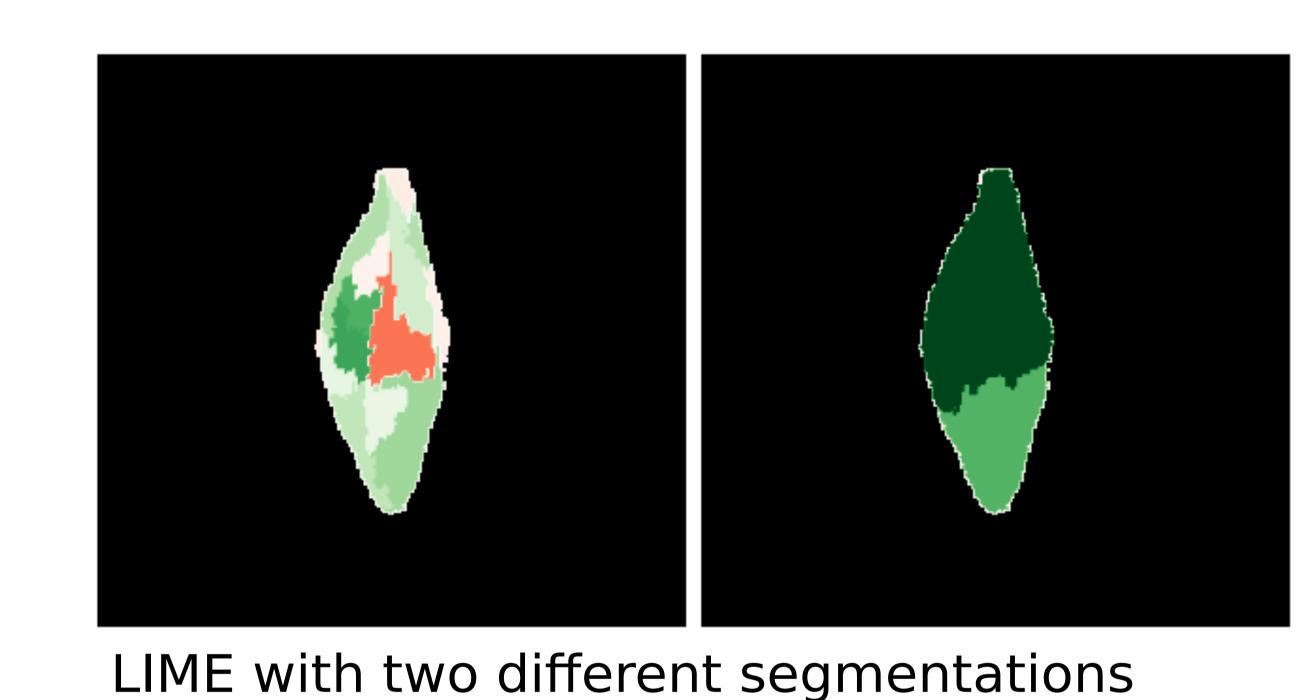


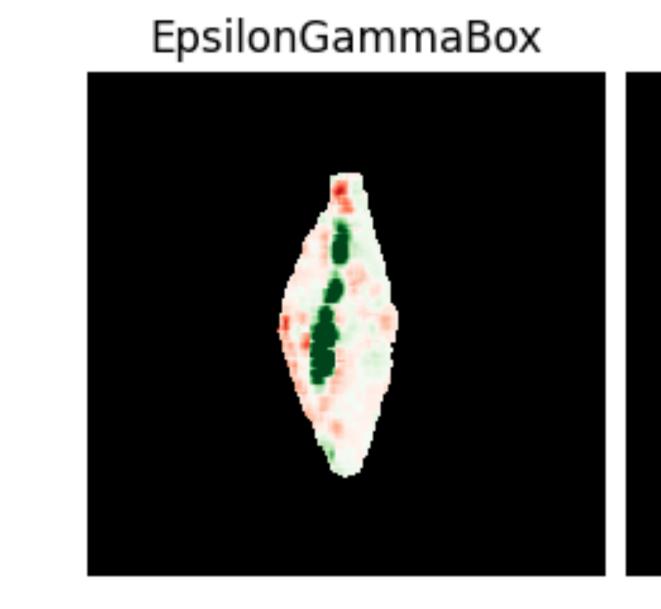


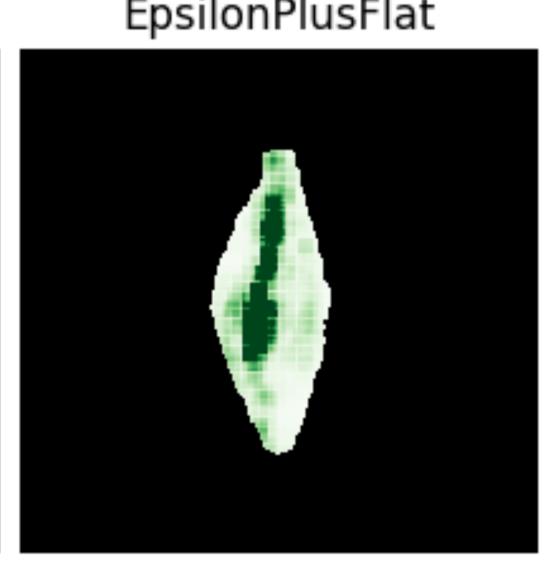
Multiple ways of normalizing the same explanation

Hyperparameters

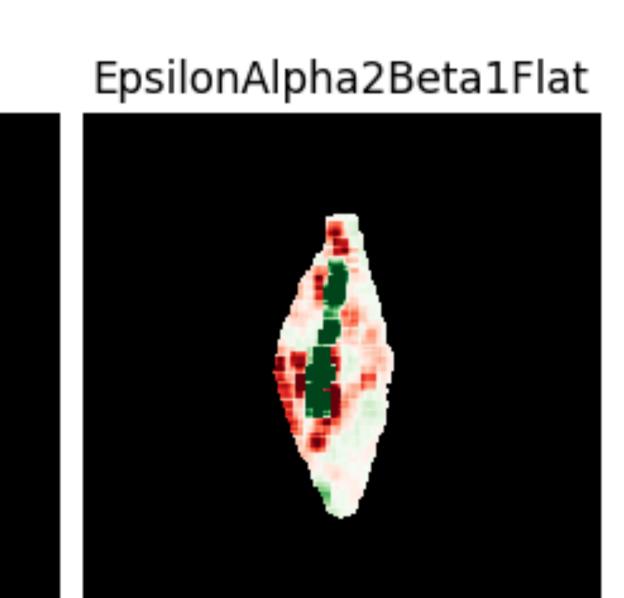
Many choices to make with big impact on the explanations.







Variants of LRP



Experiments

- Robustness to data augmentations [1]
- Quality of explanations
 - Pixel-flipping [2], IROF [3], sensitivity [4], complexity [5]
- Similarity to the ground truth
 - ROC-AUC, Relevance Mass Accuracy [6]

- Robustness to data augmentation is similar to ImageNet.
- Better robustness to invariant augmentations (e.g., changes in brightness) than equivariant augmentations (e.g., rotation).
- Some methods are more robust to channel pooling than others.
- Explanations are quite aligned with the human annotations.
- No explainability method is good at all the aspects.
- Mean of explanations performs better than a single method.

lake-away

Many choices to make with big impacts.

Aloout me

PhD student

