

Manipulating explainability and fairness in machine learning

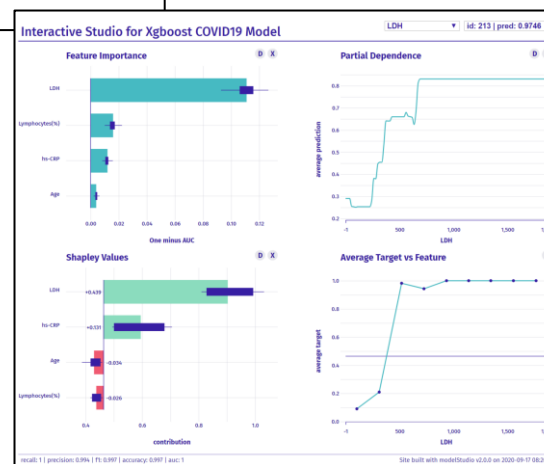
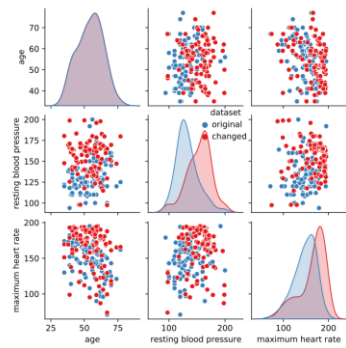
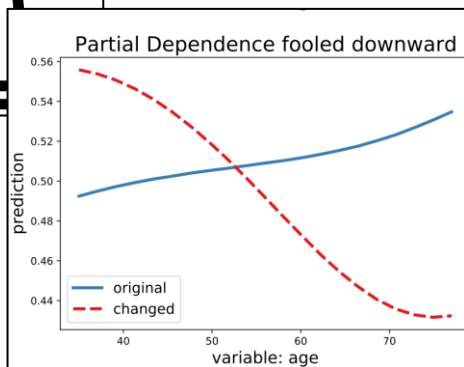
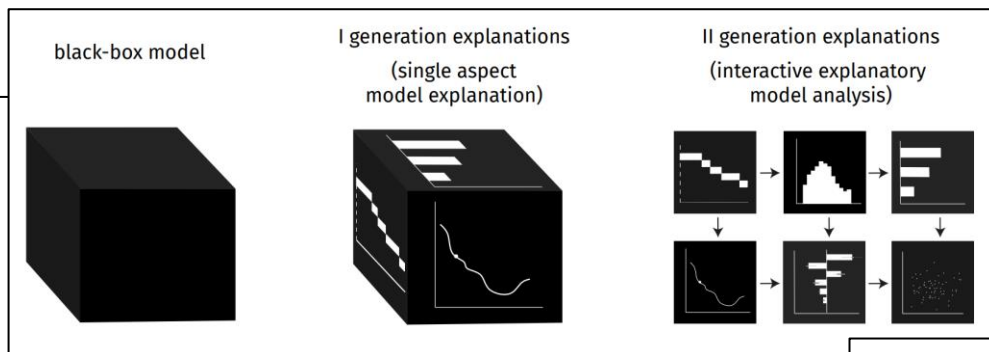
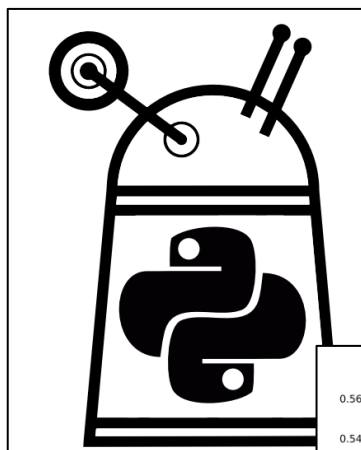
Hubert Baniecki
@ MI2DataLab, Warsaw University of Technology

ML in PL Conference, Poland
November 2021



```
hbaniecki:~$ whoami
```

- Reseracher and Data Science student at Warsaw University of Technology
- Interested in explainable machine learning, developing open source software
- Specifically joint with adversarial machine learning, and evaluation



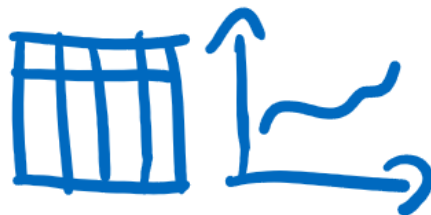
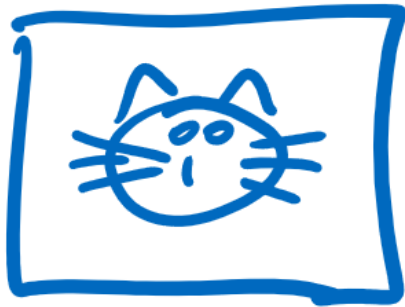
DISCLAIMER

1. Presenting interesting work of others (fair use)
2. Selective survey based on a live list of related work:
<https://github.com/hbaniecki/adversarial-explainable-ai> (contributions are welcomed)
3. Omitting technical details and math

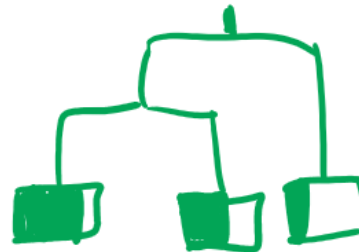
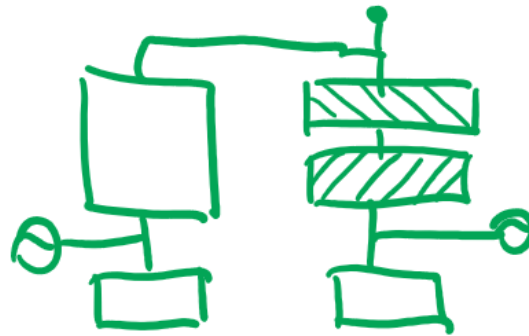
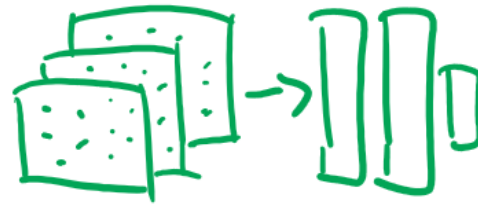


explainability and fairness

DATA



MODEL

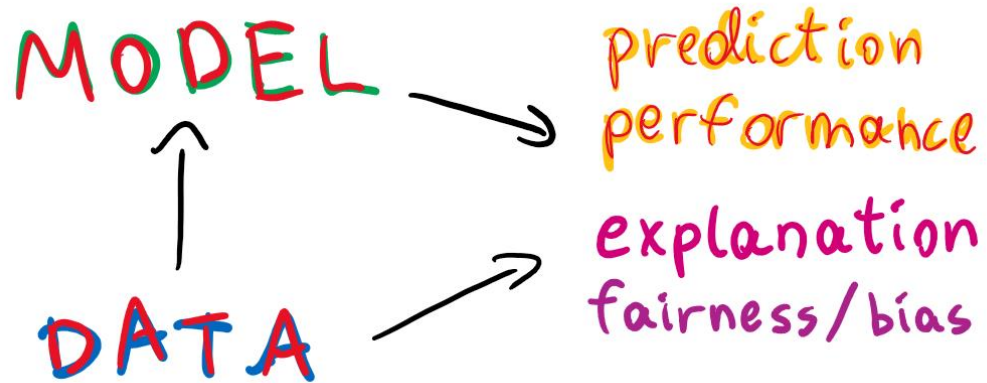


OUTCOME

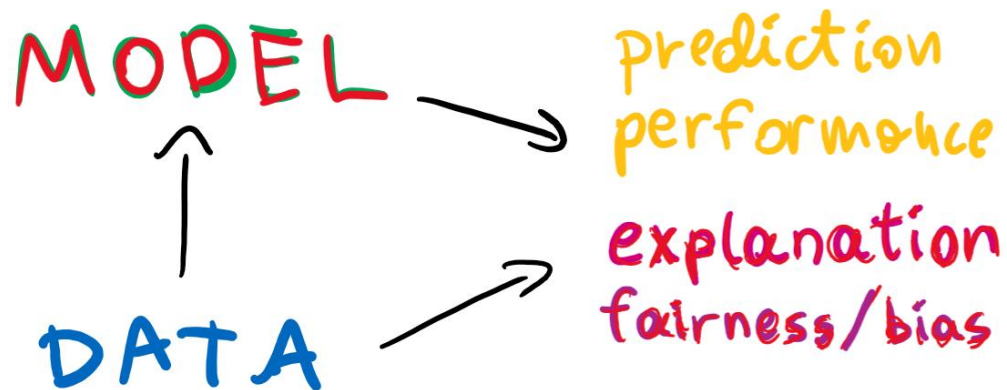
prediction
performance
decision
discovery
automation
forecast

explanation
fairness/bias

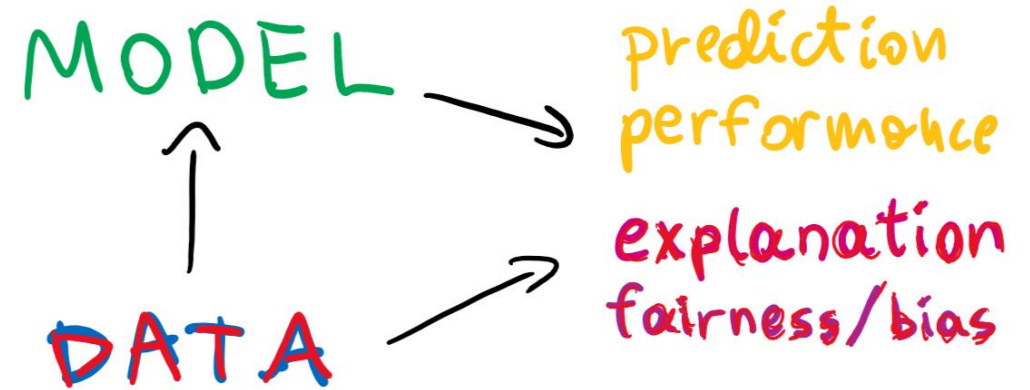
Manipulating explainability and fairness



1. change the model



2. change the data



Manipulating explainability via data change

Target: change explanation/saliency maps for deep neural network image classification

Method: perturb an image with gradient optimization; aim for an arbitrary target map

Loss ~

*distance(manipulated explanation, target explanation) + γ * distance(manipulated prediction, original prediction)*

Result: It is possible to manipulate explanations (ImageNet + VGG, ResNet, DenseNet)

Explanations can be manipulated and geometry is to blame

Ann-Kathrin Dombrowski¹, Maximilian Alber⁵, Christopher J. Anders¹,
Marcel Ackermann², Klaus-Robert Müller^{1,3,4}, Pan Kessel¹

¹Machine Learning Group, Technische Universität Berlin, Germany

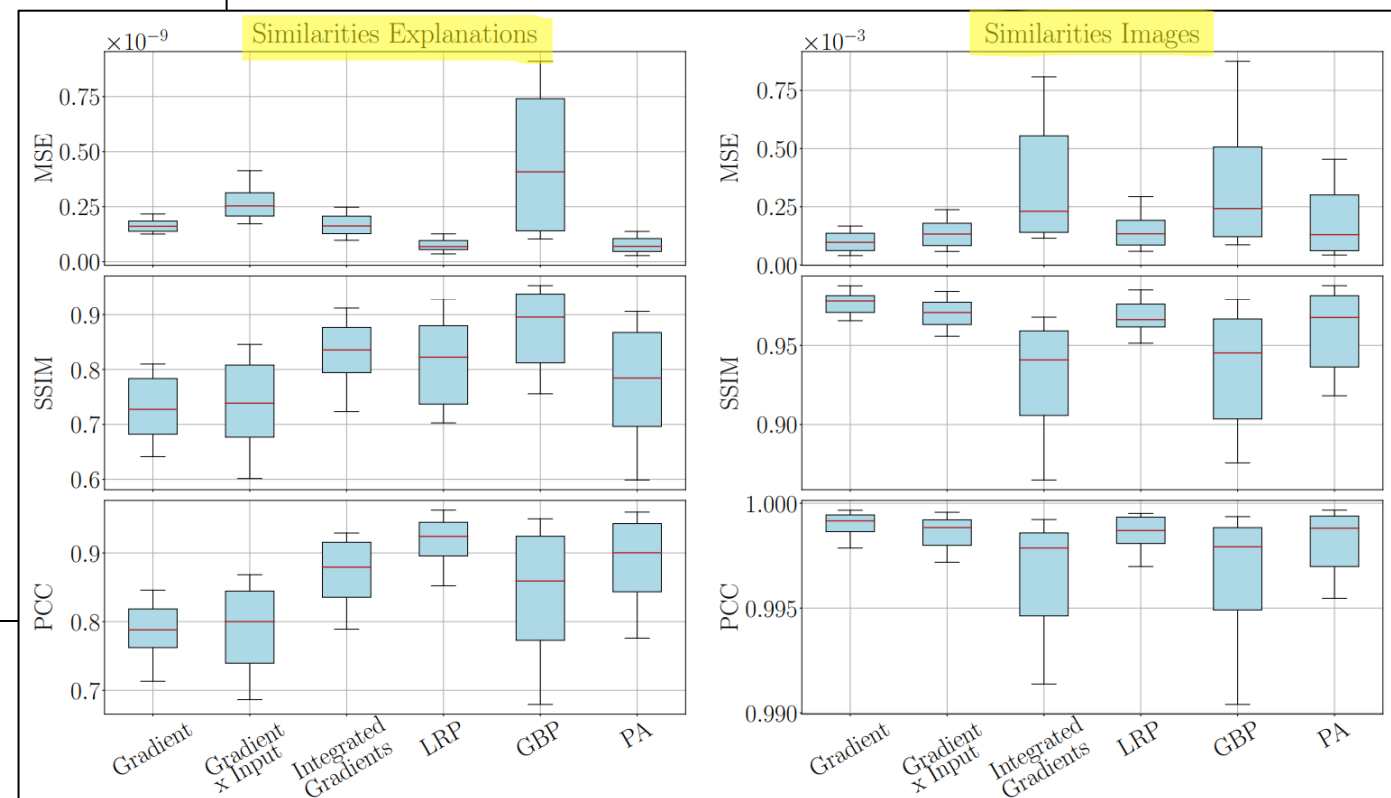
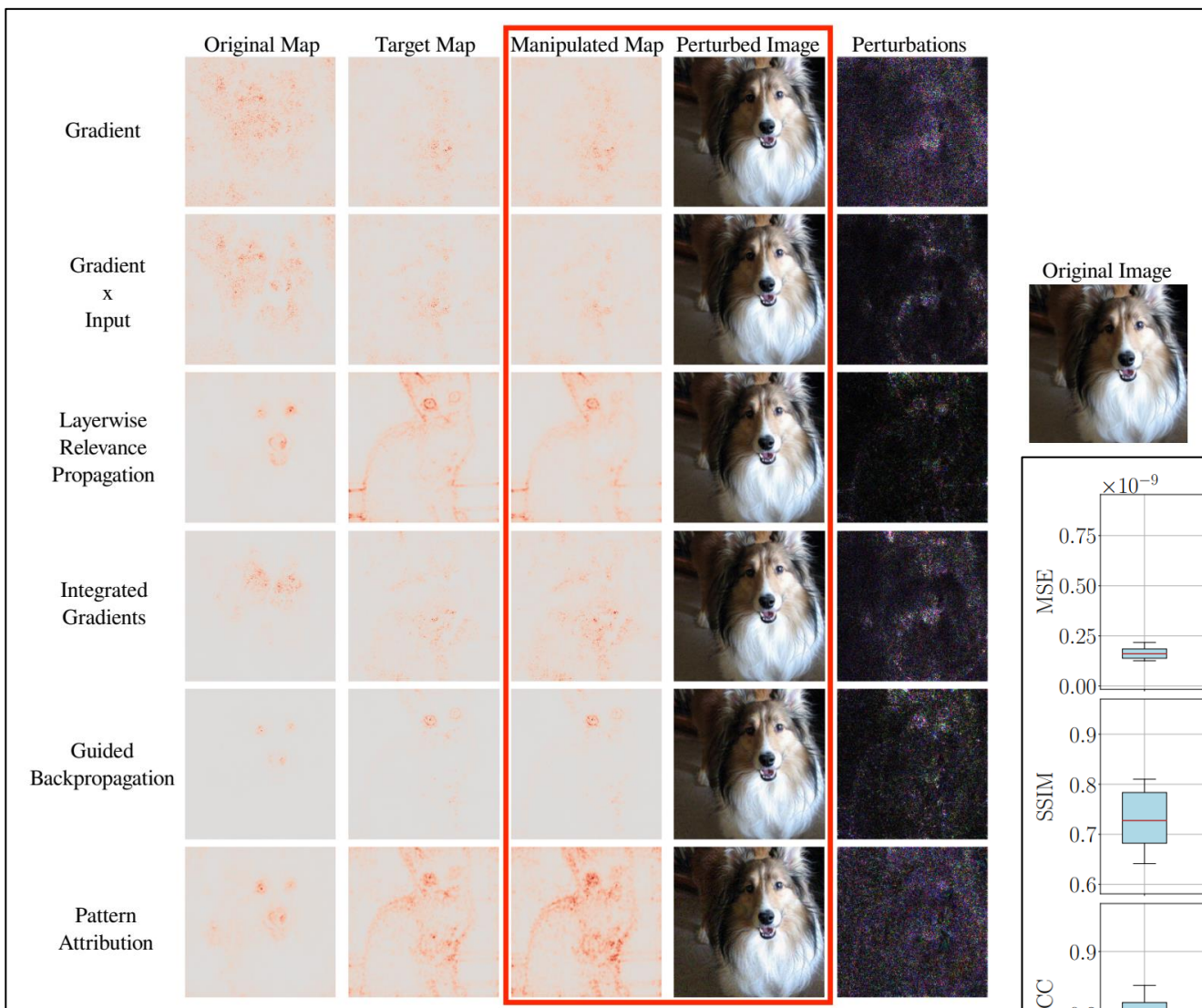
²Department of Video Coding & Analytics, Fraunhofer Heinrich-Hertz-Institute, Berlin, Germany

³Max-Planck-Institut für Informatik, Saarbrücken, Germany

⁴Department of Brain and Cognitive Engineering, Korea University, Seoul, Korea

⁵Charité Berlin, Berlin, Germany





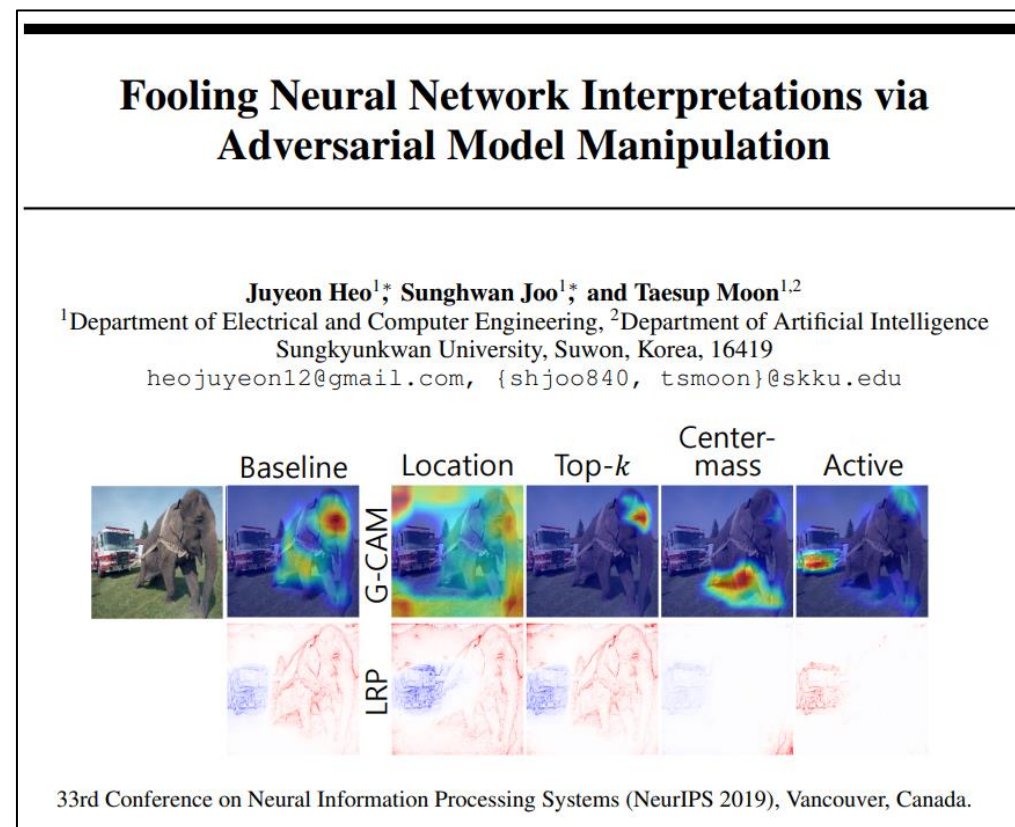
Manipulating explainability via model change

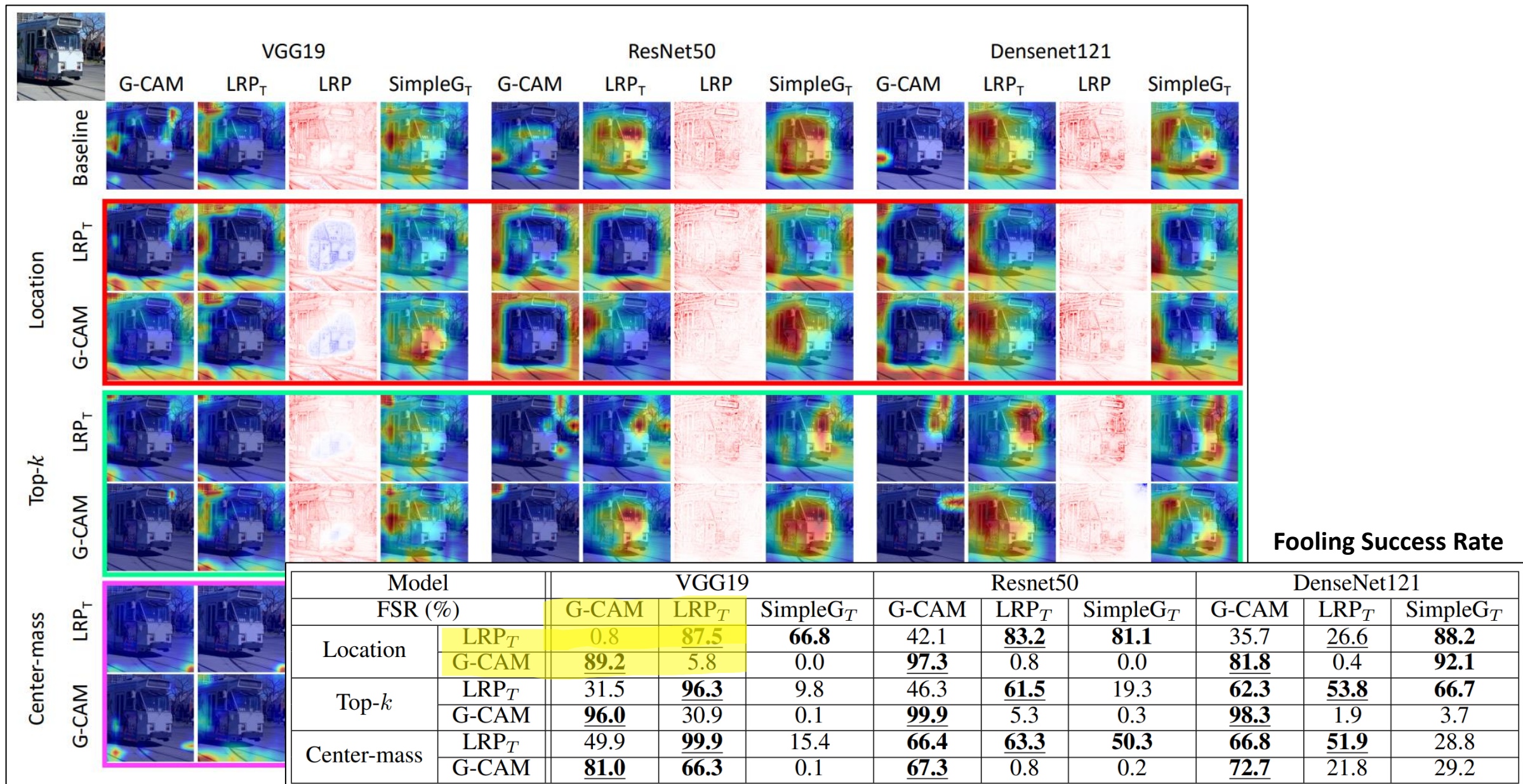
Target: change explanation/saliency maps for deep neural network image classification

Method: fine-tune a model; change parameters to generate uninformative or false explanations

Loss \sim
performance of the manipulated model +
 $\lambda * \text{distance}(\text{manipulated explanations}, \text{target explanations})$

Result: It is possible to manipulate explanations - globally (ImageNet + VGG, ResNet, DenseNet)





2018– evaluate explanation maps

Ancona et al. **Towards better understanding of gradient-based attribution methods for Deep Neural Networks.**

International Conference on Learning Representations (ICLR). 2018.

Theoretical unification of Grad*Input, IG, LRP & DeepLIFT methods and evaluating them with an introduced *sensitivity-n* property.

Alvarez-Melis & Jaakkola. **Towards Robust Interpretability with Self-Explaining Neural Networks.**

Neural Information Processing Systems (NeurIPS). 2018.

Introduce a *self-explaining neural network* with native concept-based interpretability and evaluating its faithfulness and stability against the explanation maps.

Adebayo et al. **Sanity Checks for Saliency Maps.**

Neural Information Processing Systems (NeurIPS). 2018.

Develop model randomization and data *randomization tests* to evaluate explanations.

2018~ evaluate explanation maps

Kindermans et al. **The (Un)reliability of Saliency Methods.**

Explainable AI: Interpreting, Explaining and Visualizing Deep Learning. Springer. 2019.

Analyze *input invariance* of saliency maps and the choice of the reference point.

Ghorbani et al. **Interpretation of Neural Networks Is Fragile.**

AAAI Conference on Artificial Intelligence (AAAI). 2019.

Adversarial attack on explanations via gradient-based data perturbations.
(See Domrowski et al. 2019 for differences)

Manipulating explainability and fairness via model change

Slack et al. **Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods.**

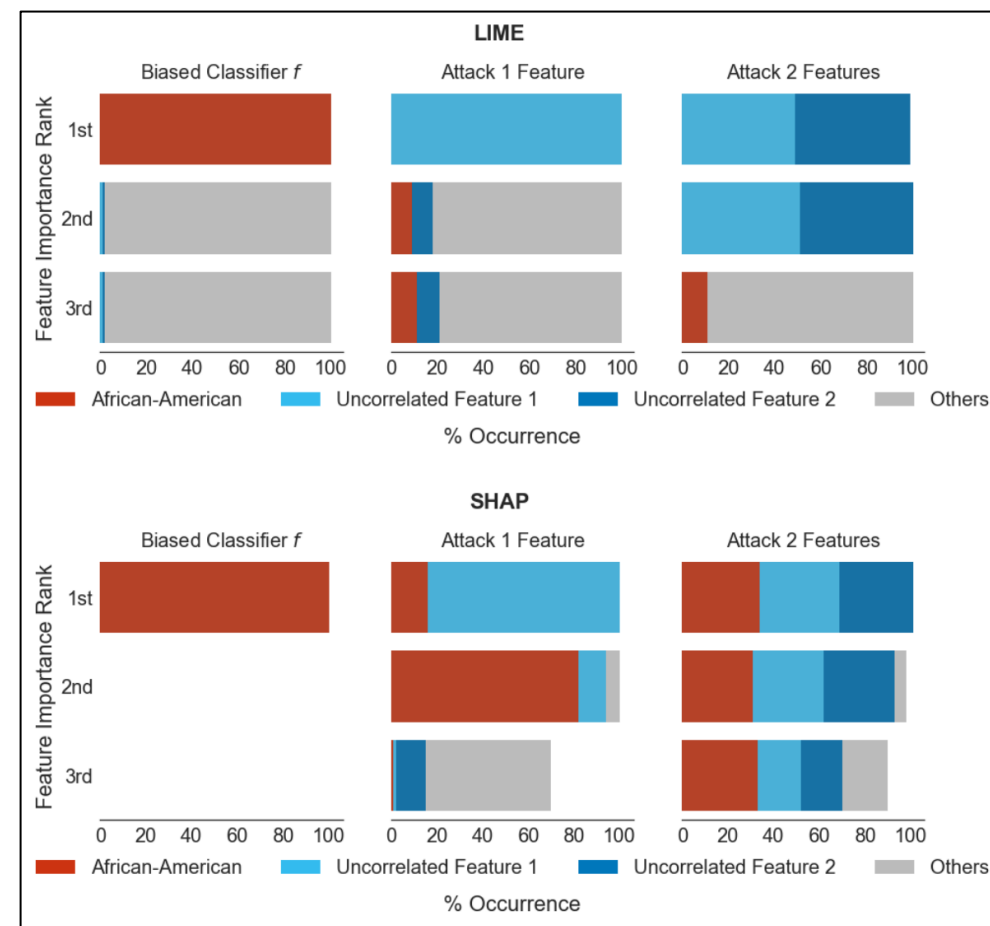
AAAI/ACM Conference on AI, Ethics, and Society (AIES). 2020.

Target: change feature attributions of a (biased) black-box model

Framework:

1. build a surrogate model that predicts the same as black-box in-distribution, but arbitrarily out-of-distribution (as needed)
2. support it with another classifier of out-of-distribution samples
3. black-box predictions in-distribution don't change, but LIME and SHAP attributions change as they use perturbed data

Result: Fooled LIME and SHAP attributions



Manipulating fairness via model change

Aivodji et al. **Fairwashing: the risk of rationalization.**
International Conference on Machine Learning (ICML). 2019.

Target: change the value of fairness parity measure of a black-box

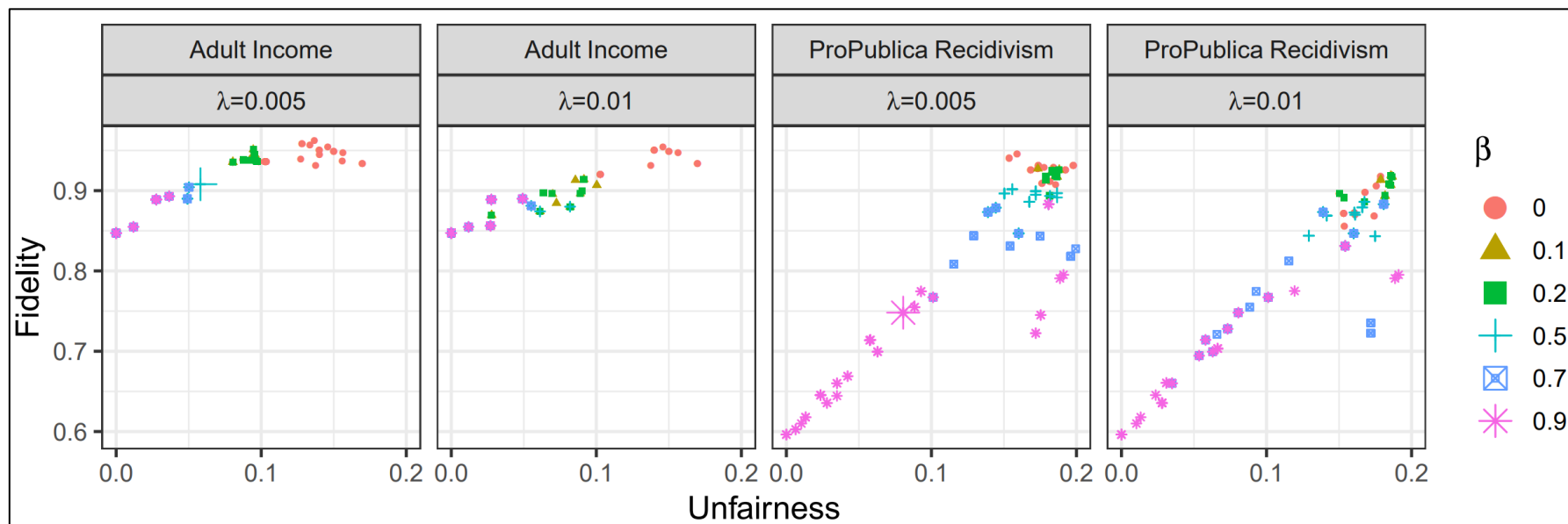
Method: build a fair surrogate model that predicts the same as black-box

$$\text{Loss} \sim (1 - \beta) * (1 - \text{fidelity}) + \beta * \text{unfairness} + \lambda * \text{complexity}$$

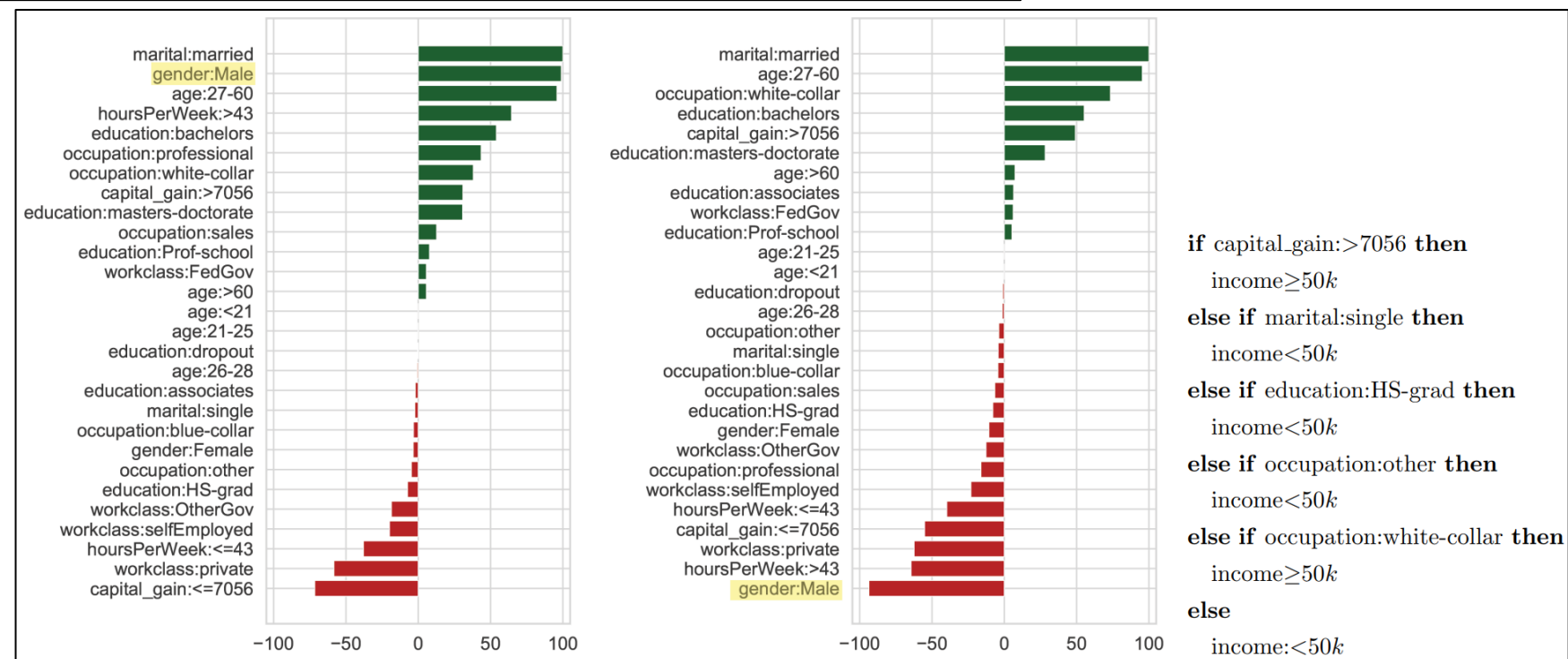
1. fidelity: accuracy of predicting the same outcomes as the black-box
2. unfairness: demographic parity measure of bias with respect to the sensitive attribute, e.g. gender or race
3. complexity: a number of rules in a list (serves as a regularization)

In practice: a classification rule list approximates a random forest model

Result: A fairer surrogate model of high fidelity



Relative feature dependence



Manipulating fairness via data change

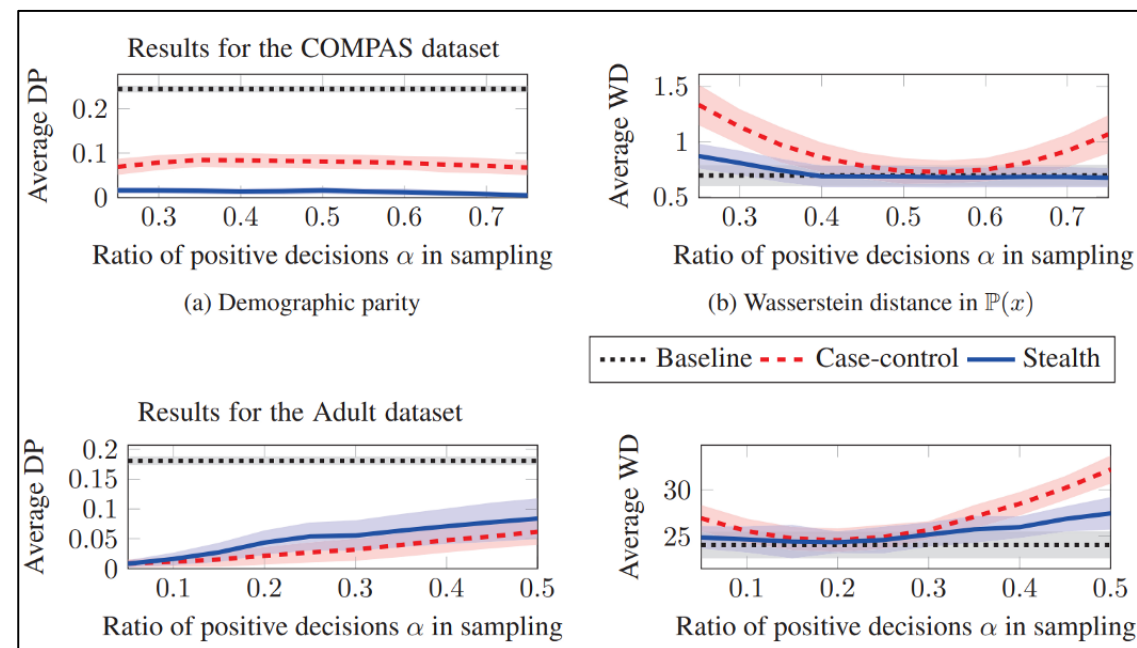
Fukuchi et al. **Faking Fairness via Stealthily Biased Sampling.**
AAAI Conference on Artificial Intelligence (AAAI). 2020.

Target: change the value of fairness parity measure of a black-box

Method: subset a sample of dataset on which a model *appears to* be fair; minimize:

1. similarity of data distributions, specifically Wasserstein distance
2. unfairness: demographic parity measure of bias with respect to the sensitive attribute, e.g. gender or race

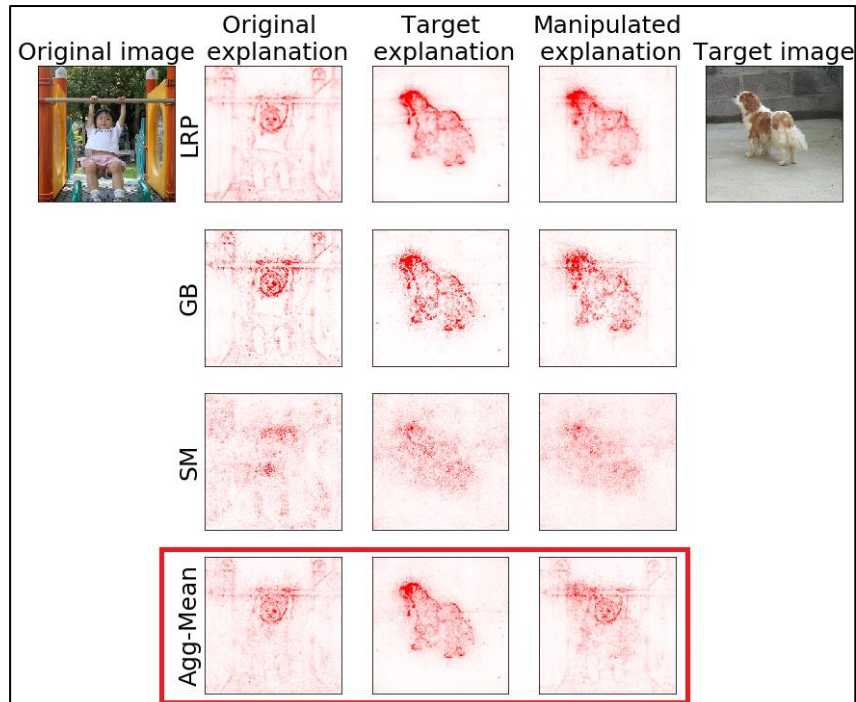
Result: A benchmark dataset, which proves model fairness



Summary & future work

- It is possible to manipulate explanations and fairness measures
- Evaluation approaches are required for a trustworthy adoption of these methods
- Develop defense mechanisms

Vilone & Longo. **Notions of explainability and evaluation approaches for explainable artificial intelligence.** *Information Fusion*. 2021.



Rieger & Hansen. **A simple defense against adversarial attacks on heatmap explanations.** *Workshop on Human Interpretability in Machine Learning (ICML WHI)*. 2020.

Summary & future work (2022)

Investigate critical vulnerabilities in novel explanations

Brittle interpretations: The Vulnerability of TCAV and Other Concept-based Explainability Tools to Adversarial Attack

Anonymous

29 Sept 2021 (modified: 06 Oct 2021) ICLR 2022 Conference Blind Submission Readers: 

Everyone [Show Bibtex](#) [Show Revisions](#)

Keywords: interpretability, adversarial attack

One-sentence Summary: We identify a novel vulnerability in the deep learning interpretability pipeline, and design attacks that mislead model explanations for two well known interpretability tools.

Develop robust explanations

Dombrowski et al. **Explanations can be manipulated and geometry is to blame.** *NeurIPS*. 2019.

Dombrowski et al. **Towards robust explanations for deep neural networks.** *Pattern Recognition*. 2022.

Feedback appreciated!

Contact: <https://www.linkedin.com/in/hbaniecki>

Resources: <https://github.com/hbaniecki/adversarial-explainable-ai>