# Manipulating Model Explanations: How to fool what tries to make sense of

Verena Heusser
verena.heusser@student.kit.edu
Karlsruhe Institute of Technology (KIT)

## ABSTRACT

This paper reviews state-of-the-art approaches to model explanations with a focus on those techniques that try to fool these methods.

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

## KEYWORDS

neural networks, model explanations, adversarial training

## 1  INTRODUCTION

## 2  METHODS AND INTERPRETERS

## 3  EXPERIMENTS

### 3.1  Models

### 3.2  Datasets

#### 3.2.1  *ImageNet.*

#### 3.2.2  *Recidivism Dataset.*

#### 3.2.3  *German dataset of..*

### 3.3  Measures

#### 3.3.1  *Fooling Success Rate.*

## 4  RESULTS