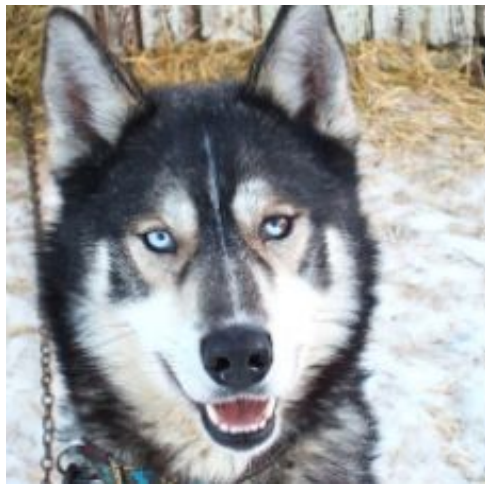# Why Should I Trust You?

**L**ocal **I**nterpretable **M**odel-agnostic **E**xplanations
**LIME**

"... machine learning models remain mostly **black boxes**."

**HUSKY** classified as **WOLF**

**THE WHY**



**LIME's Explanation**

# Explanations' Desiderata

**Meaningful**



**Locally Faithful**



Ribeiro, M , et. al 2016

Global

Complex Non-linear

Local

Simple Linear

**Model Agnostic**

# Friendly Data Representation



**Images Space**
**m x m x 3 dimensions**

**SuperPixels Space**
**[100 dimensions]**

*Biecek, P and Burzykowski, T 2020*

# Friendly Data Representation



| SP | ON |
|----|----|
| 1 | **1** |
| 2 | **1** |
| 3 | **1** |
| ... | **1** |
| 98 | **1** |
| 99 | **1** |
| 100 | **1** |

**vector**

*Biecek, P and Burzykowski, T 2020*

# Friendly Data Representation



| SP | ON |
|----|----|
| 1 | **1** |
| 2 | **1** |
| 3 | **0** |
| ... | **…** |
| 98 | **0** |
| 99 | **1** |
| 100 | **0** |

**vector**

*Biecek, P and Burzykowski, T 2020*

# Locally Faithful



Global        Local

Complex Non-linear        Simple Linear

Ribeiro, M , et. al 2016

# Model Agnostic

EXPLANATIONS

TRUST?

https://github.com/ed-ortizm/**AID-explain-me-why** **?**

# I guess there is no option…