

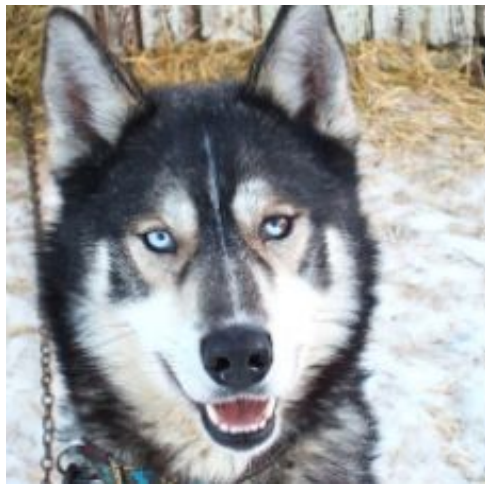
# Why Should I Trust You?

Local Interpretable Model-agnostic Explanations  
**LIME**

“... machine learning models remain mostly **black boxes**.”

Ribeiro, M , et. al 2016

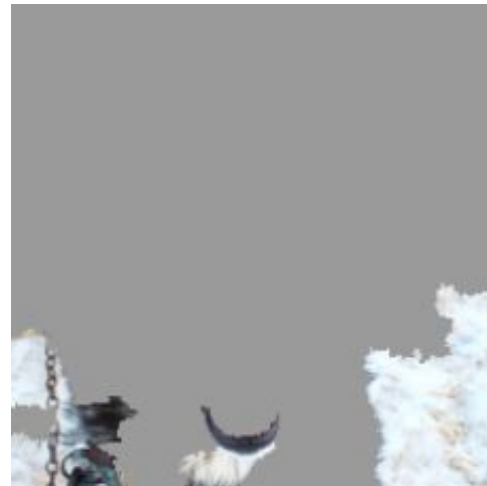
**HUSKY** classified as **WOLF**



**LIME's  
Explanation**

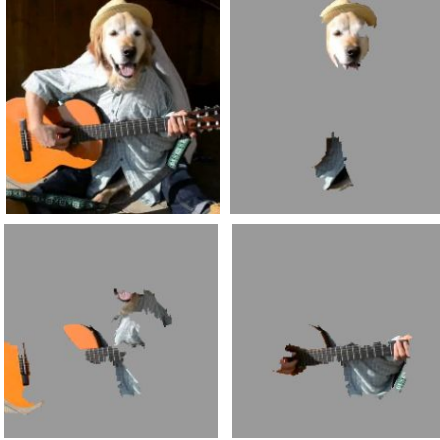


**THE WHY**

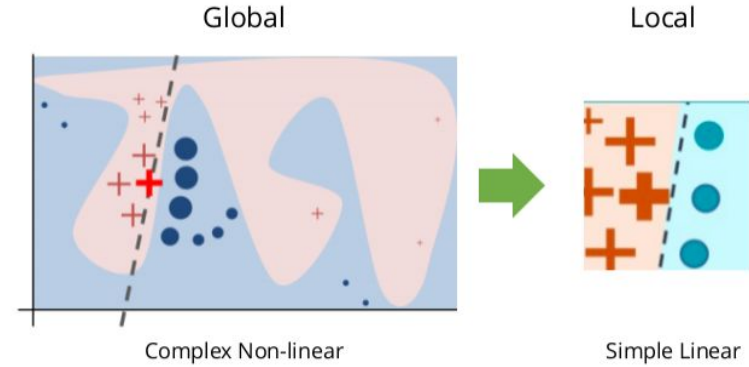


# Explanations' Desiderata

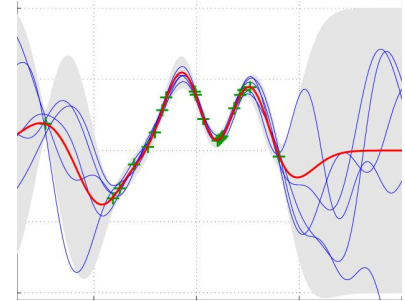
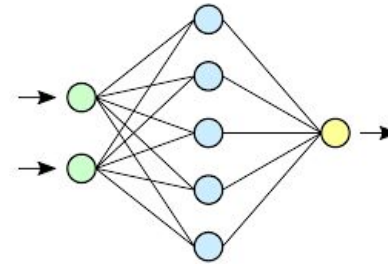
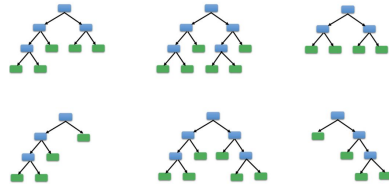
**Meaningful**



**Locally Faithful**



**Model Agnostic**



# Friendly Data Representation



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



**SuperPixels Space**  
[100 dimensions]

# Friendly Data Representation



vector



SP	ON
1	<b>1</b>
2	<b>1</b>
3	<b>1</b>
...	<b>1</b>
98	<b>1</b>
99	<b>1</b>
100	<b>1</b>

# Friendly Data Representation

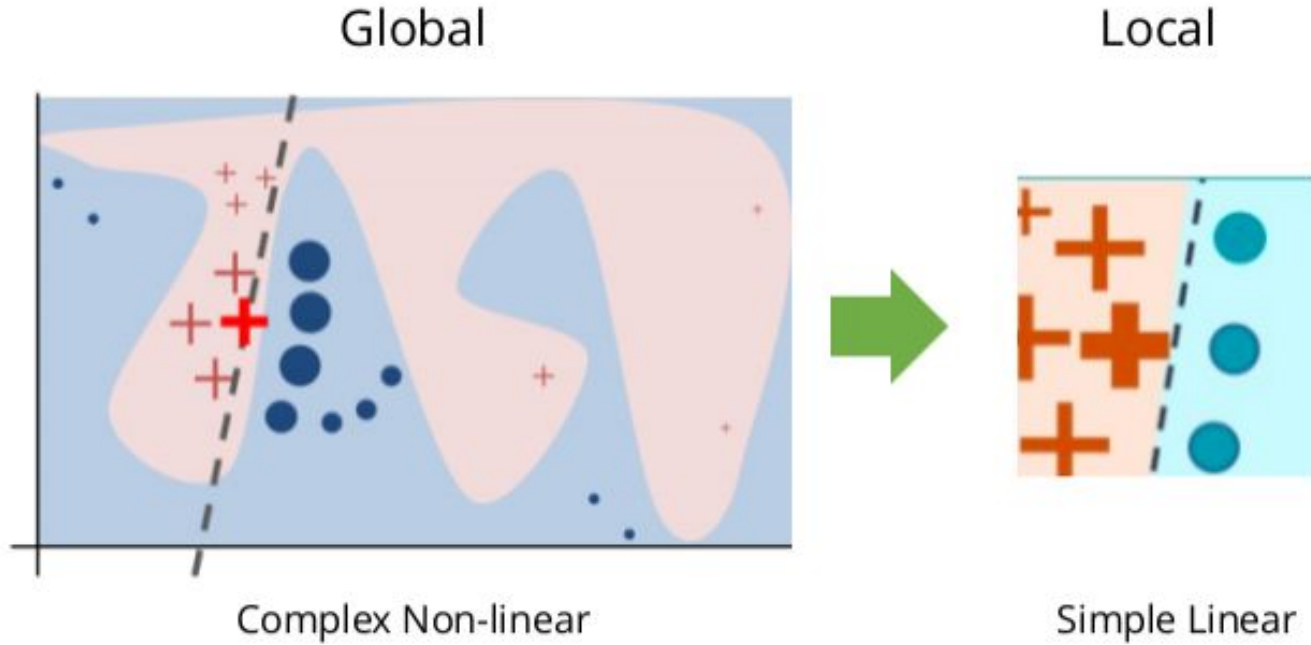


vector



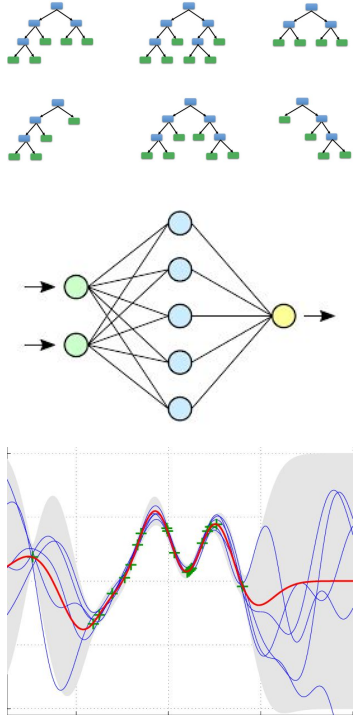
SP	ON
1	<b>1</b>
2	<b>1</b>
3	<b>0</b>
...	<b>...</b>
98	<b>0</b>
99	<b>1</b>
100	<b>0</b>

# Locally Faithful

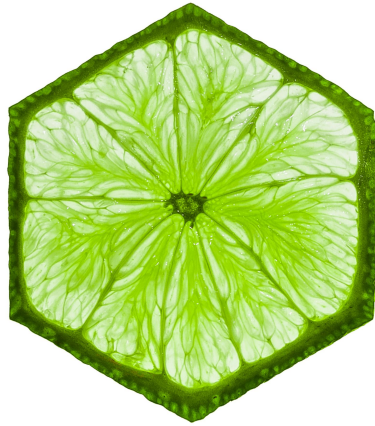




# Model Agnostic



EXPLANATIONS



TRUST?



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



**I guess there is no option...**

