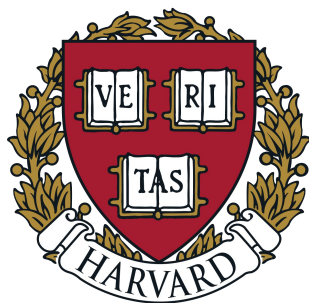
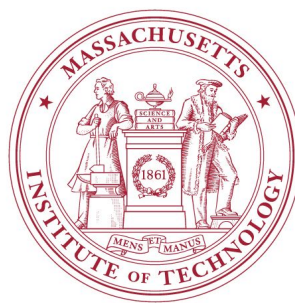


# Explaining Machine Learning Predictions: State-of-the-art, Challenges, Opportunities

Hima Lakkaraju



Julius Adebayo



Sameer Singh





**Julius Adebayo**  
MIT



**Hima Lakkaraju**  
Harvard University



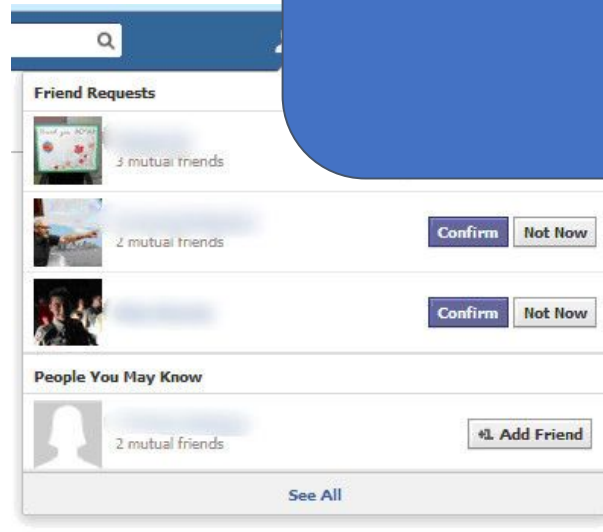
**Sameer Singh**  
UC Irvine

**Slides and Video:** [explainml-tutorial.github.io](https://explainml-tutorial.github.io)

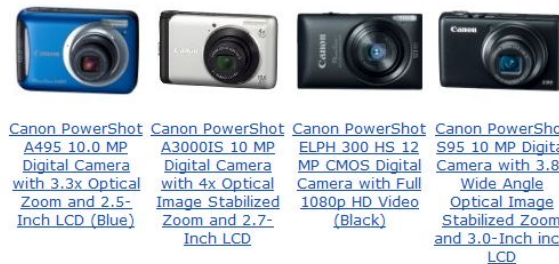
# Motivation



Machine Learning is EVERYWHERE!!



this week's bestselling models.



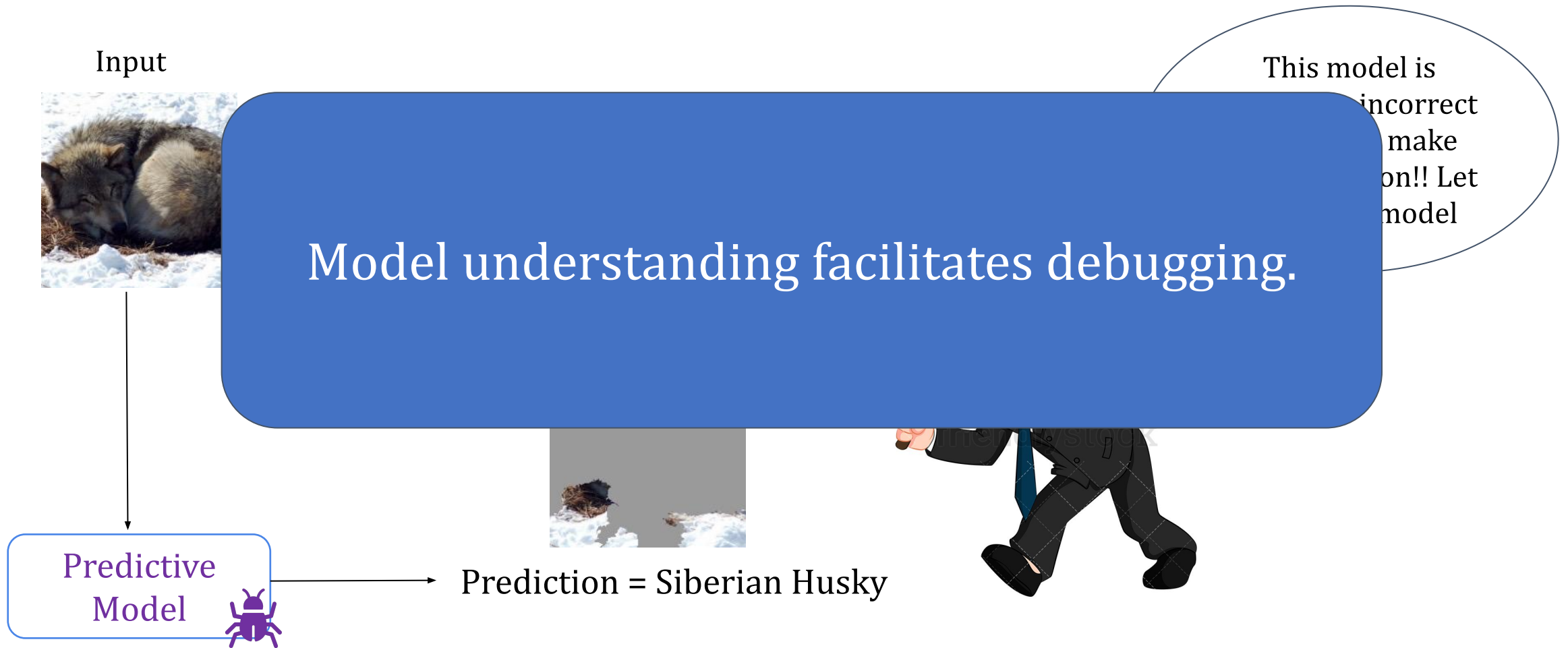
# Motivation

Model understanding is absolutely critical in several domains -- particularly those involving *high stakes decisions*!





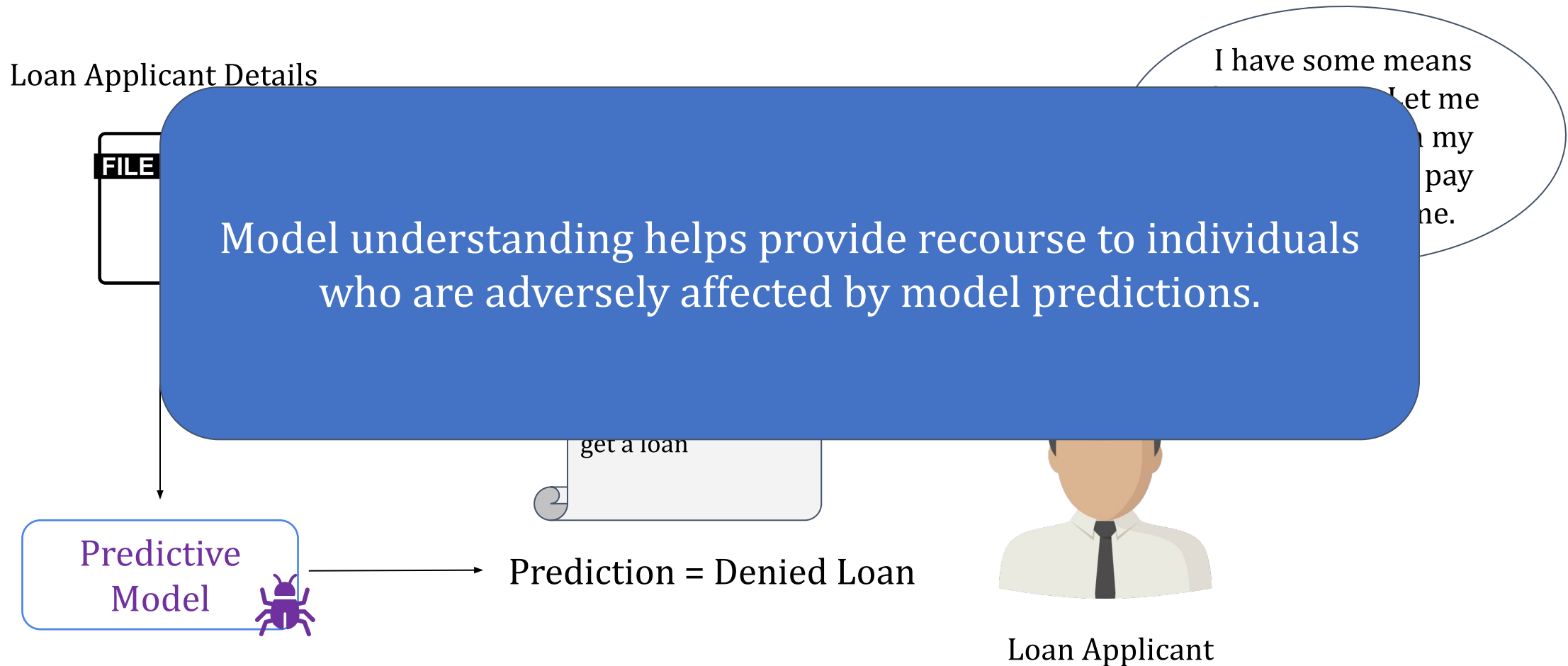
# Motivation: Why Model Understanding?



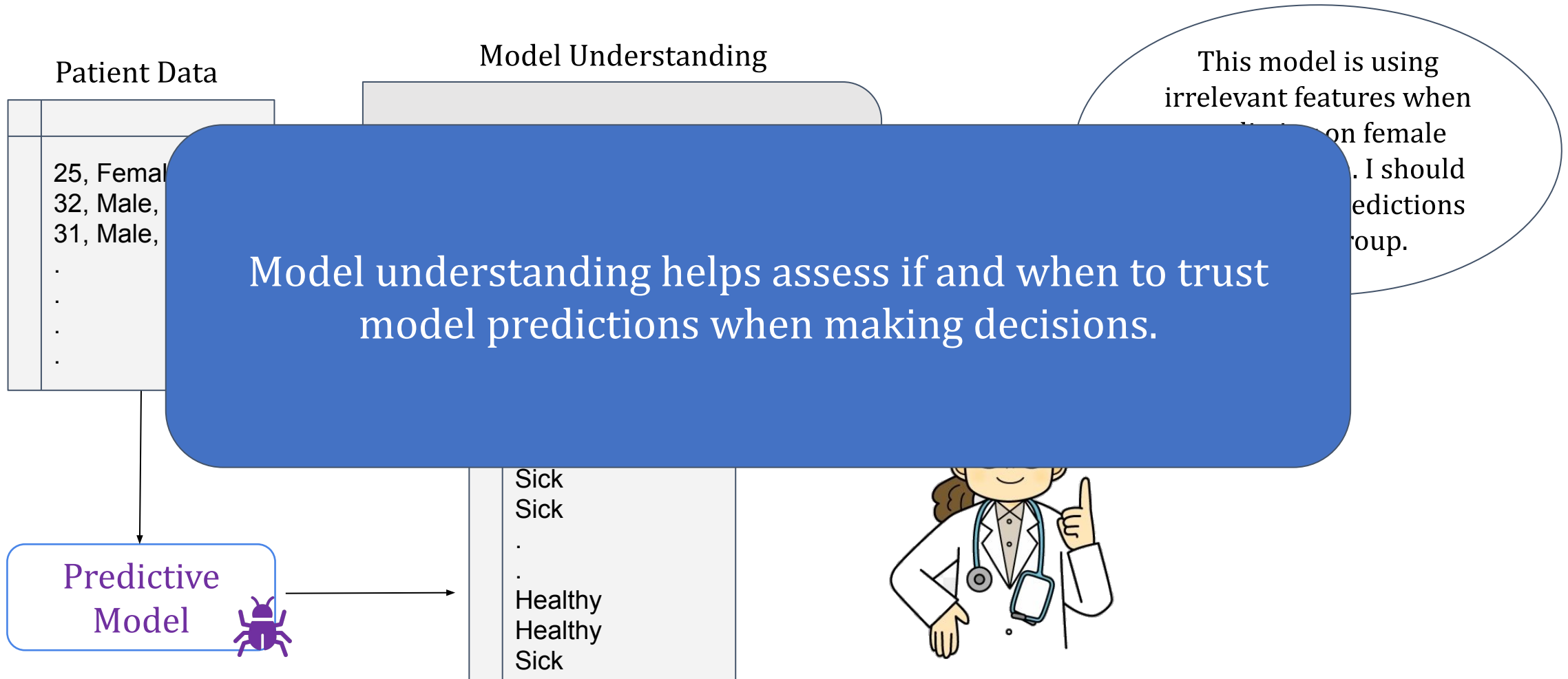
# Motivation: Why Model Understanding?



# Motivation: Why Model Understanding?

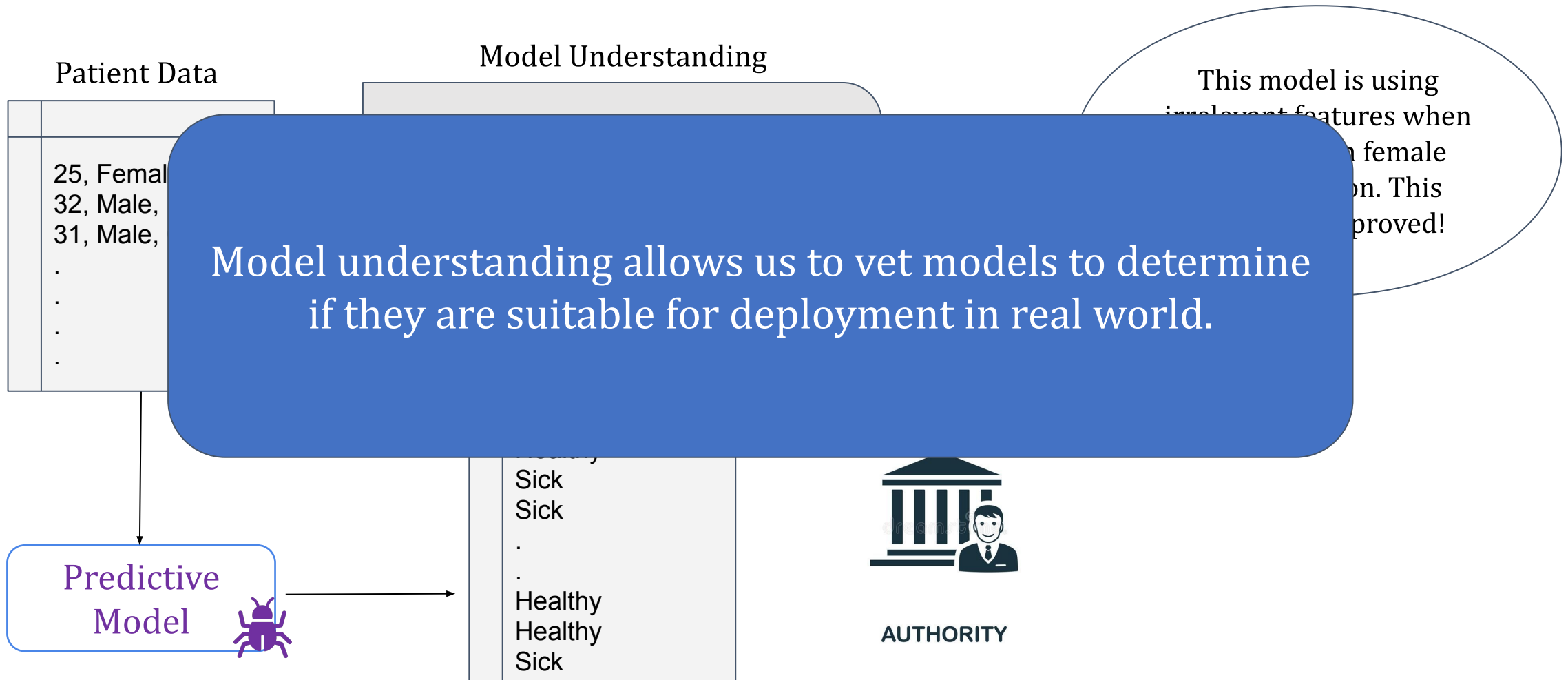


# Motivation: Why Model Understanding?





# Motivation: Why Model Understanding?



# Motivation: Why Model Understanding?

## Utility

Debugging

Bias Detection

Recourse

If and when to trust model predictions

Vet models to assess suitability for deployment

## Stakeholders

End users (e.g., loan applicants)

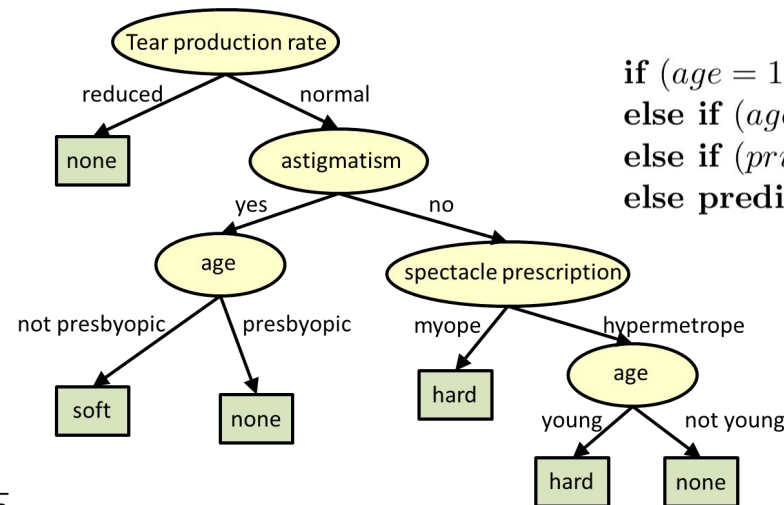
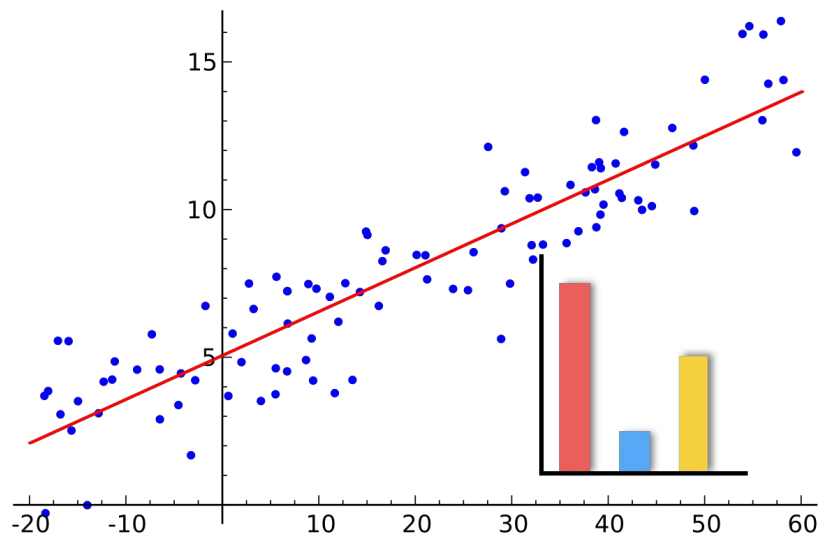
Decision makers (e.g., doctors, judges)

Regulatory agencies (e.g., FDA, European commission)

Researchers and engineers

# Achieving Model Understanding

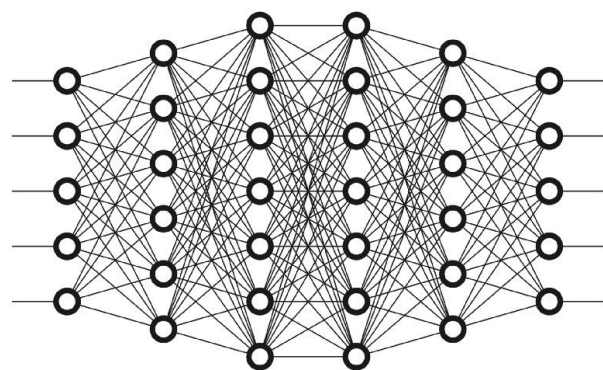
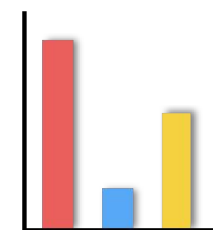
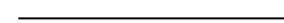
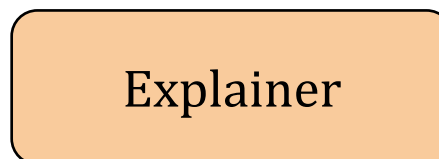
**Take 1:** Build *inherently interpretable* predictive models



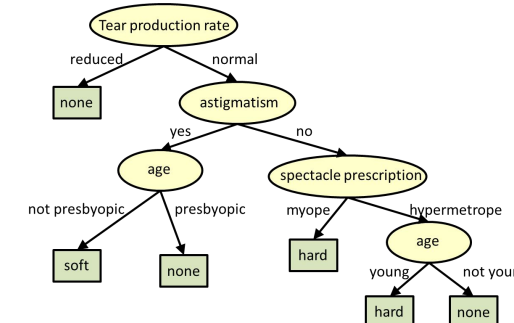
if ( $age = 18 - 20$ ) and ( $sex = male$ ) then predict *yes*  
 else if ( $age = 21 - 23$ ) and ( $priors = 2 - 3$ ) then predict *yes*  
 else if ( $priors > 3$ ) then predict *yes*  
 else predict *no*

# Achieving Model Understanding

Take 2: *Explain* pre-built models *in a post-hoc manner*



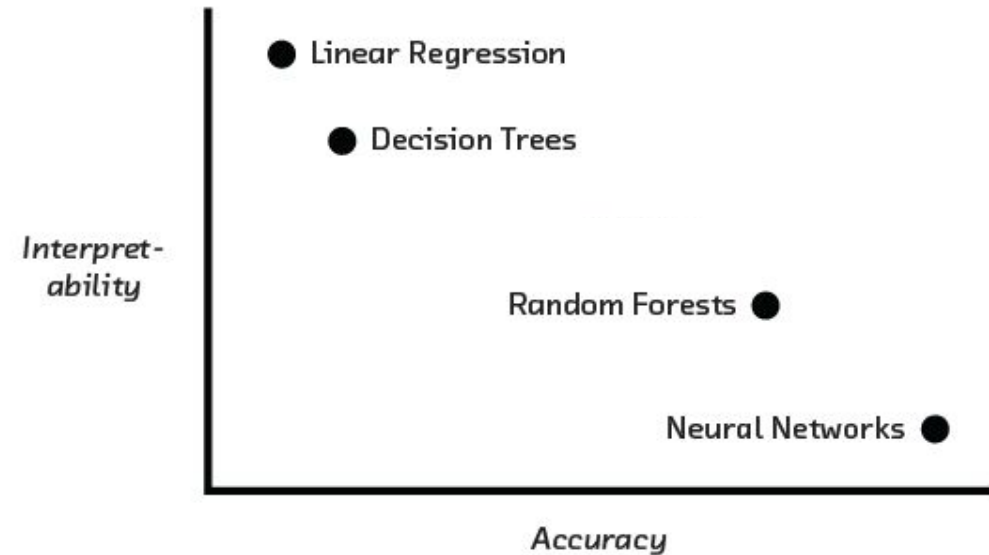
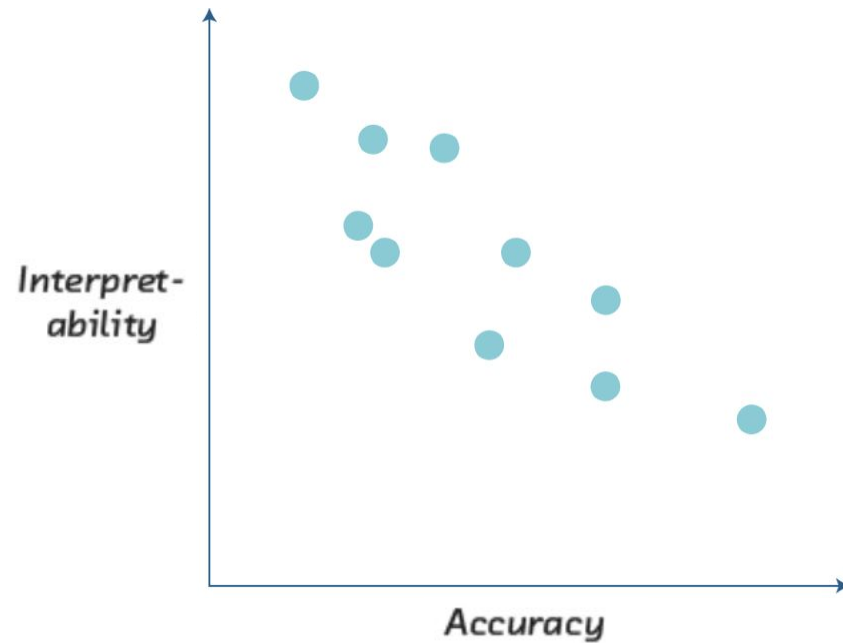
if ( $age = 18 - 20$ ) and ( $sex = male$ ) then predict *yes*  
 else if ( $age = 21 - 23$ ) and ( $priors = 2 - 3$ ) then predict *yes*  
 else if ( $priors > 3$ ) then predict *yes*  
 else predict *no*





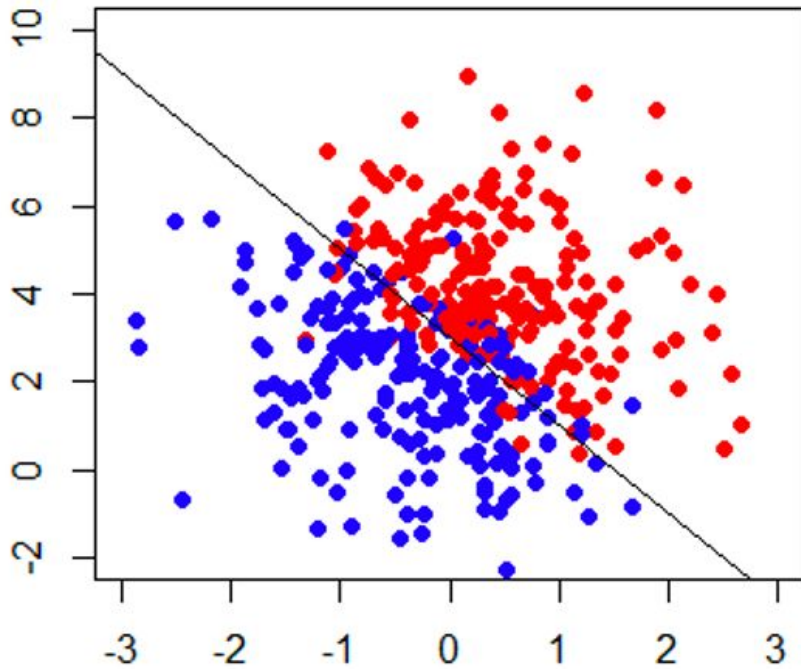
# Inherently Interpretable Models vs. Post hoc Explanations

## Example

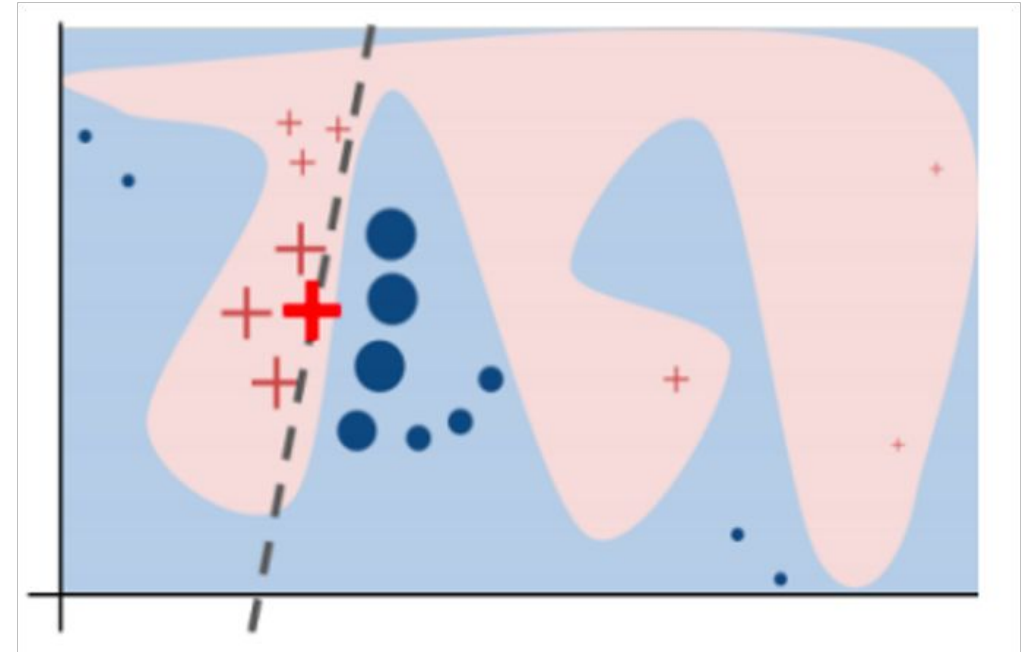


In ***certain*** settings, *accuracy-interpretability trade offs* may exist.

# Inherently Interpretable Models vs. Post hoc Explanations



can build interpretable +  
accurate models



complex models might  
achieve higher accuracy

# Inherently Interpretable Models vs. Post hoc Explanations

Sometimes, you don't have enough data to build your model from scratch.

And, all you have is a (proprietary) black box!



# Inherently Interpretable Models vs. Post hoc Explanations

If you *can build* an interpretable model which is also adequately accurate for your setting, DO IT!

Otherwise, *post hoc explanations* come to the rescue!

*This tutorial will focus on post hoc explanations!*

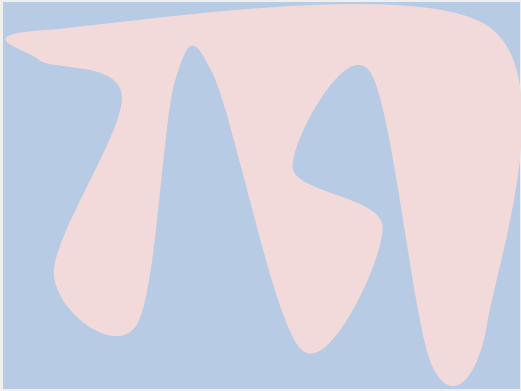


# What is an Explanation?

# What is an Explanation?

**Definition:** Interpretable description of the model behavior

Classifier



Faithful

Explanation

Understandable

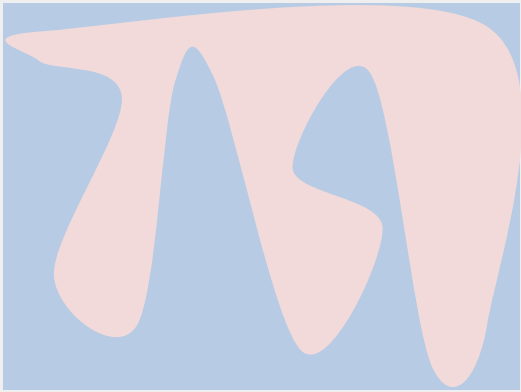
User



# What is an Explanation?

**Definition:** Interpretable description of the model behavior

Classifier



Send all the model parameters  $\theta$ ?

Send many example predictions?

Summarize with a program/rule/tree

Select most important features/points

Describe how to *flip* the model prediction

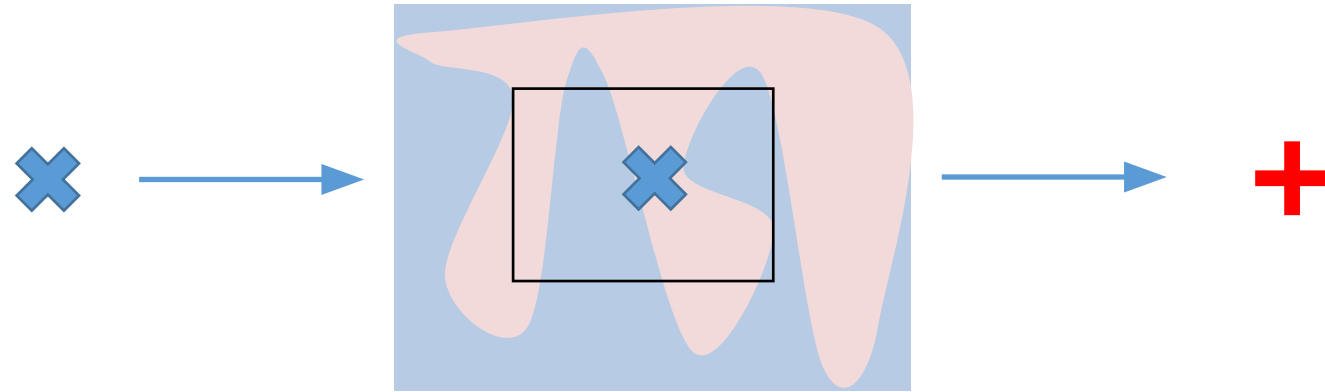
...

User



# Local versus Global Explanations

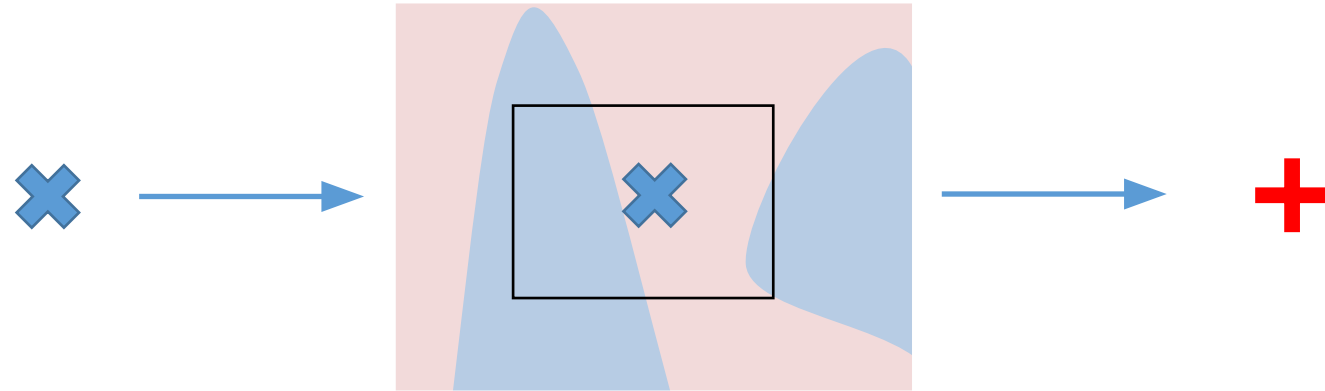
Global explanation may be too complicated





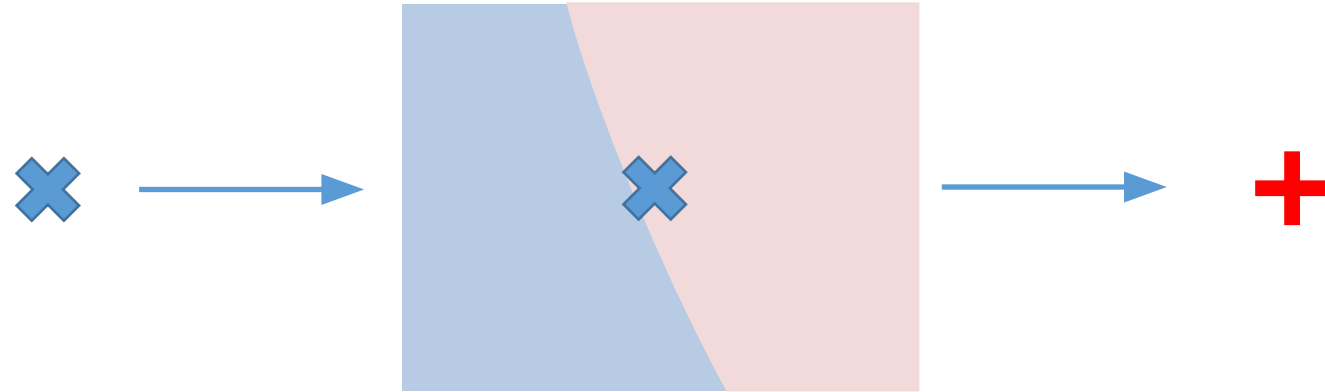
# Local versus Global Explanations

Global explanation may be too complicated



# Local versus Global Explanations

Global explanation may be too complicated

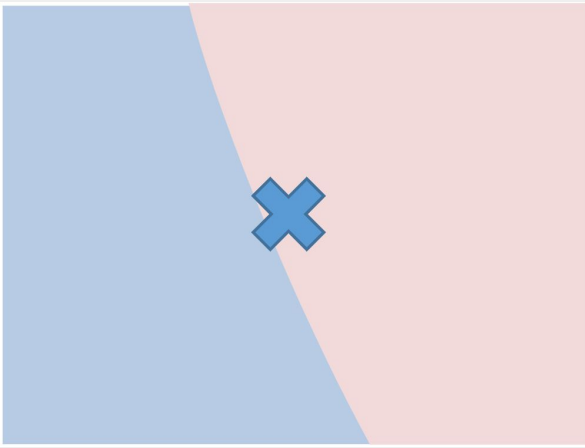


**Definition:** Interpretable description of the model behavior *in a target neighborhood*.

# Local Explanations

**Definition:** Interpretable description of the model behavior *in a target neighborhood*.

Classifier



Send many example predictions?

Summarize with a program/rule/tree

Select most important features/points

Describe how to *flip* the model prediction

...

User



# Local Explanations vs. Global Explanations

Explain individual predictions

Help unearth biases in the *local neighborhood* of a given instance

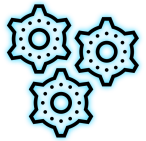
Help vet if individual predictions are being made for the right reasons

Explain complete behavior of the model

Help shed light on *big picture biases* affecting larger subgroups

Help vet if the model, at a high level, is suitable for deployment

# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations

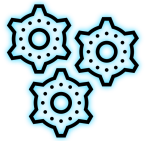


**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability

# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations

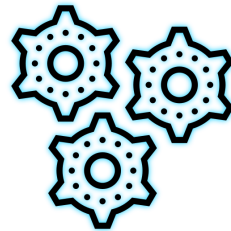


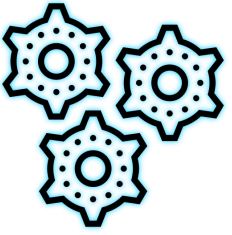
**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability

# Approaches for Post hoc Explainability





# Approaches for Post hoc Explainability

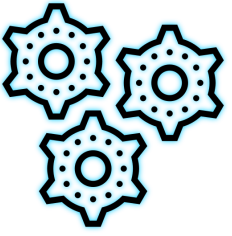
## Local Explanations

- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based





# Approaches for Post hoc Explainability

## Local Explanations

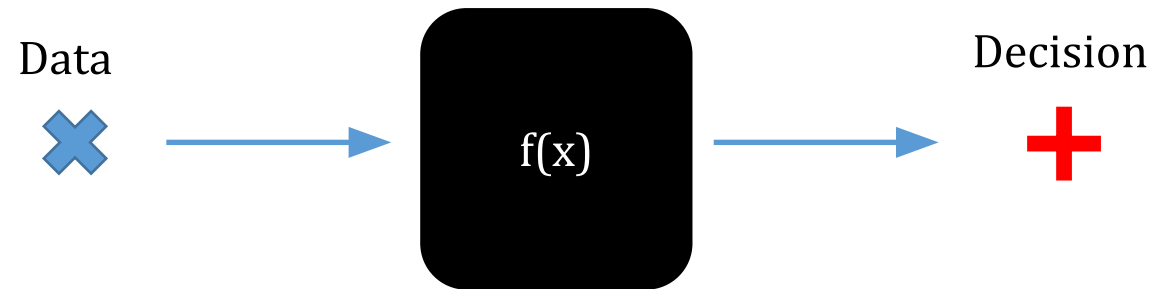
- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Being Model-Agnostic...

No access to the internal structure...



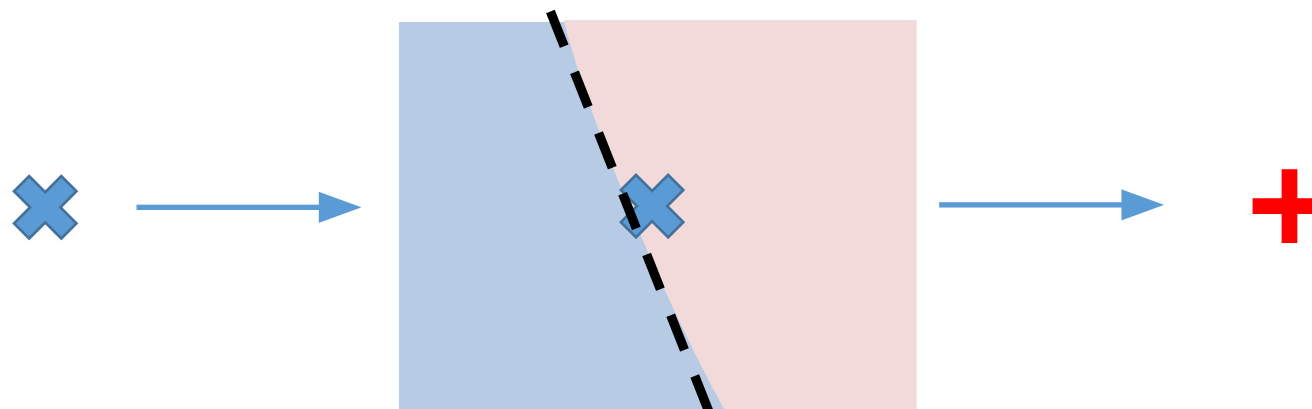
Not restricted to specific models

Practically easy: not tied to PyTorch, Tflow, etc.

Study models that you don't have access to!

# LIME: Sparse, Linear Explanations

Identify the important dimensions,  
and present their relative importance





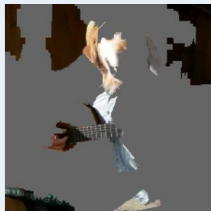



# LIME Example - Images



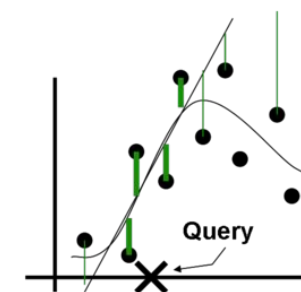
Original Image

$P(\text{labrador}) = 0.21$

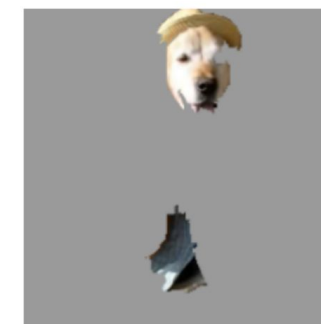


Perturbed Instances	P(Labrador)
	 0.92
	 0.001
	 0.34

Maybe to a fault?



Locally weighted  
regression



Explanation

LIME is quite customizable:

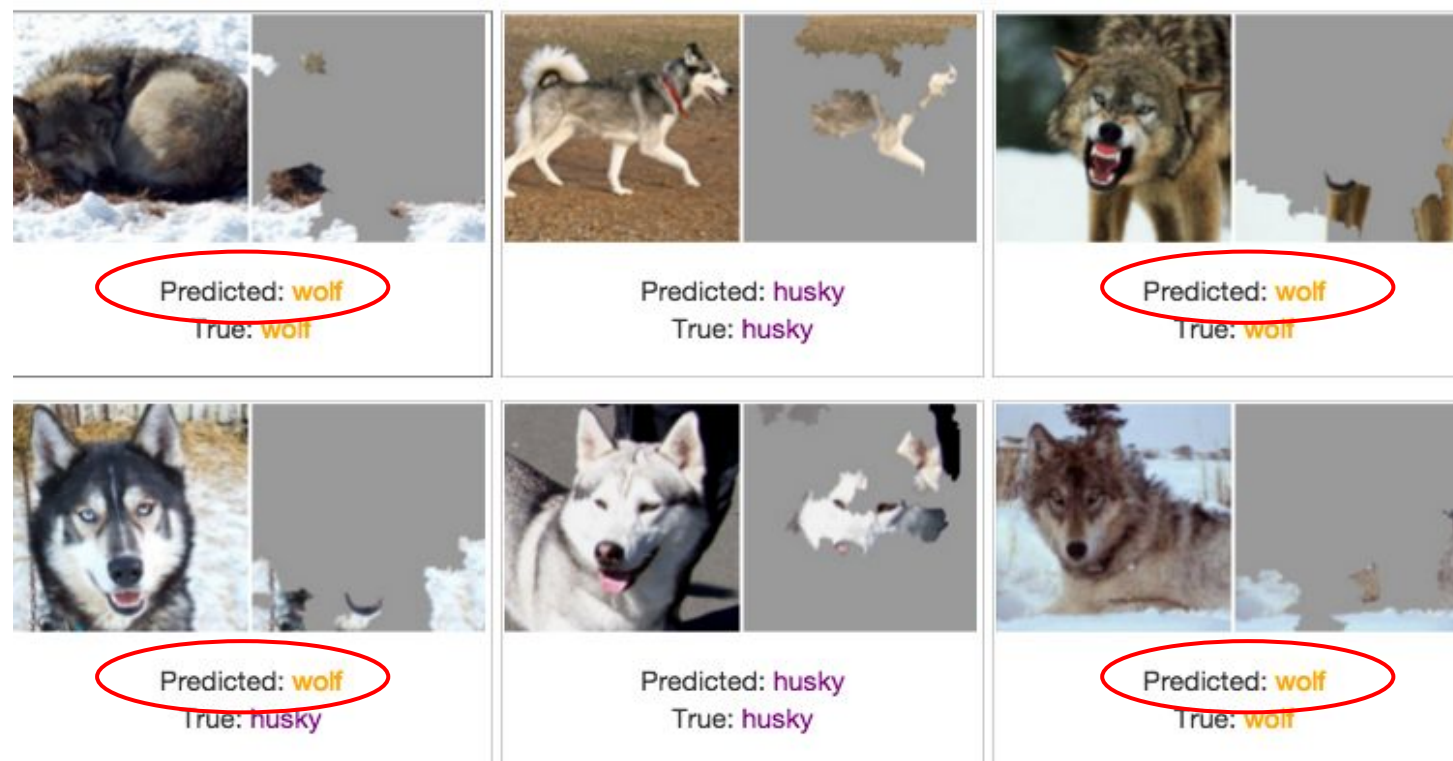
- How to perturb?
- Distance/similarity?
- How *local* you want it to be?
- How to express explanation

# Predict Wolf vs Husky

Only 1 mistake!



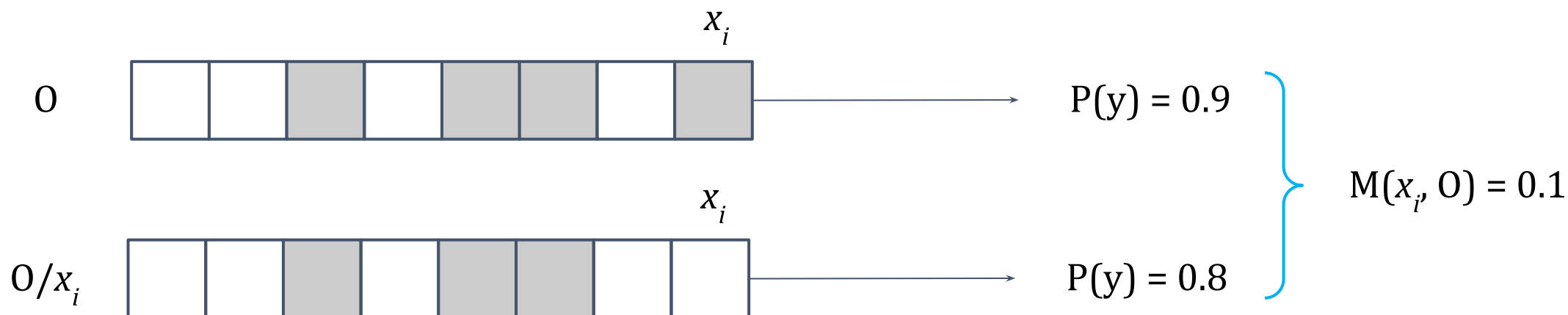
# Predict Wolf vs Husky



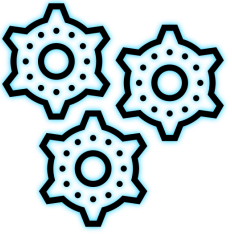
We've built a great snow detector...

# SHAP: Shapley Values as Importance

**Marginal contribution** of each feature towards the prediction, averaged over all possible permutations.



**Fairly attributes** the prediction to all the features.



# Approaches for Post hoc Explainability

## Local Explanations

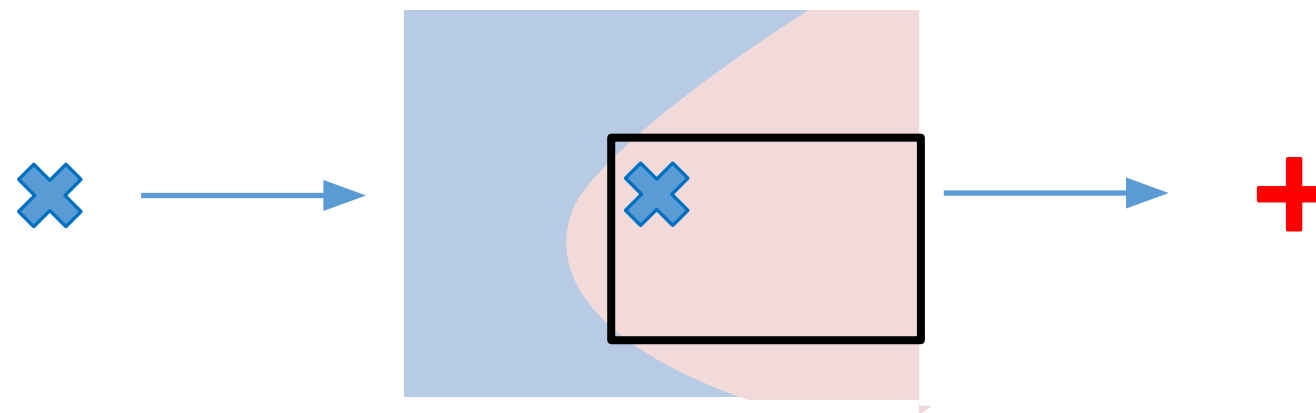
- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based



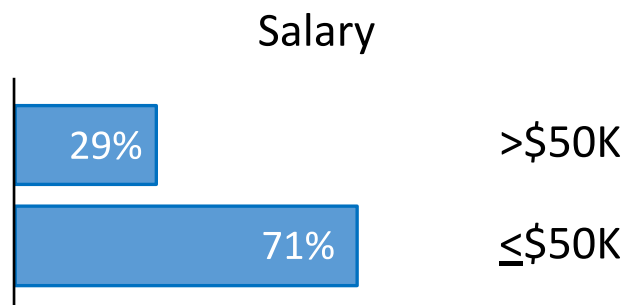
# Anchors: Sufficient Conditions



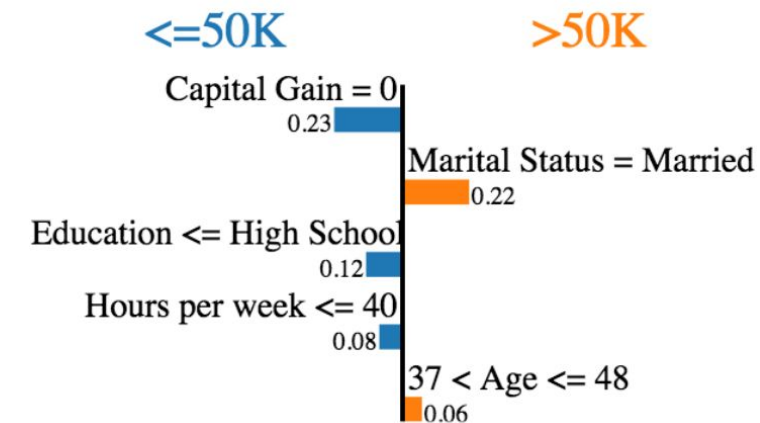
Identify the conditions under which the classifier has the same prediction

# Salary Prediction

Feature	Value
Age	37 $< \text{Age} \leq 48$
Workclass	Private
Education	$\leq$ High School
Marital Status	Married
Occupation	Craft-repair
Relationship	Husband
Race	Black
Sex	Male
Capital Gain	0
Capital Loss	0
Hours per week	$\leq 40$
Country	United States

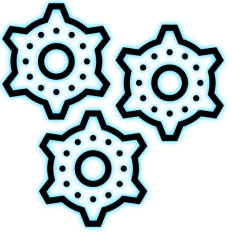


LIME



Anchors

**IF Education  $\leq$  High School  
Then Predict Salary  $\leq 50K$**



# Approaches for Post hoc Explainability

## Local Explanations

- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Saliency Map Overview

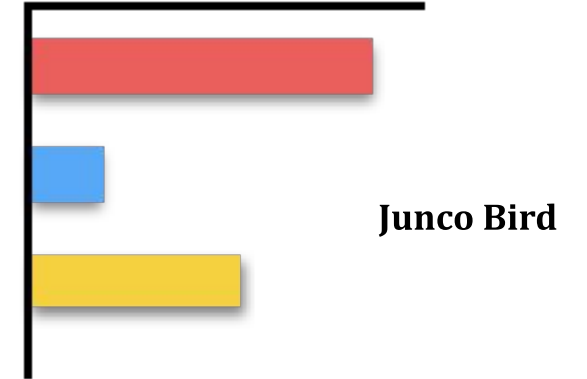
Input



Model



Predictions



# Saliency Map Overview

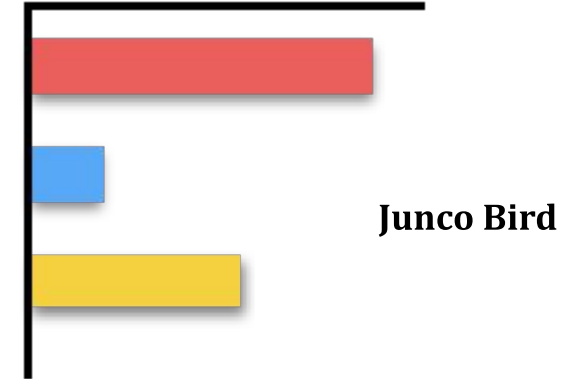
Input



Model



Predictions



What parts of the input are most relevant for the model's prediction: **'Junco Bird'**?

# Saliency Map Overview

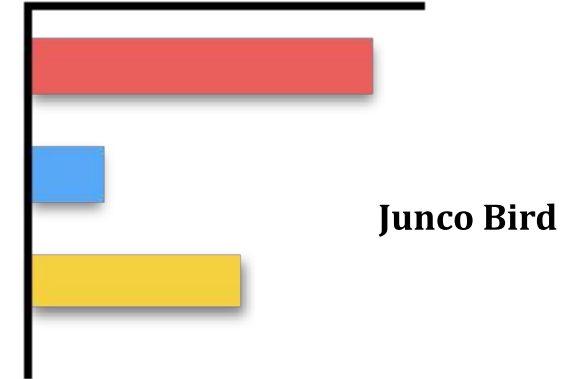
Input



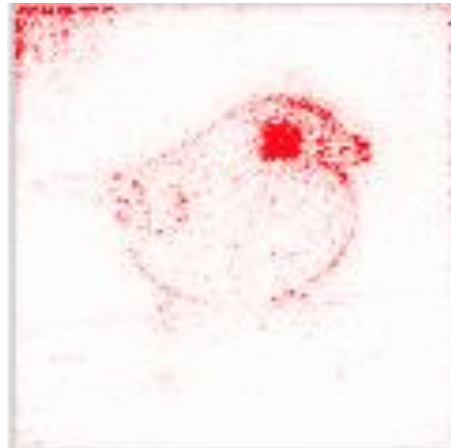
Model



Predictions

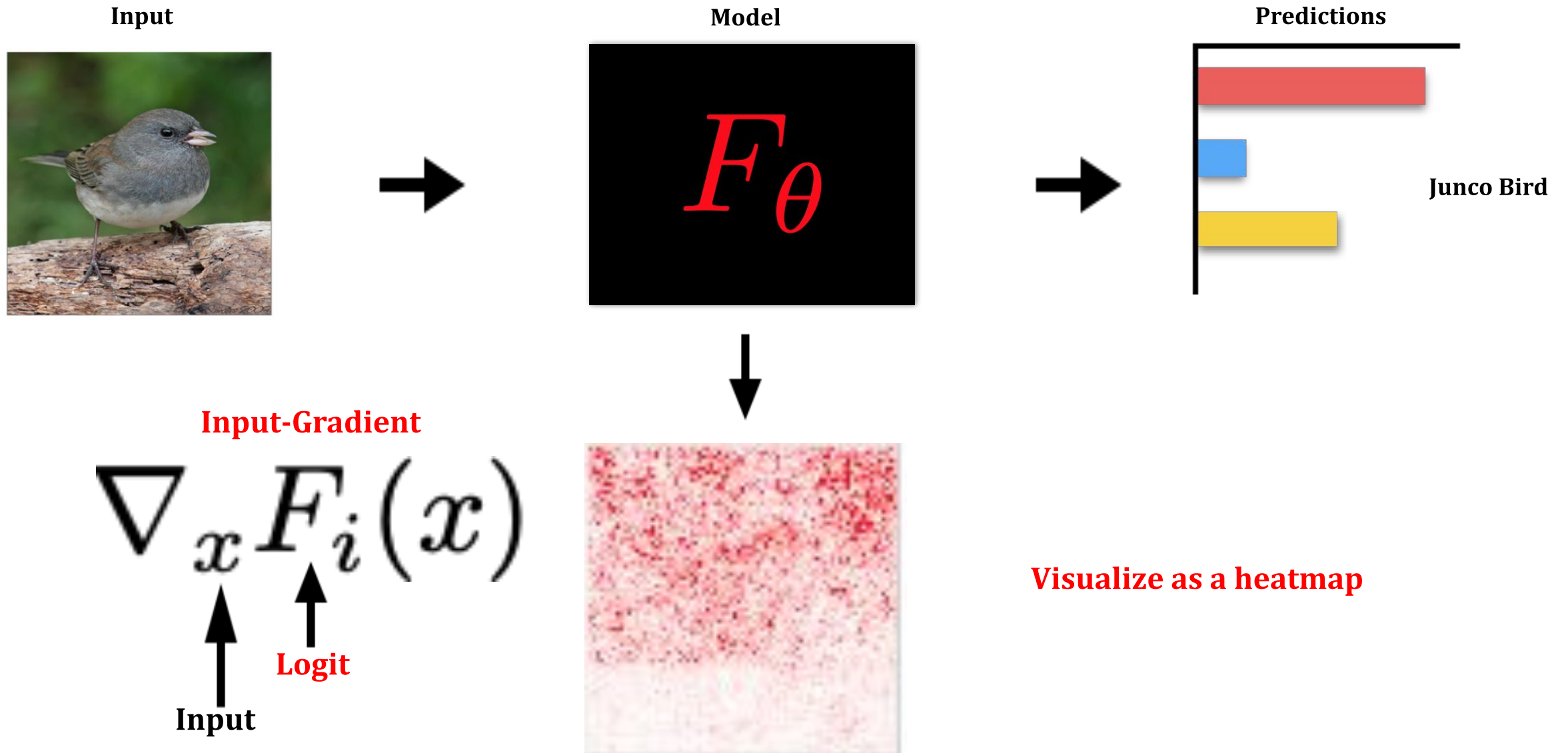


What parts of the input are most relevant for the model's prediction: **'Junco Bird'**?

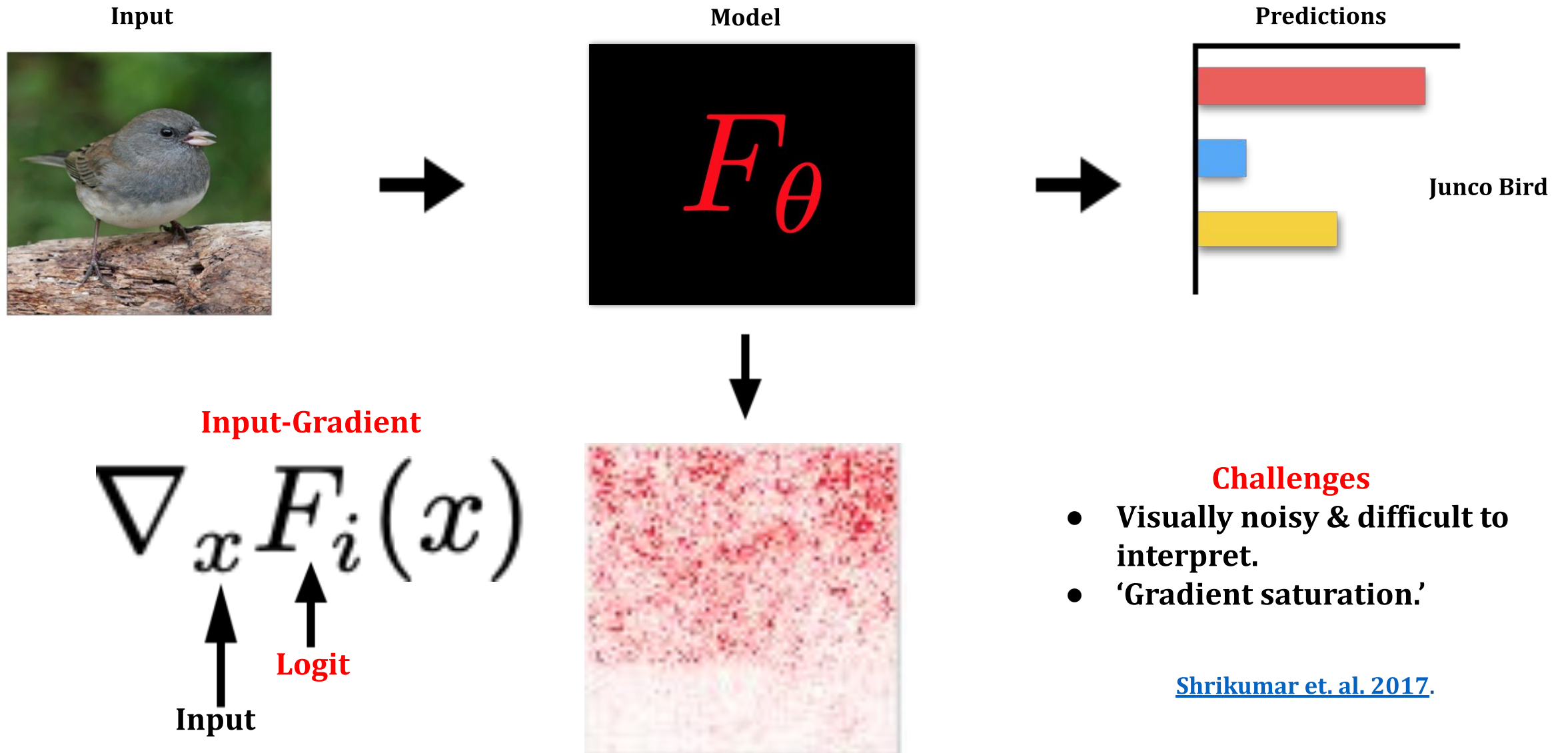


- Feature Attribution
- 'Saliency Map'
- Heatmap

# Input-Gradient



# Input-Gradient





# SmoothGrad

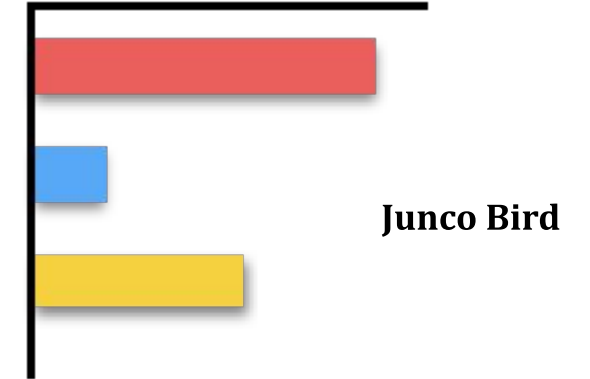
Input



Model



Predictions

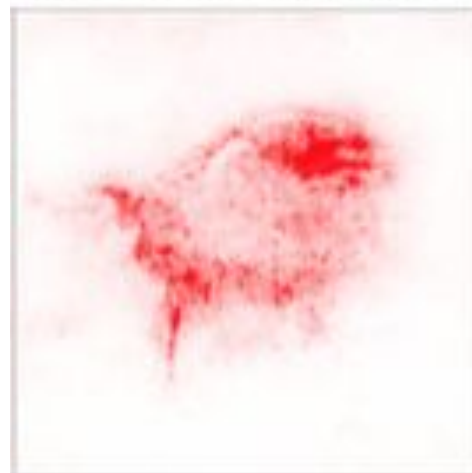


SmoothGrad

$$\frac{1}{N} \sum_i^N \nabla_{(x+\epsilon)} F_i(x + \epsilon)$$



Gaussian noise



Average Input-gradient of  
'noisy' inputs.

# Integrated Gradients

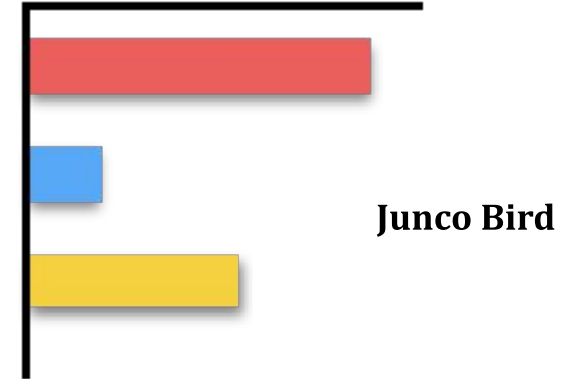
Input




Model

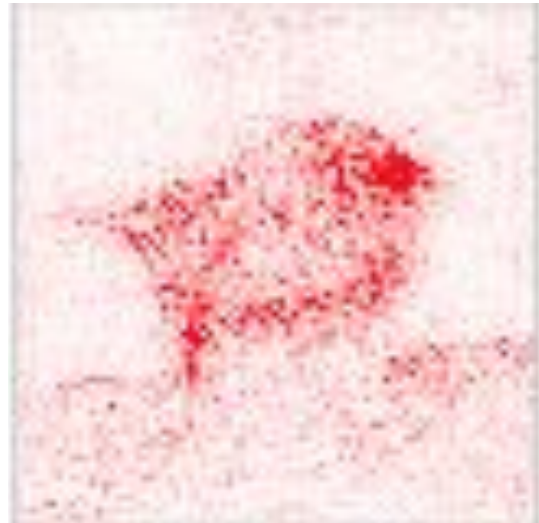


Predictions



$$(x - \tilde{x}) \times \int_{\alpha=0}^1 \frac{\partial F(\tilde{x} + \alpha \times (x - \tilde{x}))}{\partial x}$$


Baseline input



Path integral: 'sum' of interpolated gradients

# 'Modified Backprop' Approaches

Compute feature relevance by modifying the backpropagation via **positive aggregation**.

# ‘Modified Backprop’ Approaches: **Guided BackProp**

Compute feature relevance by modifying the backpropagation via **positive aggregation**.

activation:  $f_i^{l+1} = \text{relu}(f_i^l) = \max(f_i^l, 0)$

backpropagation:  $R_i^l = (f_i^l > 0) \cdot R_i^{l+1}$ , where  $R_i^{l+1} = \frac{\partial f^{out}}{\partial f_i^{l+1}}$

guided  
backpropagation:  $R_i^l = (f_i^l > 0) \cdot (R_i^{l+1} > 0) \cdot R_i^{l+1}$

# Attribution: Guided BackProp

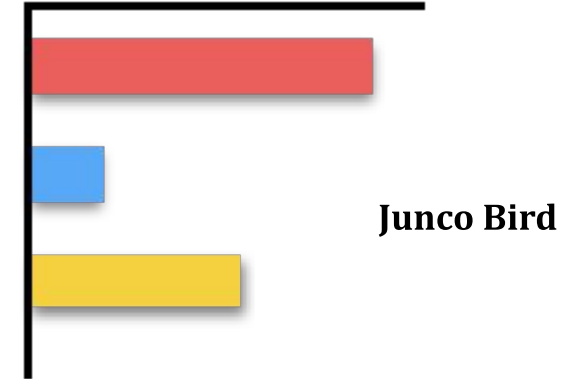
Input



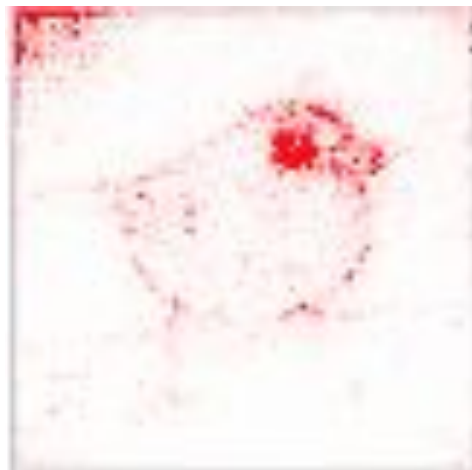
Model



Predictions



Guided BackProp



# Attribution: Guided BackProp

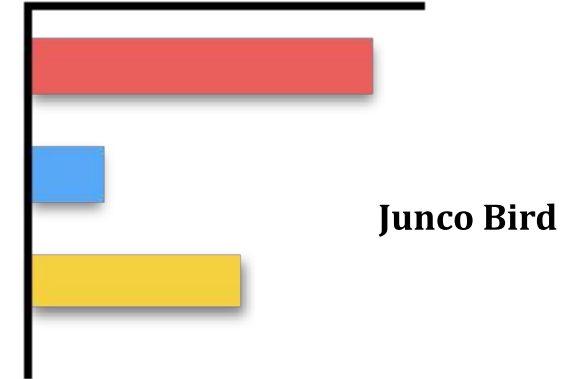
Input



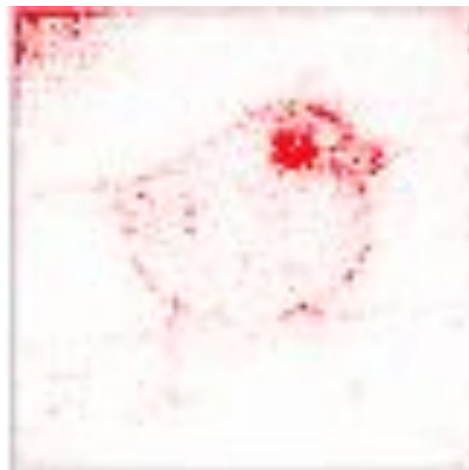
Model



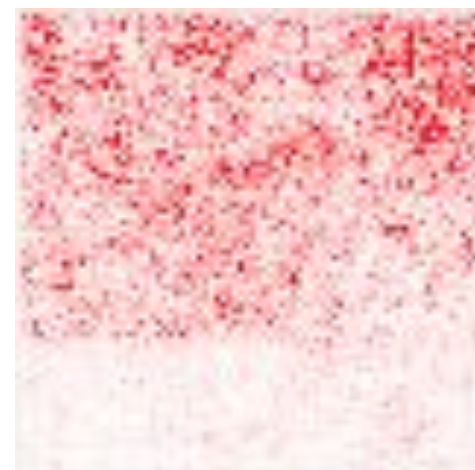
Predictions



Guided BackProp

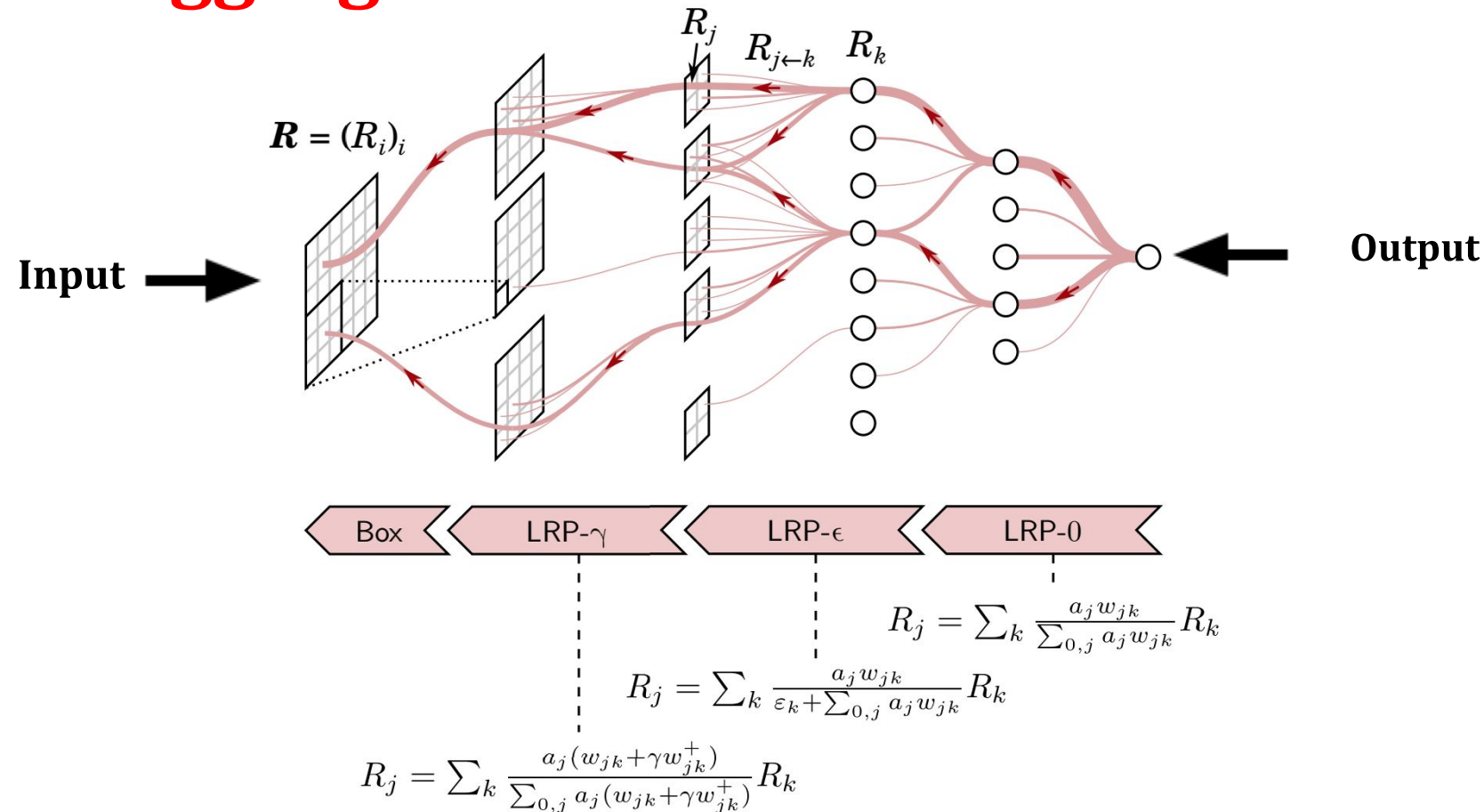


Input Gradient

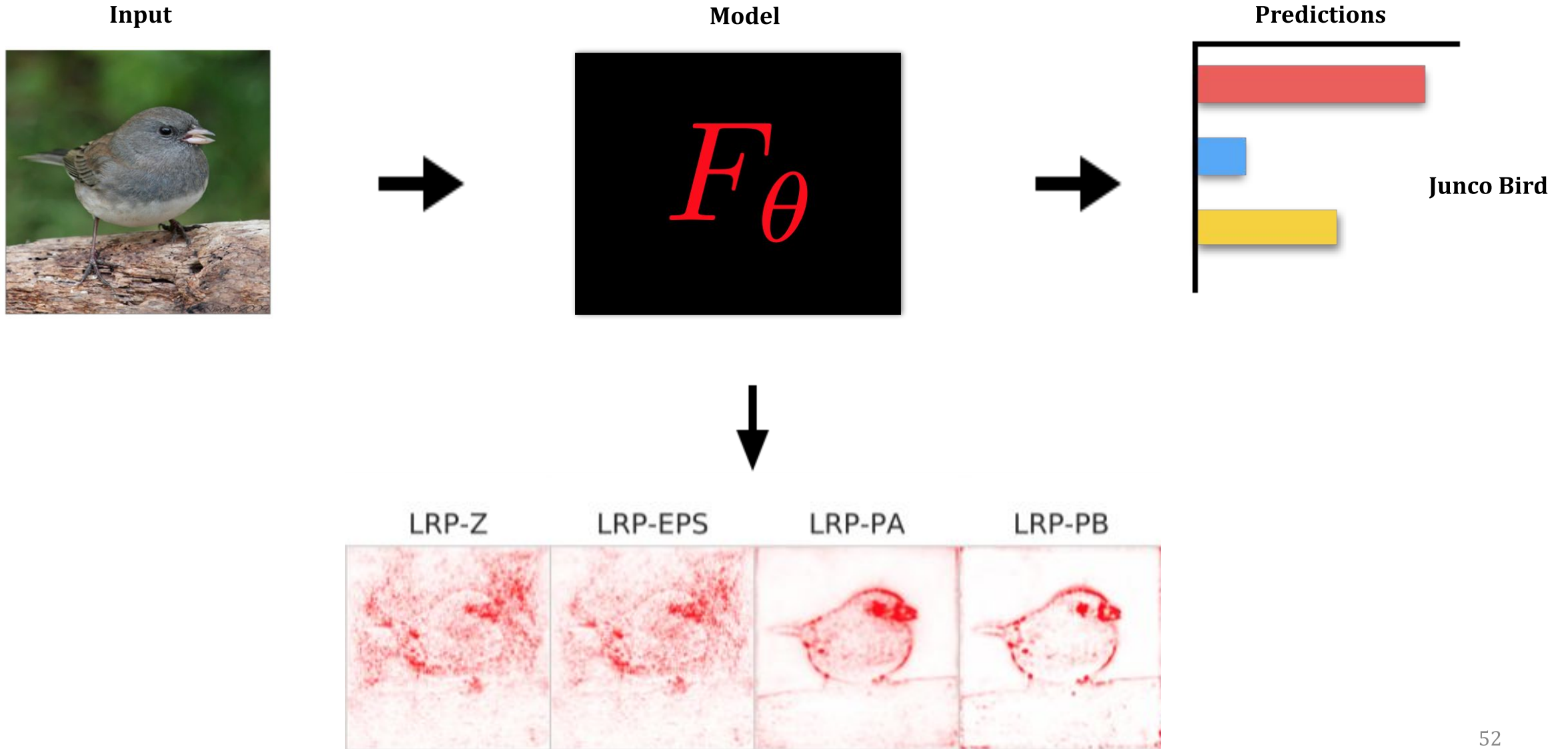


# ‘Modified Backprop’ Approaches: **LRP**

Compute feature relevance by modifying the backpropagation via **positive aggregation**.



# Layer Relevance Propagation (LRP)





# Recap

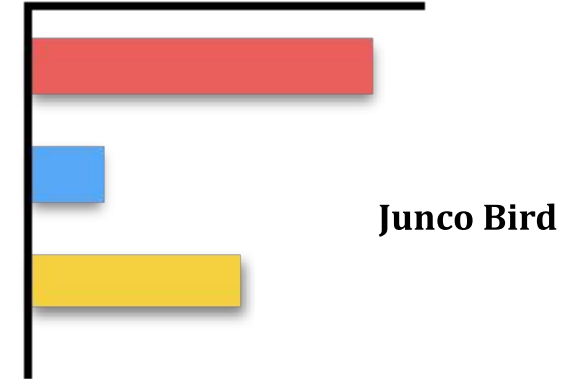
Input



Model



Predictions



LIME



SHAP



# Recap

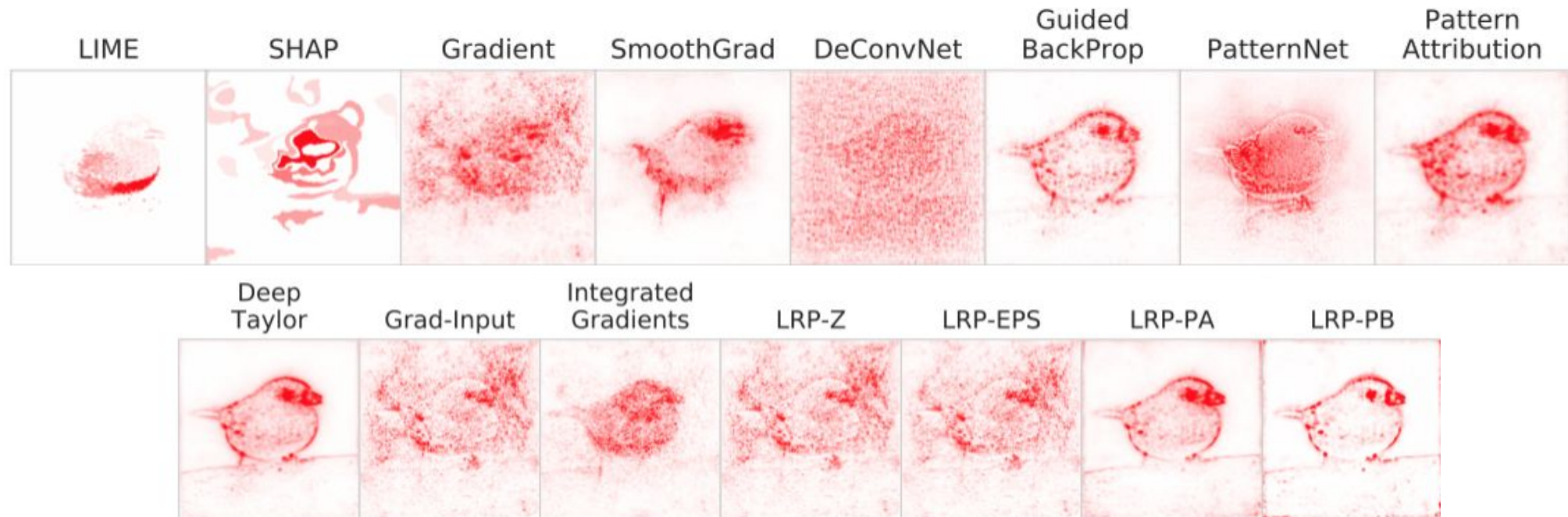
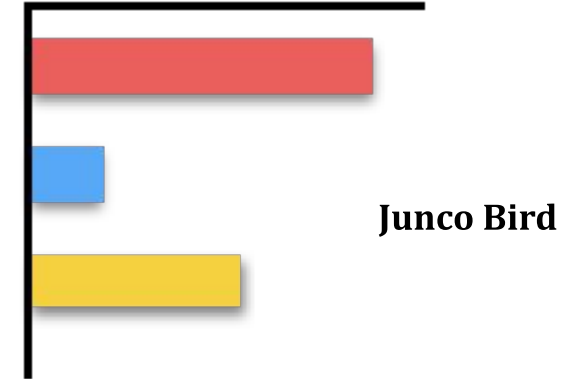
Input



Model



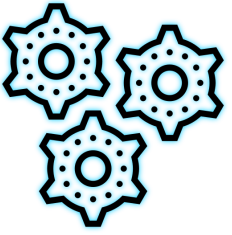
Predictions



# Additional Methods

- **Class Activation Mapping** (Zhou et. al. 2016).
- **Meaningful Perturbation** (Fong et. al. 2017).
- **RISE** (Petsuik et. al. 2018).
- **Extremal Perturbations** (Fong & Patrick 2019).
- **DeepLift** (Shrikumar et. al. 2018).
- **Expected Gradients** (Erion et. al. 2019)
- **Excitation Backprop** (Zhang et. al. 2016)
- **GradCAM** (Selvaraju et. al. 2016)
- **Guided GradCAM** (Selvaraju et. al. 2016)
- **Occlusion** (Zeiler et. al. 2014).
- **Prediction Difference Analysis** (Gu. et. al. 2019).
- **Internal Influence** (Leino et. al. 2018).

**See for additional methods:** [Samek & Montavon et. al. 2020](#)



# Approaches for Post hoc Explainability

## Local Explanations

- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Prototype Approaches

Explain a model with synthetic or natural input **‘examples’**.

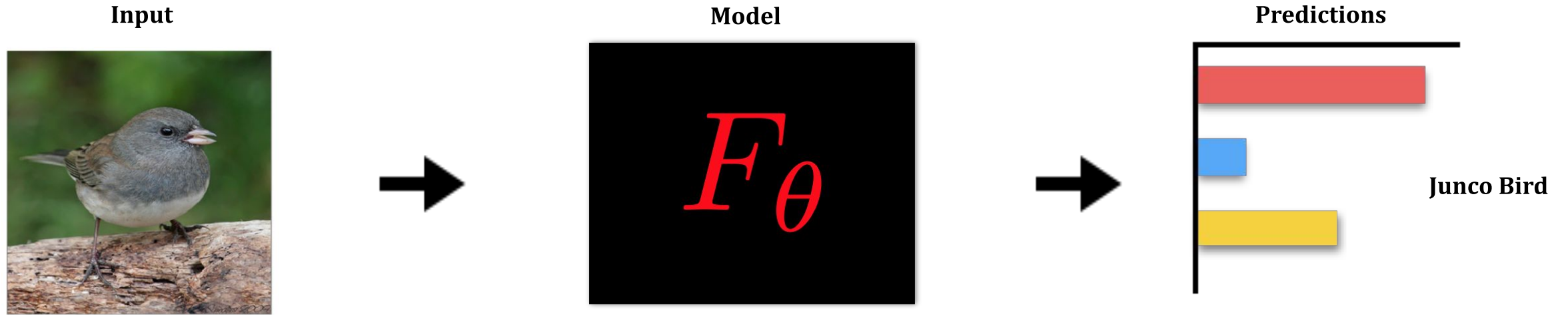
# Prototype Approaches

Explain a model with synthetic or natural input **‘examples’**.

## Insights

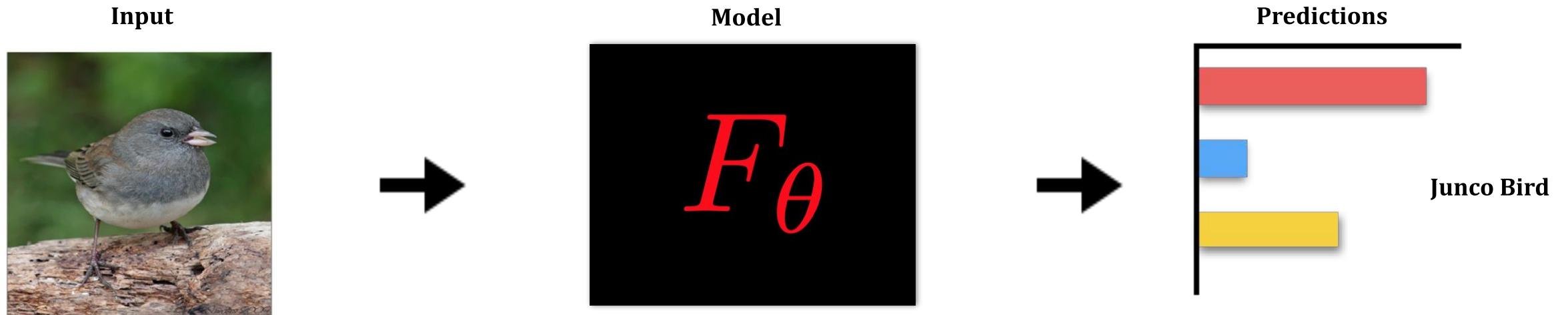
- What kind of input is the model **most likely to misclassify**?
- Which training samples are **mislabelled**?
- Which input **maximally activates** an intermediate neuron?

# Training Point Ranking via **Influence Functions**



**Which training points have the most ‘influence’ on test input’s loss?**

# Training Point Ranking via **Influence Functions**



**Which training points have the most ‘influence’ on test input’s loss?**





# Training Point Ranking via **Influence Functions**

**Influence Function:** classic tool used in robust statistics for assessing the effect of a sample on regression parameters ([Cook & Weisberg, 1980](#)).

**Influence of Training Point on Parameters**

$$\mathcal{I}_{z_j} = \left. \frac{d\hat{\theta}_{\epsilon, z_j}}{d\epsilon} \right|_{\epsilon=0} = -H_{\hat{\theta}}^{-1} \nabla_{\theta} \ell(z_j, \hat{\theta})$$

**Influence of Training Point on Test-Input's loss**

$$\mathcal{I}_{z_j, z_{\text{test}}, \text{loss}} = -\nabla_{\theta} \ell(z_{\text{test}}, \hat{\theta})^{\top} H_{\hat{\theta}}^{-1} \nabla_{\theta} \ell(z_j, \hat{\theta})$$

# Challenges and Other Approaches

## Influence function Challenges:

1. **scalability:** computing hessian-vector products can be tedious in practice.
2. **non-convexity:** possibly loose approximation for deeper networks ([Basu et. al. 2020](#)).

# Challenges and Other Approaches

## Influence function Challenges:

1. **scalability:** computing hessian-vector products can be tedious in practice.
2. **non-convexity:** possibly loose approximation for deeper networks ([Basu et. al. 2020](#)).

## Alternatives:

- **Representer Points** ([Yeh et. al. 2018](#)).
- **TracIn** ([Pruthi et. al.](#) appearing at NeuRIPs 2020).

# ‘Activation Maximization’

These approaches identify examples, synthetic or natural, that **strongly activate a function (neuron) of interest.**

# 'Activation Maximization'

These approaches identify examples, synthetic or natural, that **strongly activate a function (neuron) of interest.**

## Implementation Flavors:

- Search for **natural examples within a specified set** (training or validation corpus) that strongly activate a neuron of interest;
- **Synthesize examples**, typically via gradient descent, that strongly activate a neuron of interest.

# Feature Visualization

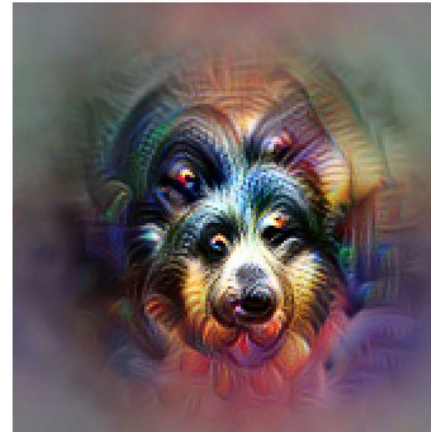
**Dataset Examples** show us what neurons respond to in practice



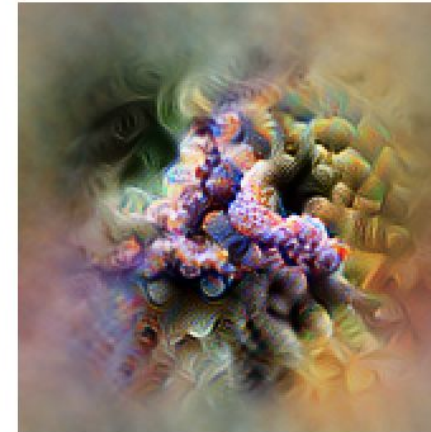
**Optimization** isolates the causes of behavior from mere correlations. A neuron may not be detecting what you initially thought.



Baseball—or stripes?  
*mixed4a, Unit 6*



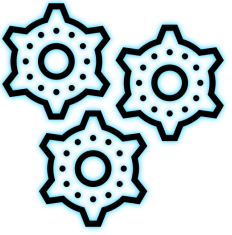
Animal faces—or snouts?  
*mixed4a, Unit 240*



Clouds—or fluffiness?  
*mixed4a, Unit 453*



Buildings—or sky?  
*mixed4a, Unit 492*



# Approaches for Post hoc Explainability

## Local Explanations

- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Counterfactual Explanations

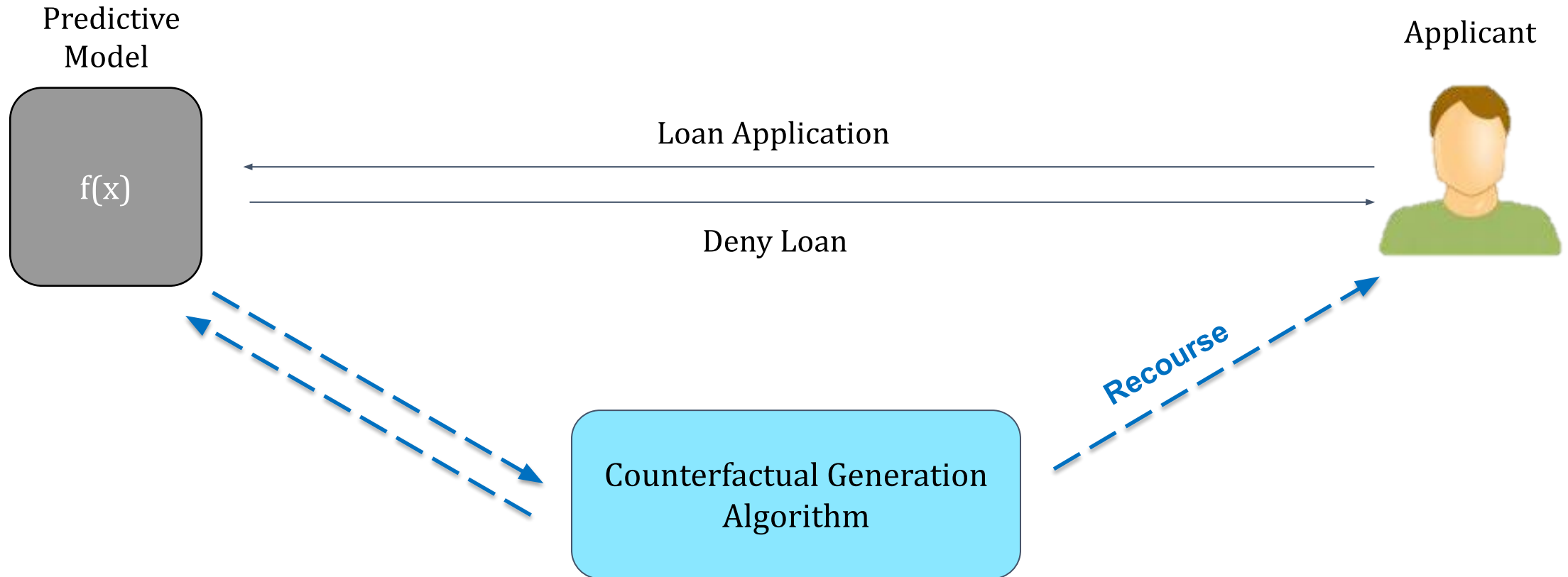
As ML models increasingly deployed to make high-stakes decisions (e.g., loan applications), it becomes important to provide **recourse** to affected individuals.

## ***Counterfactual Explanations***

*What features need to be changed and by how much to flip a model's prediction ?  
(i.e., to reverse an unfavorable outcome).*

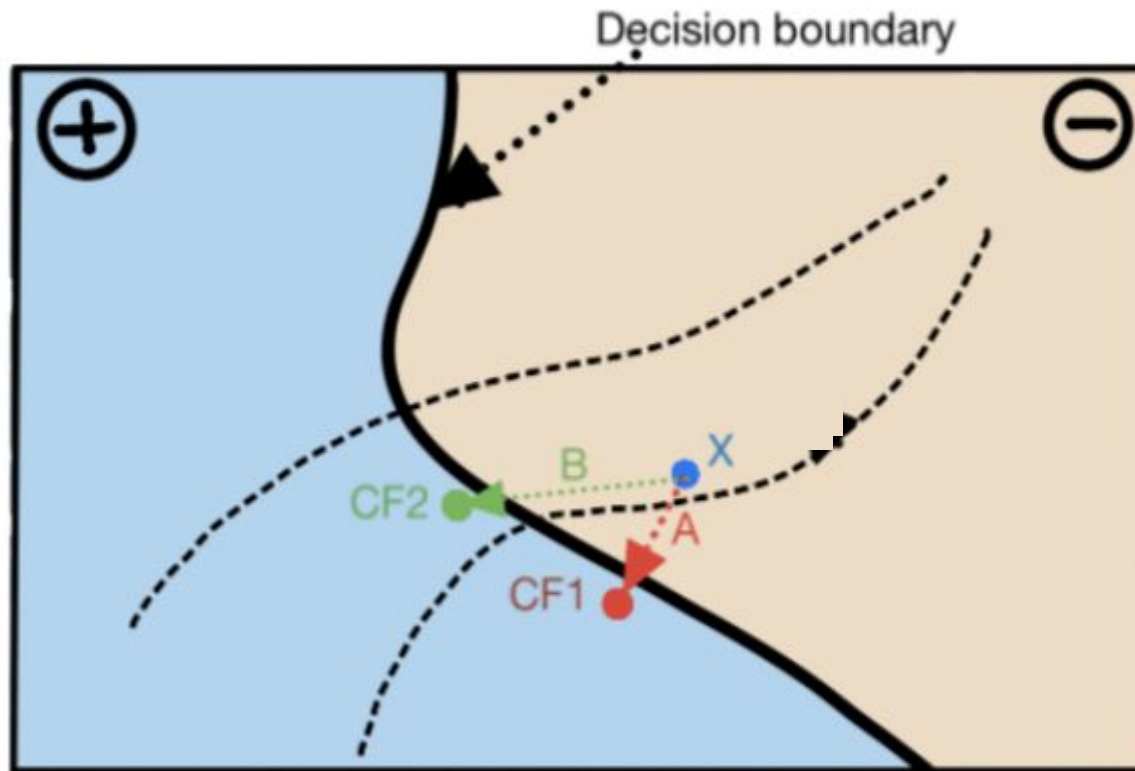


# Counterfactual Explanations



**Recourse:** Increase your salary by 50K & pay your credit card bills on time for next 3 months

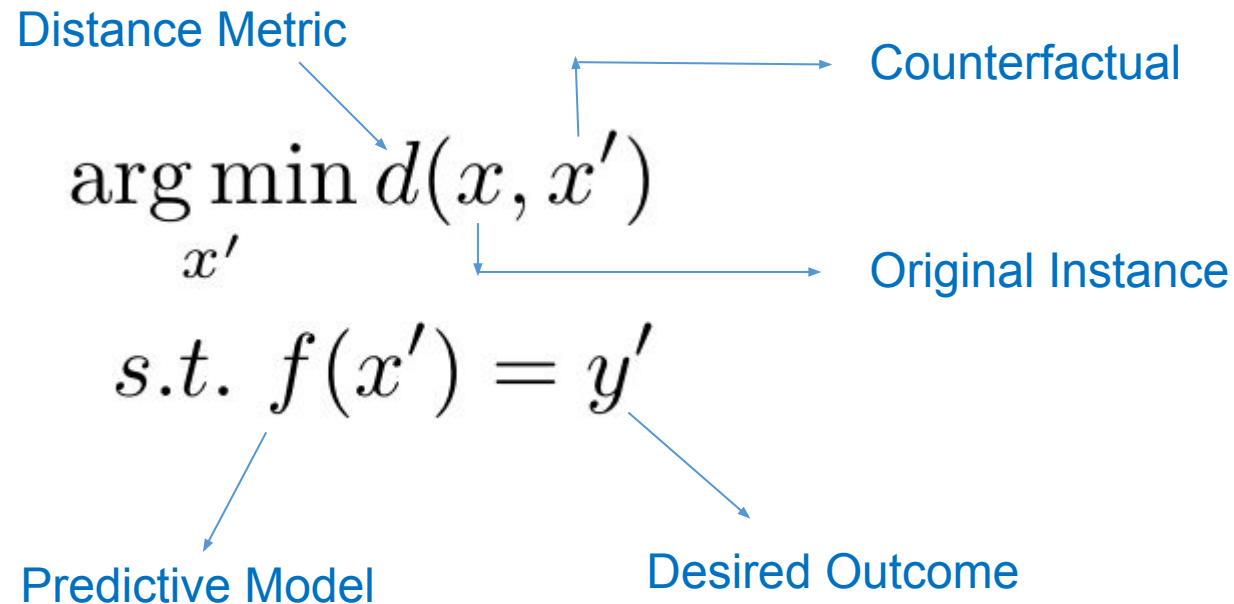
# Generating Counterfactual Explanations: Intuition



Proposed solutions differ on:

How to choose among  
candidate counterfactuals?

# Take 1: Minimum Distance Counterfactuals



Choice of distance metric dictates what kinds of counterfactuals are chosen.

Wachter et. al. use normalized Manhattan distance.

# Take 1: Minimum Distance Counterfactuals

**Person 1:** If your LSAT was 34.0, you would have an average predicted score (0).

**Person 2:** If your LSAT was 32.4, you would have an average predicted score (0).

**Person 3:** If your LSAT was 33.5, and you were 'white', you would have an average predicted score (0).

**Person 4:** If your LSAT was 35.8, and you were 'white', you would have an average predicted score (0).

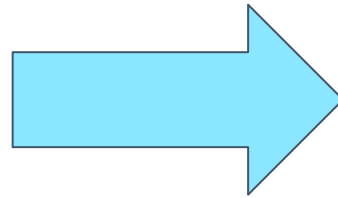
**Person 5:** If your LSAT was 34.9, you would have an average predicted score (0).



Not feasible to act upon these features!

## Take 2: Feasible and Least Cost Counterfactuals

$$\begin{aligned} \arg \min_{x'} d(x, x') \\ s.t. f(x') = y' \end{aligned}$$



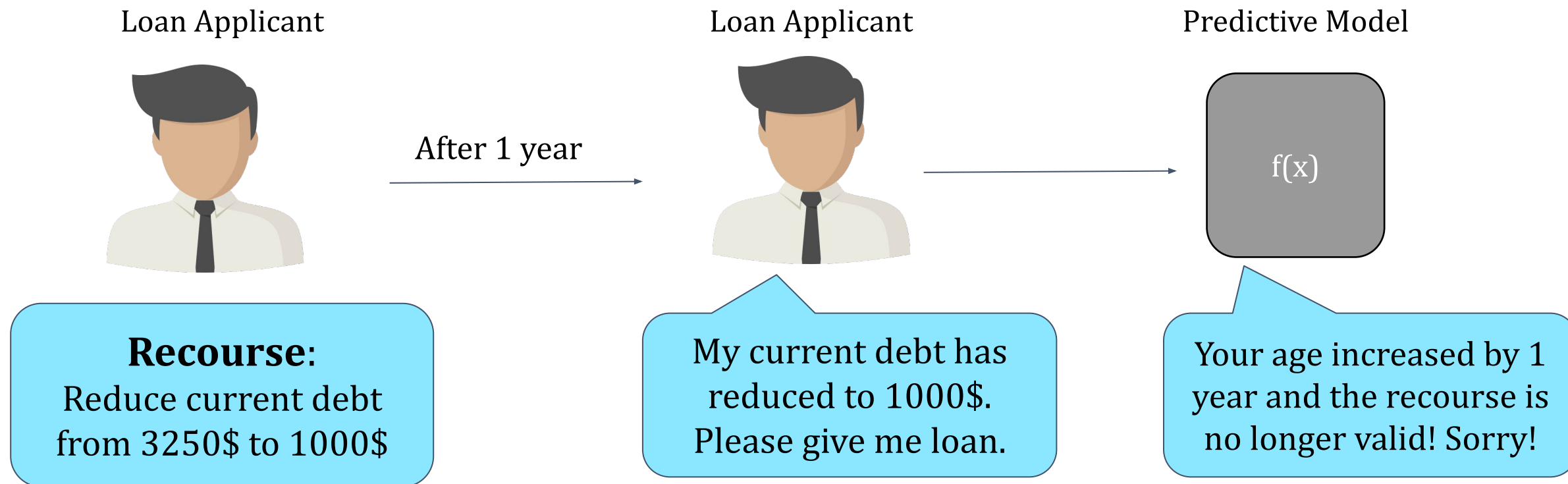
$$\begin{aligned} \arg \min_{x' \in \mathcal{A}} \text{cost}(x, x') \\ s.t. f(x') = y' \end{aligned}$$

- $\mathcal{A}$  is the set of **feasible** counterfactuals (input by end user)
  - E.g., changes to race, gender are not feasible
- **Cost** to capture how hard it is to go from  $x$  to  $x'$

# Take 2: Feasible and Least Cost Counterfactuals

FEATURES TO CHANGE	CURRENT VALUES		REQUIRED VALUES
<i>n_credit_cards</i>	5	→	3
<i>current_debt</i>	\$3,250	→	\$1,000
<i>has_savings_account</i>	FALSE	→	TRUE
<i>has_retirement_account</i>	FALSE	→	TRUE

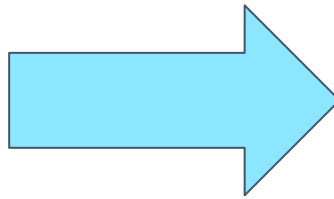
# Take 3: Causally Feasible Counterfactuals



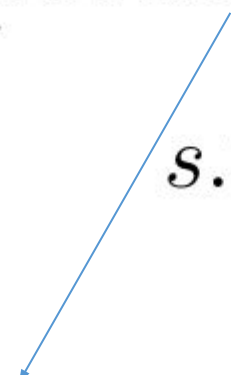
Important to account for *feature interactions* when generating counterfactuals!  
***But how?!***

# Take 3: Causally Feasible Counterfactuals

$$\begin{aligned} \arg \min_{x'} d(x, x') \\ \text{s.t. } f(x') = y' \end{aligned}$$



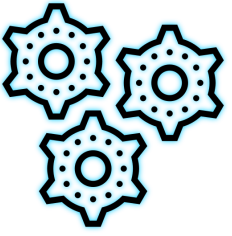
$$\begin{aligned} \arg \min_{x'} d_{\text{causal}}(x, x') \\ \text{s.t. } f(x') = y' \end{aligned}$$



Leverage Structural Causal Model (SCM) to  
define this new distance metric

Underlying causal models capture the feature interactions





# Approaches for Post hoc Explainability

## Local Explanations

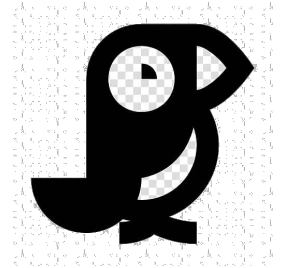
- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

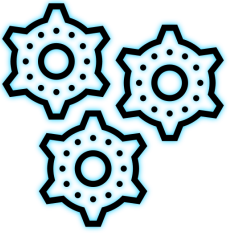
## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Global Explanations

- Explain the **complete behavior** of a given (black box) **model**
  - Provide a *bird's eye view* of model behavior
- Help **detect big picture model biases** persistent across larger subgroups of the population
  - Impractical to manually inspect local explanations of several instances to ascertain big picture biases!
- Global explanations are **complementary** to local explanations





# Approaches for Post hoc Explainability

## Local Explanations

- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Global Explanation as a Collection of Local Explanations

*How to generate a global explanation of a (black box) model?*

- Generate a local explanation for every instance in the data using one of the approaches discussed earlier
- Pick a subset of  $k$  local explanations to constitute the global explanation

# Global Explanations from Local Feature Importances: SP-LIME

LIME explains a single prediction  
local behavior for a single instance

Can't examine all explanations  
Instead pick  $k$  explanations to show to the user

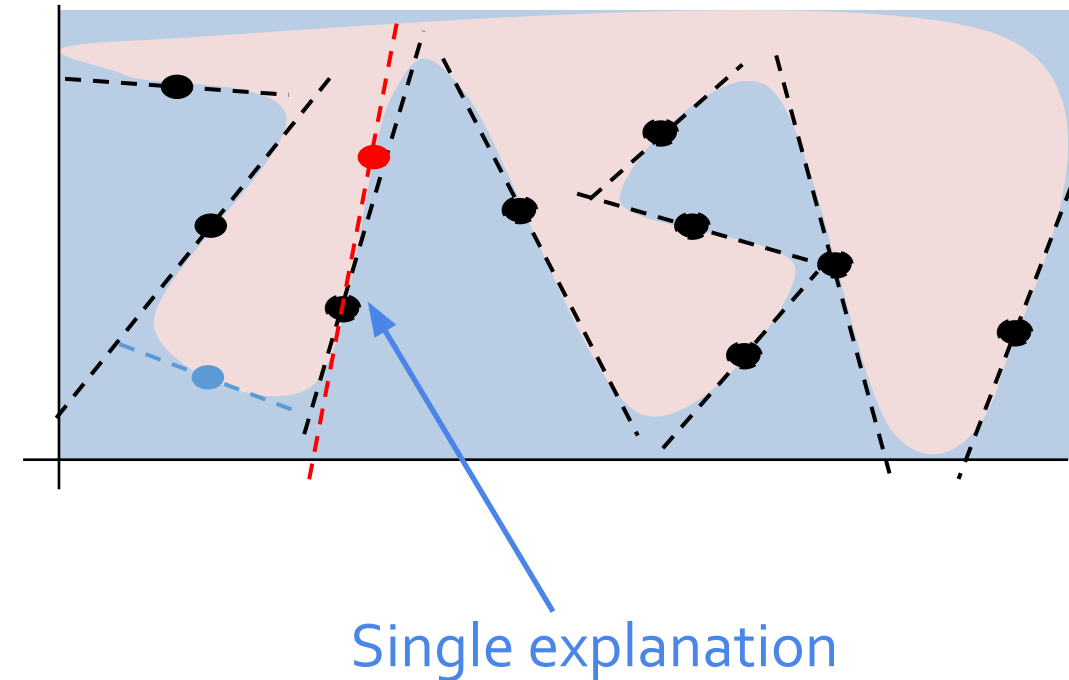
**Representative**

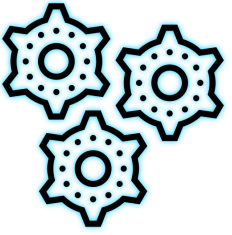
Should summarize the  
model's global behavior

**Diverse**

Should not be redundant in  
their descriptions

SP-LIME uses submodular optimization  
and *greedily* picks  $k$  explanations





# Approaches for Post hoc Explainability

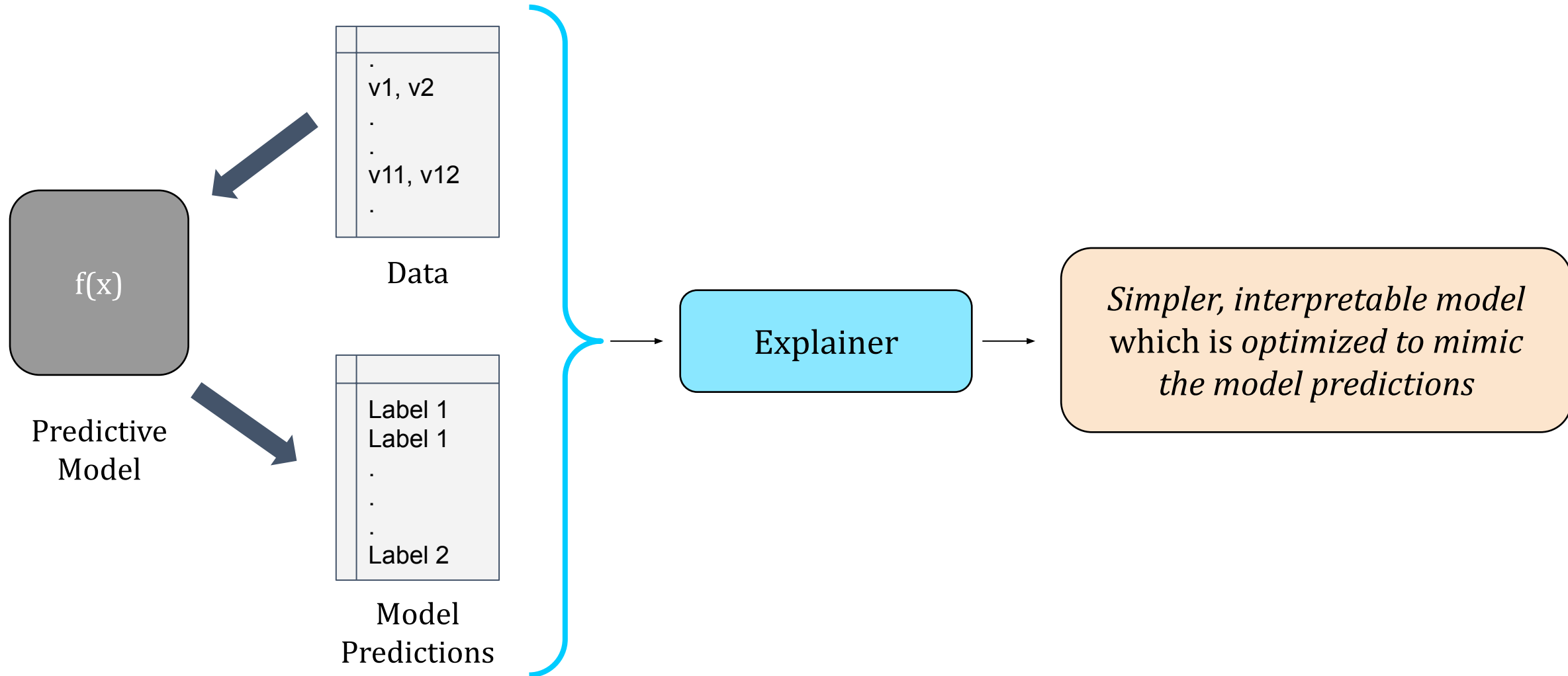
## Local Explanations

- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

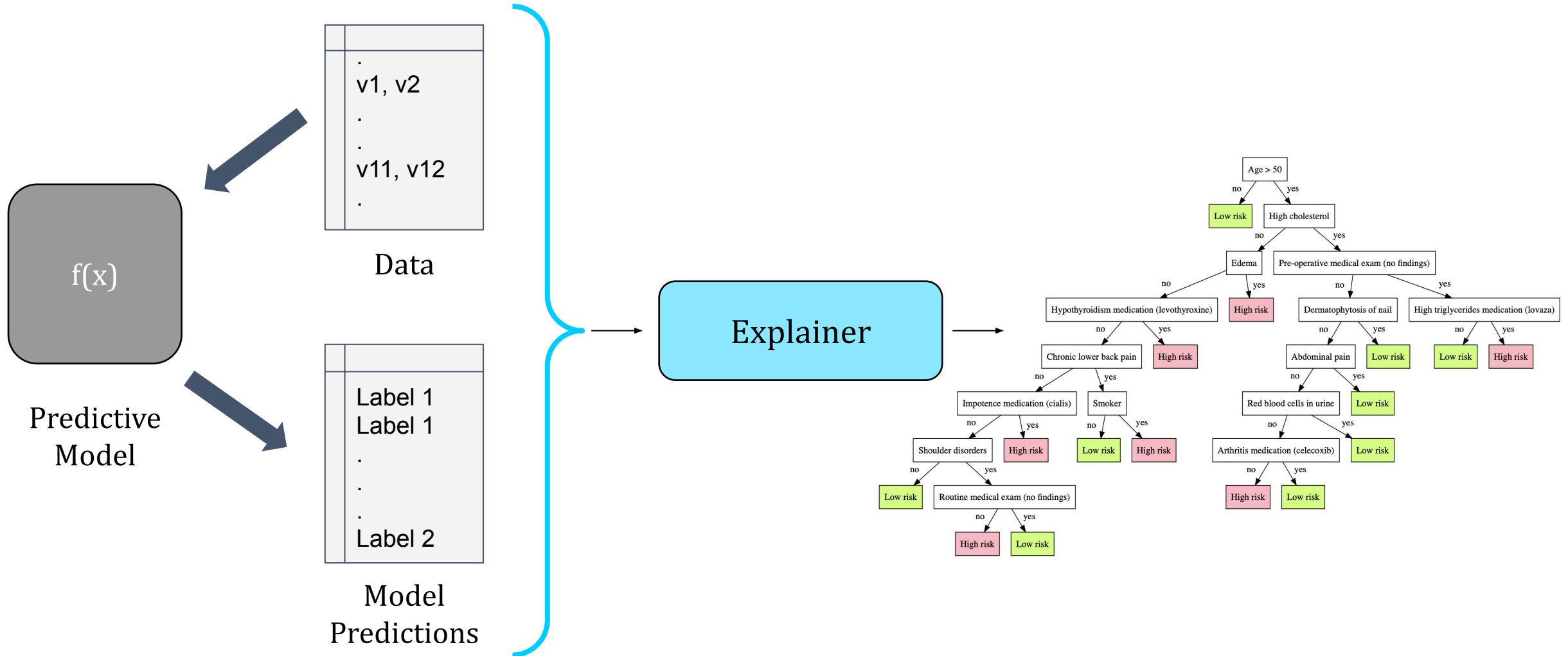
## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Model Distillation for Generating Global Explanations

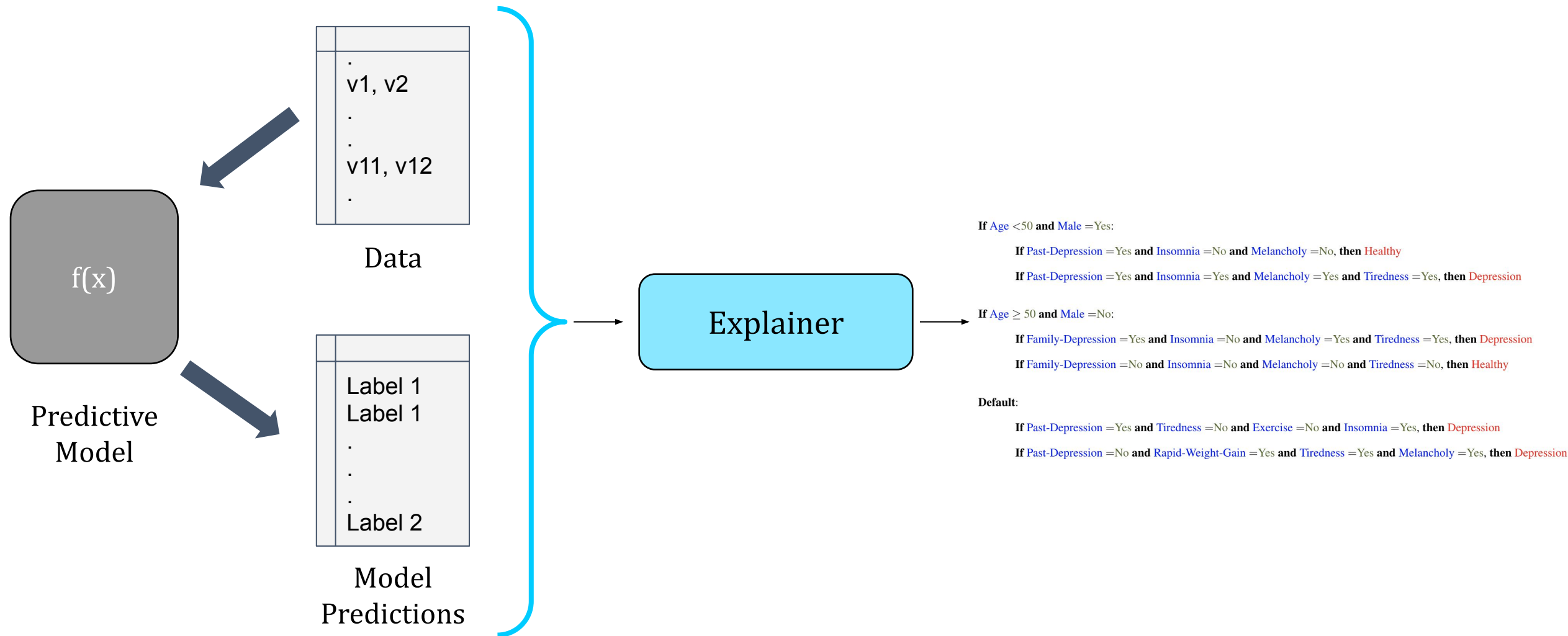


# Decision Trees as Global Explanations

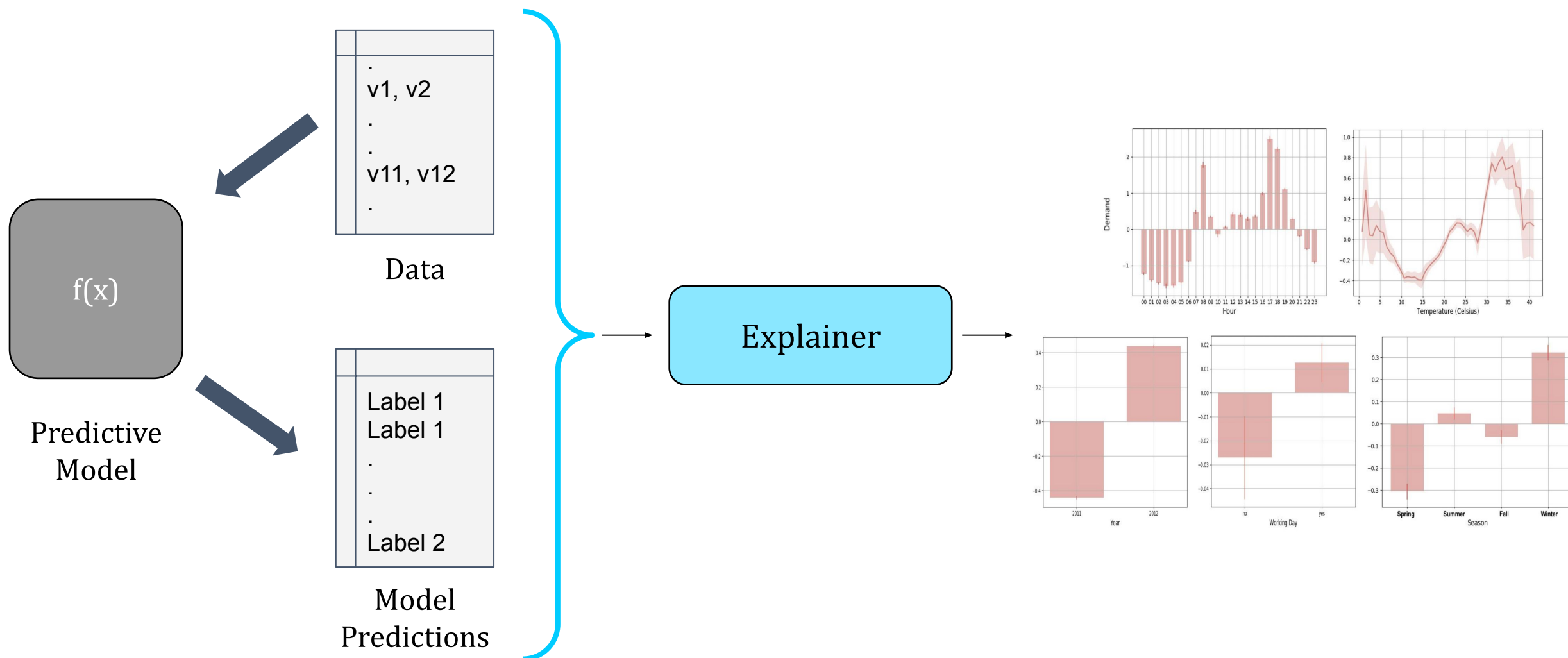


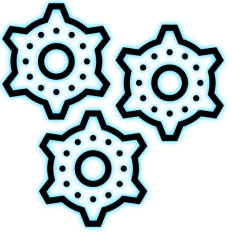


# Customizable Decision Sets as Global Explanations



# Generalized Additive Models as Global Explanations





# Approaches for Post hoc Explainability

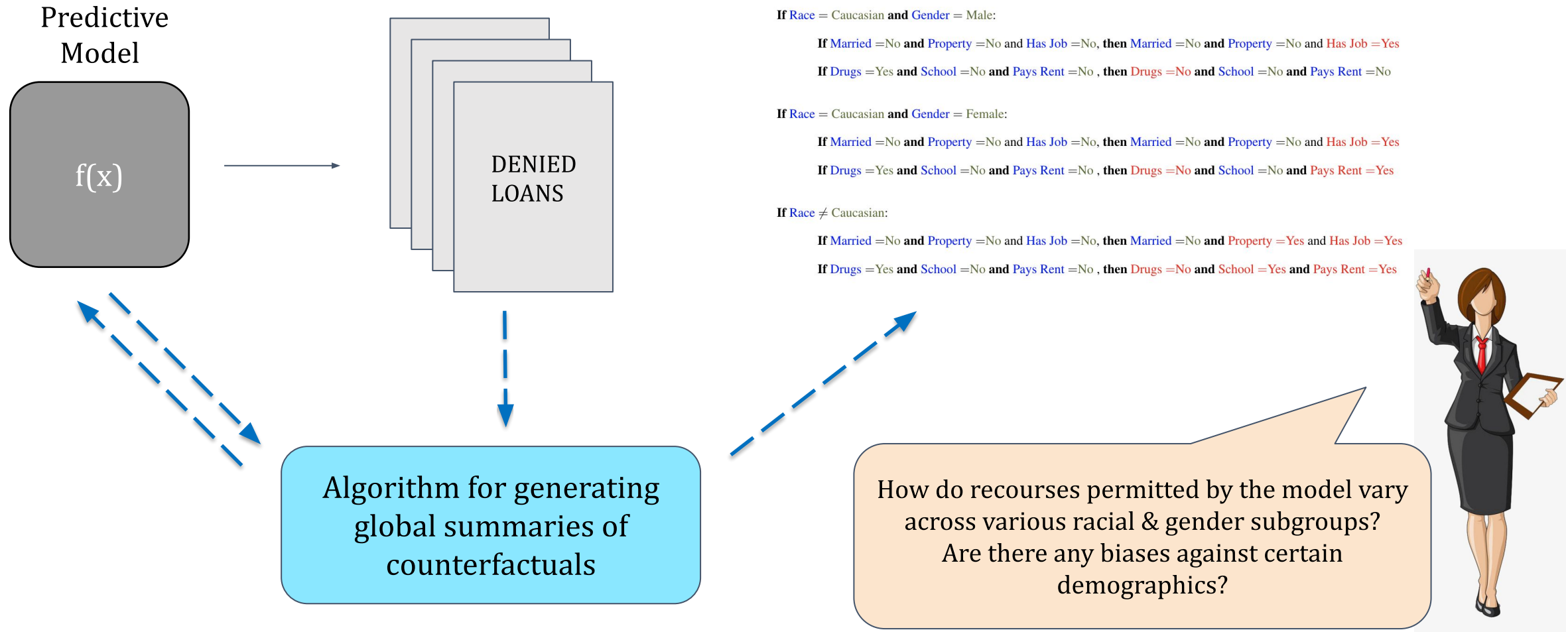
## Local Explanations

- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Customizable Global Summaries of Counterfactuals



# Customizable Global Summaries of Counterfactuals

## Subgroup Descriptor

If Race = Caucasian and Gender = Male:

If Married = No and Property = No and Has Job = No, then Married = No and Property = No and Has Job = Yes

If Drugs = Yes and School = No and Pays Rent = No, then Drugs = No and School = No and Pays Rent = No

If Race = Caucasian and Gender = Female:

If Married = No and Property = No and Has Job = No, then Married = No and Property = No and Has Job = Yes

If Drugs = Yes and School = No and Pays Rent = No, then Drugs = No and School = No and Pays Rent = Yes

If Race ≠ Caucasian:

If Married = No and Property = No and Has Job = No, then Married = No and Property = Yes and Has Job = Yes

If Drugs = Yes and School = No and Pays Rent = No, then Drugs = No and School = Yes and Pays Rent = Yes

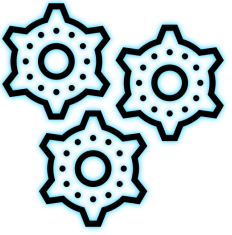
Omg! this model is biased. It requires certain demographics to “change” lot more features than others.

## Recourse Rules



# Customizable Global Summaries of Counterfactuals

- An optimization problem which is *non-negative*, *non-normal*, *non-monotone*, and *submodular* with *matroid constraints*
- Solved using the well-known *smooth local search* algorithm (Feige et. al., 2007) with best known optimality guarantees.



# Approaches for Post hoc Explainability

## Local Explanations

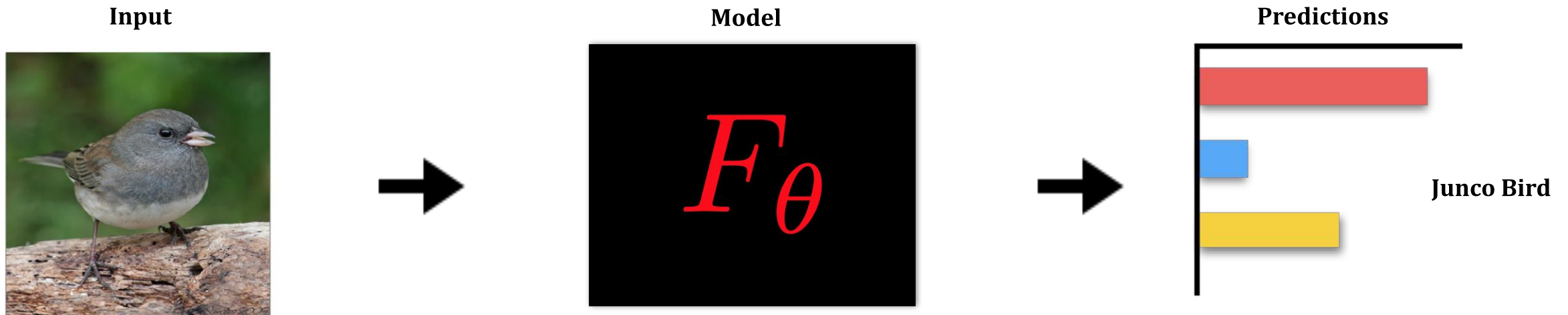
- Feature Importances
- Rule Based
- Saliency Maps
- Prototypes/Example Based
- Counterfactuals

## Global Explanations

- Collection of Local Explanations
- Model Distillation
- Summaries of Counterfactuals
- Representation Based

# Representation Based Approaches

- Derive model understanding by analyzing intermediate representations of a DNN.
- Determine model's reliance on 'concepts' that are semantically meaningful to humans.



Does the model rely on the **'green background'**?



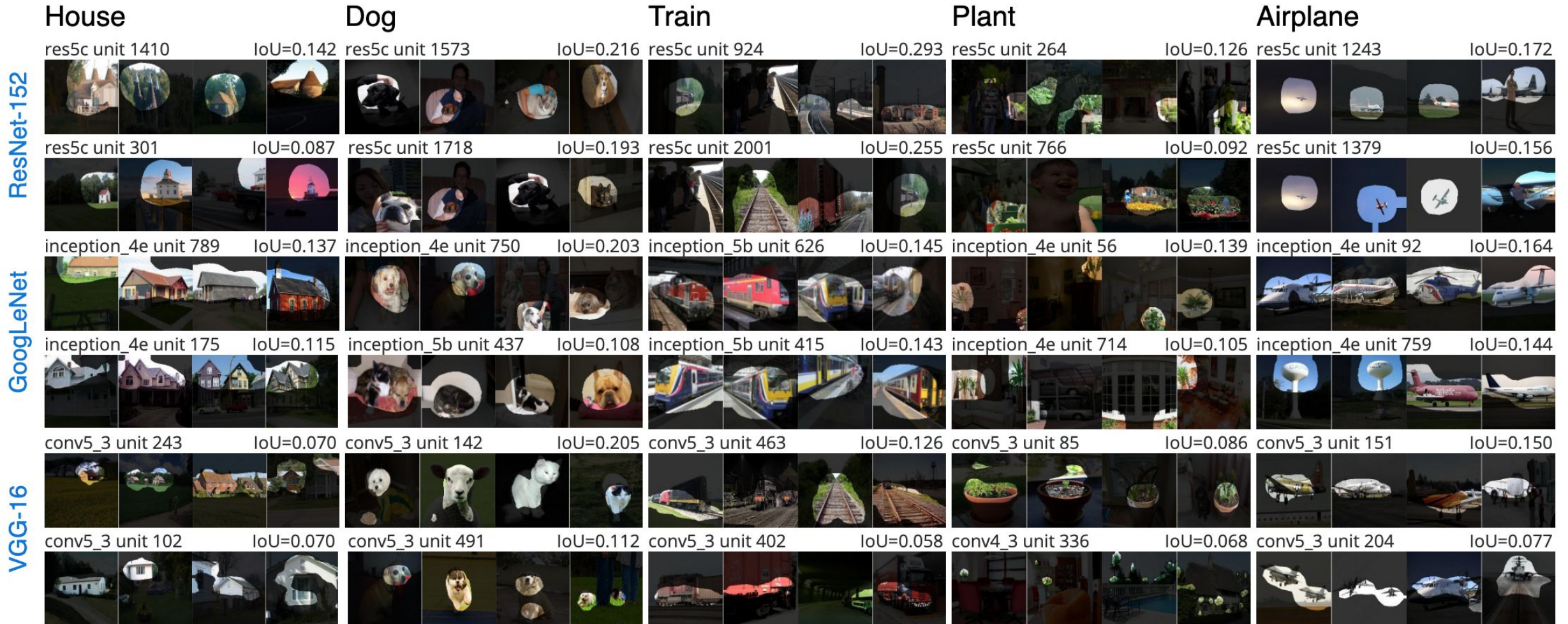
# Representation Based Approaches

- Network Dissection ([Bau & Zhou et. al. 2017](#)).
- TCAV ([Kim et. al. 2018](#)).

## Process

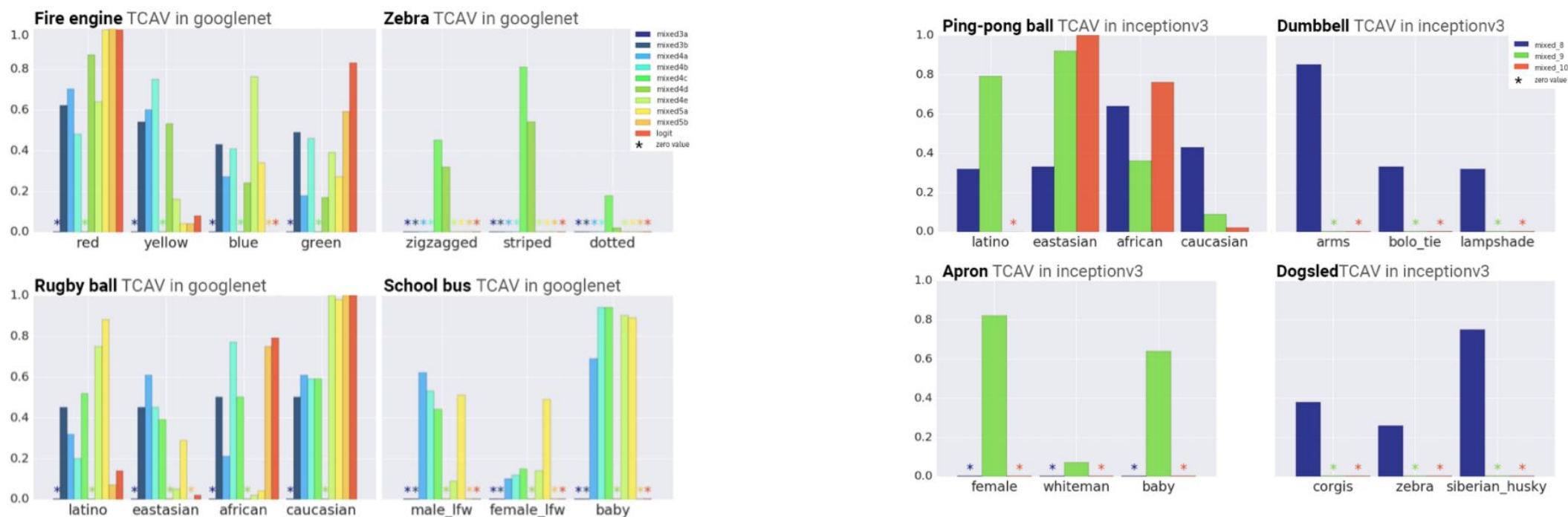
1. Identify human-labeled concepts.
2. Gather the response of hidden variables (convolutional filters) to known concepts.
3. Quantify alignment of hidden variable-concept pairs

# Network Dissection



# Quantitative Testing with Concept Activation Vectors (TCAV)

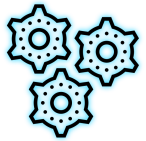
## Insights from Googlenet and Inception-v3



### Additional Variants:

- Regression problems in medical domain ([Graziani et. al. 2019](#)).
- Automatic extraction of visual concepts ([Ghorbani et. al. 2019](#)).

# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations

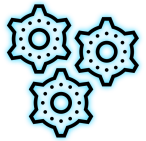


**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability

# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations



**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability

# Evaluation of Post hoc Explanations





# How we evaluate explanations?





# Evaluating Post hoc Explanations

Understand the Behavior

Help make decisions

Useful for Debugging





# Evaluating Post hoc Explanations

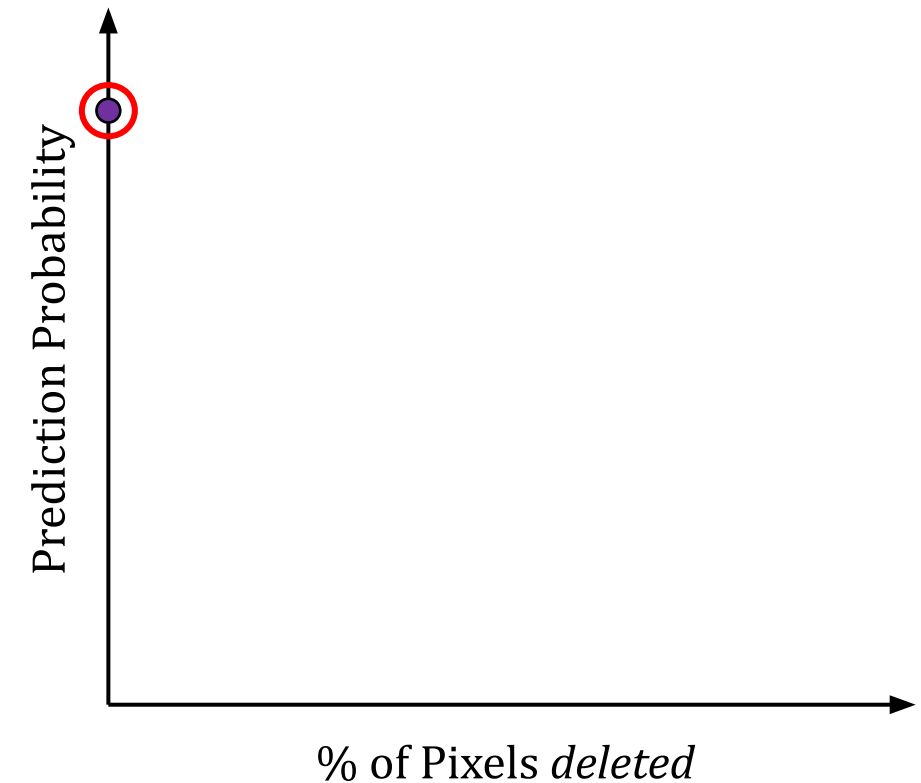
Understand the Behavior

Help make decisions

Useful for Debugging

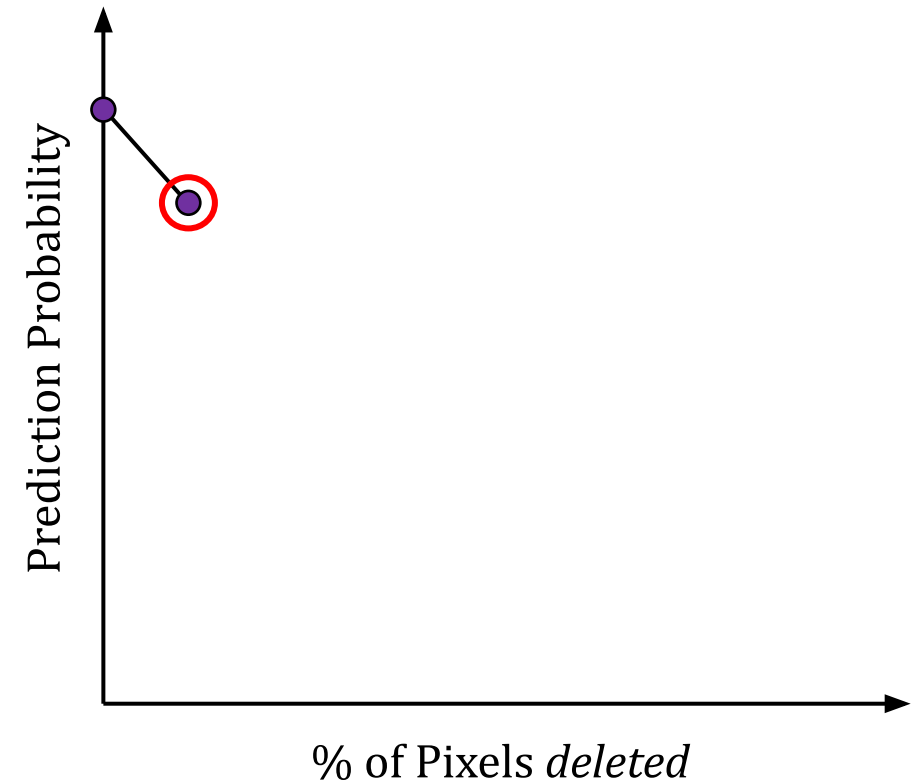
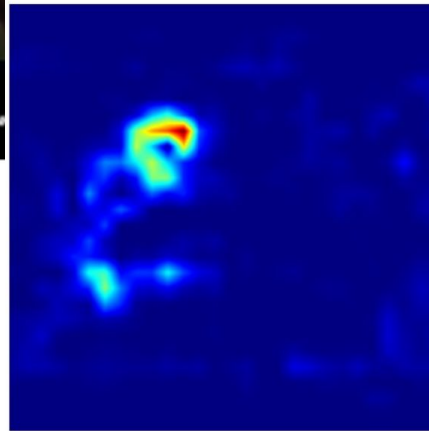
# How important are selected features?

- **Deletion**: remove important features and see what happens..



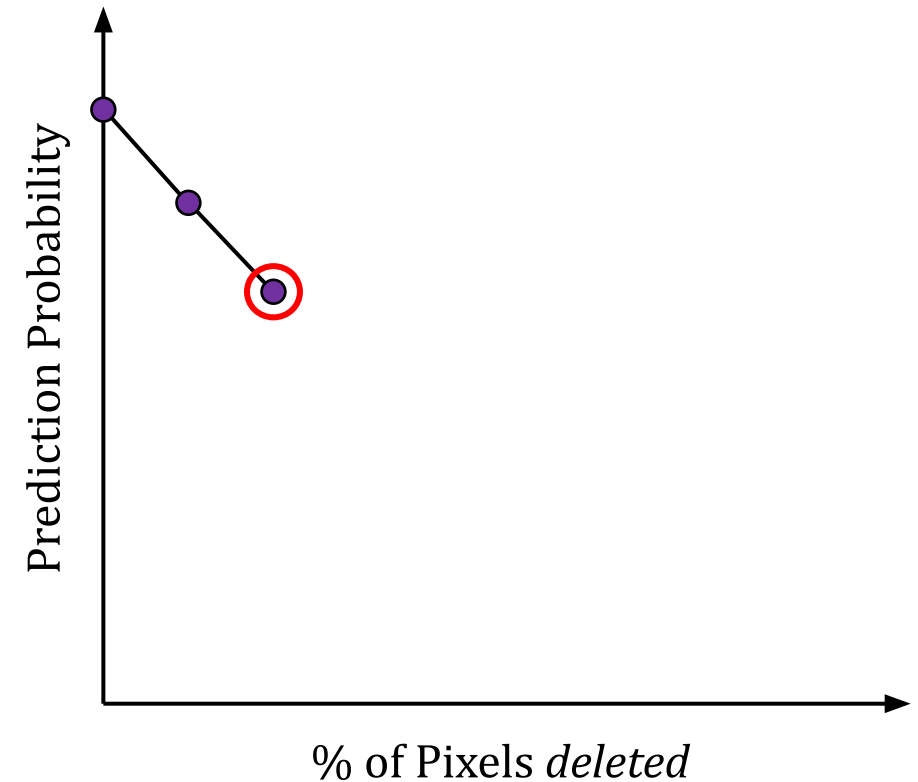
# How important are selected features?

- **Deletion:** remove important features and see what happens..



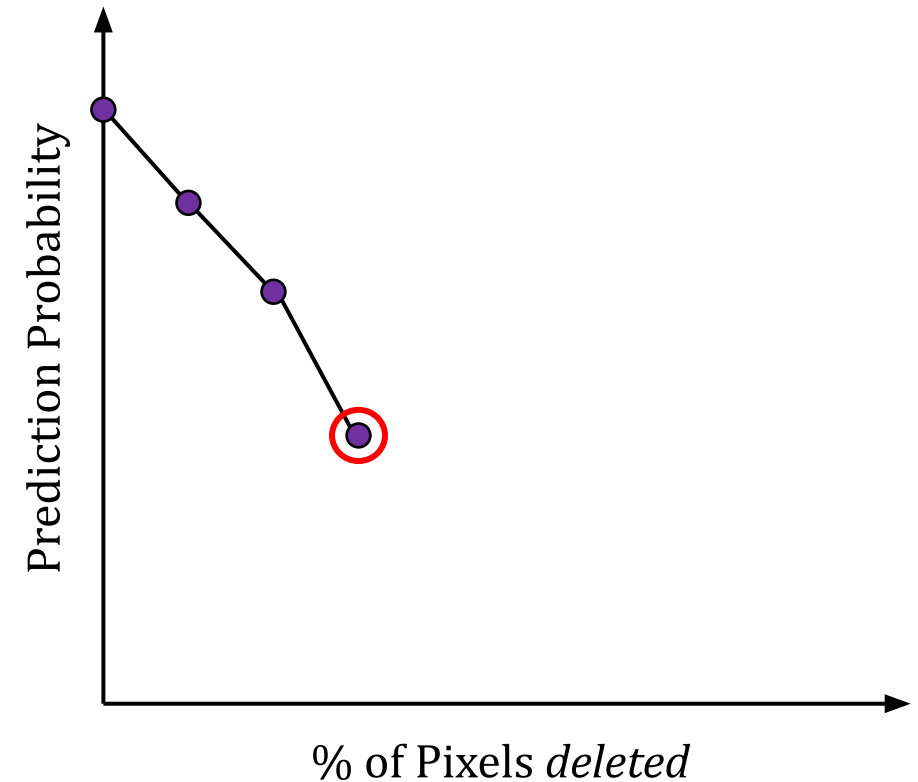
# How important are selected features?

- **Deletion:** remove important features and see what happens..



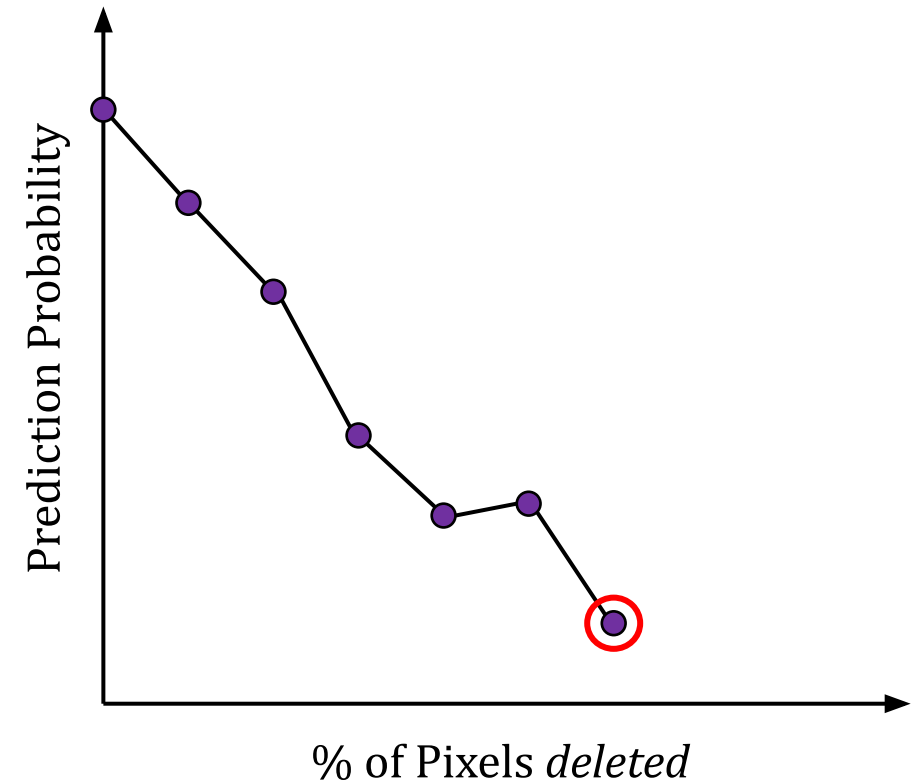
# How important are selected features?

- **Deletion:** remove important features and see what happens..



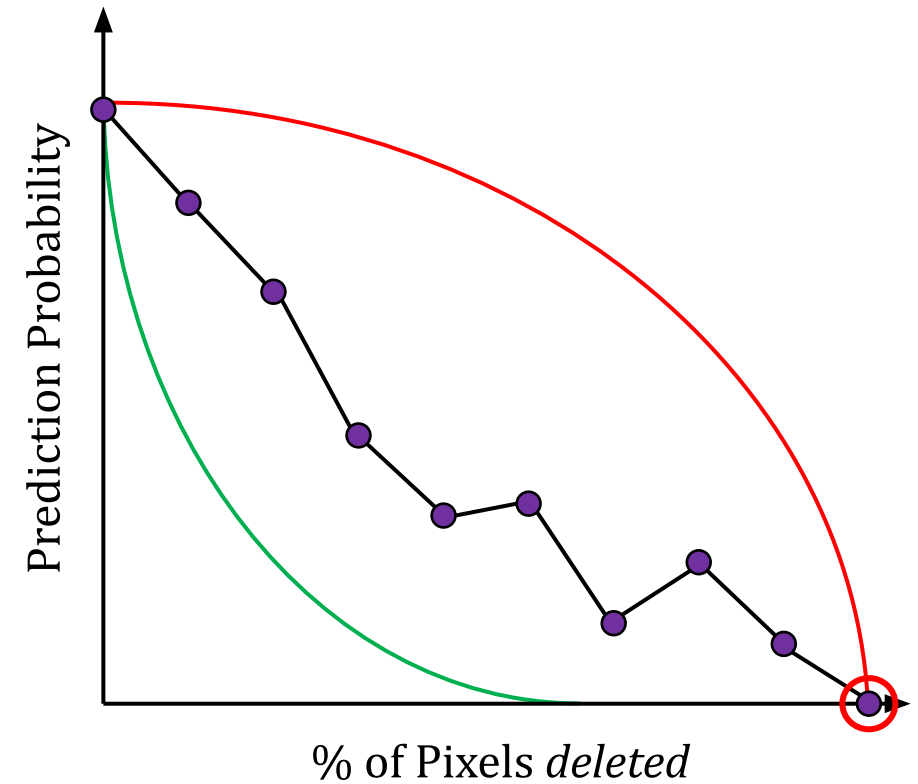
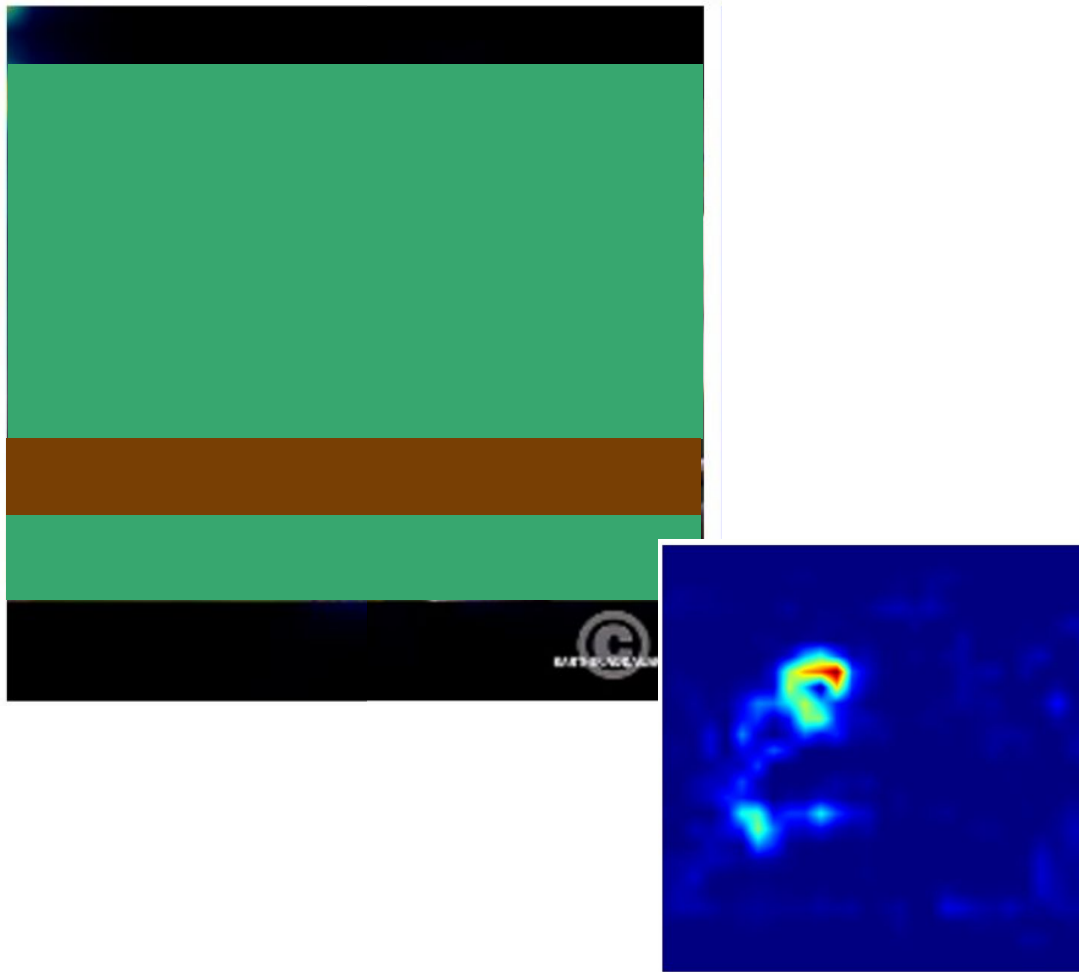
# How important are selected features?

- **Deletion:** remove important features and see what happens..



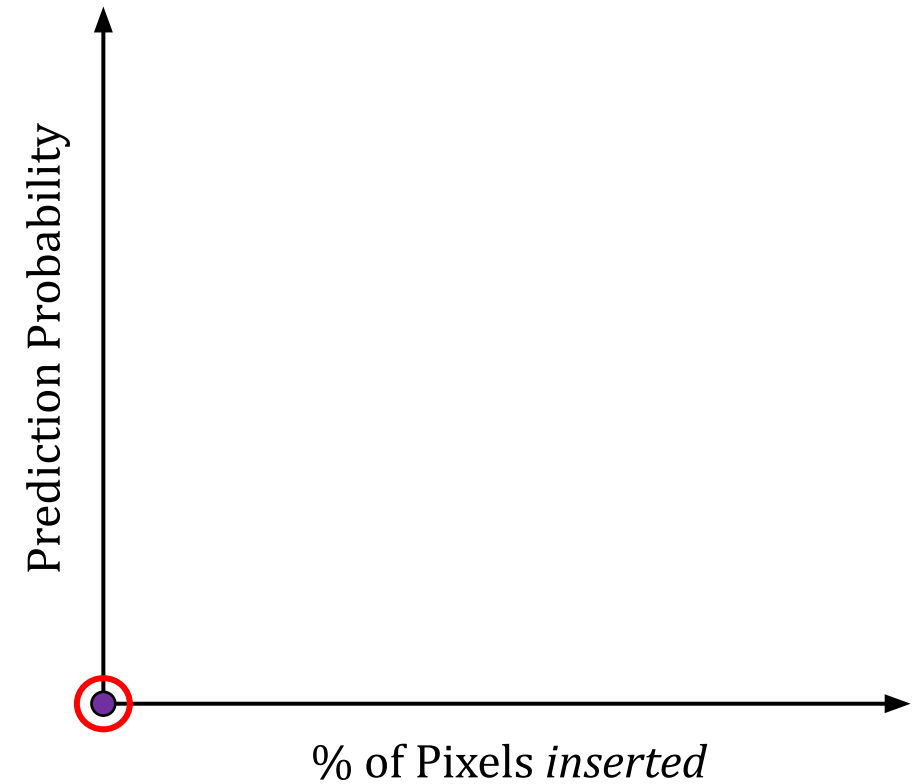
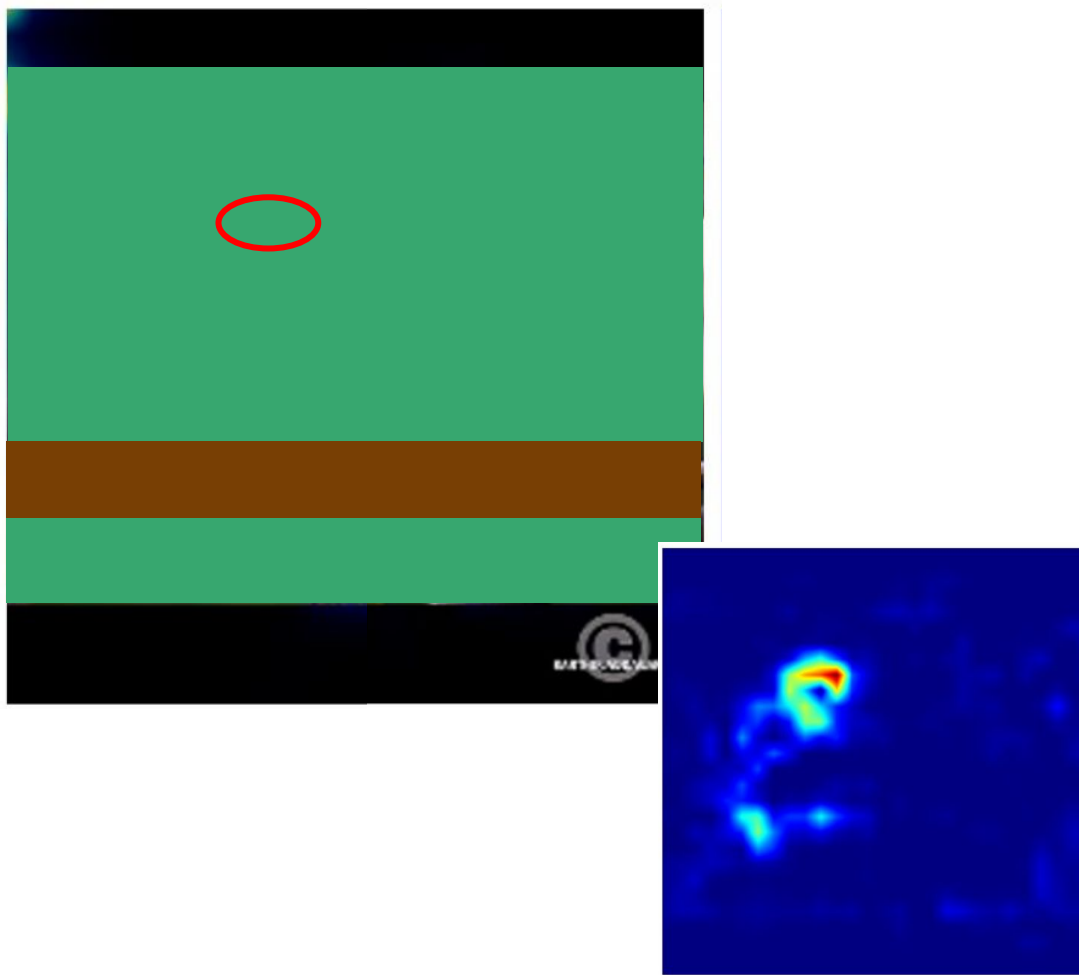
# How important are selected features?

- **Deletion:** remove important features and see what happens..



# How important are selected features?

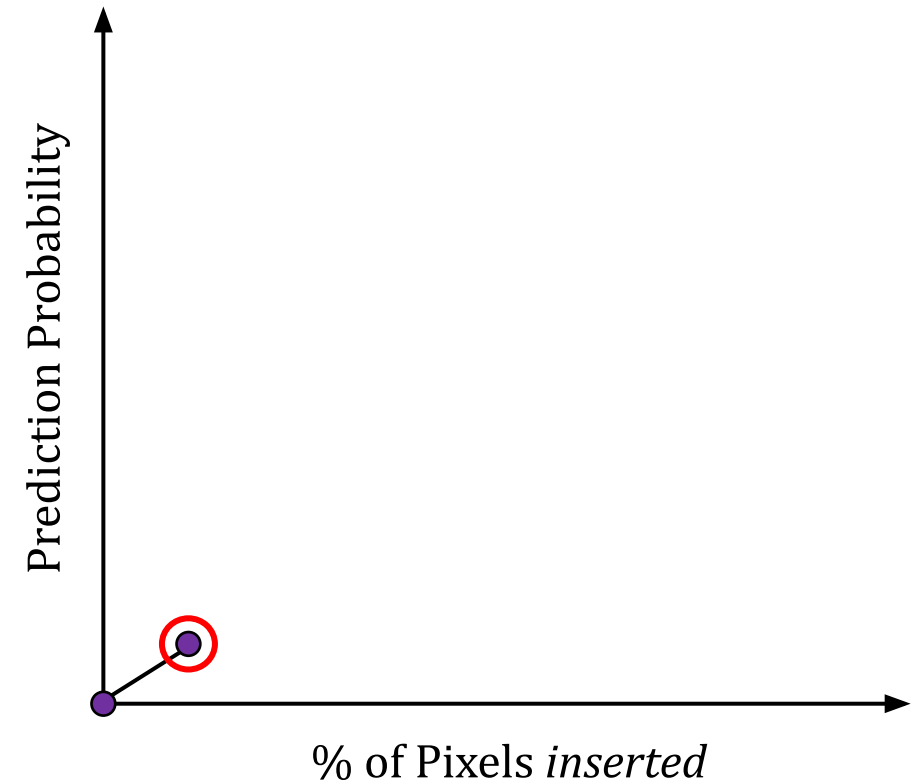
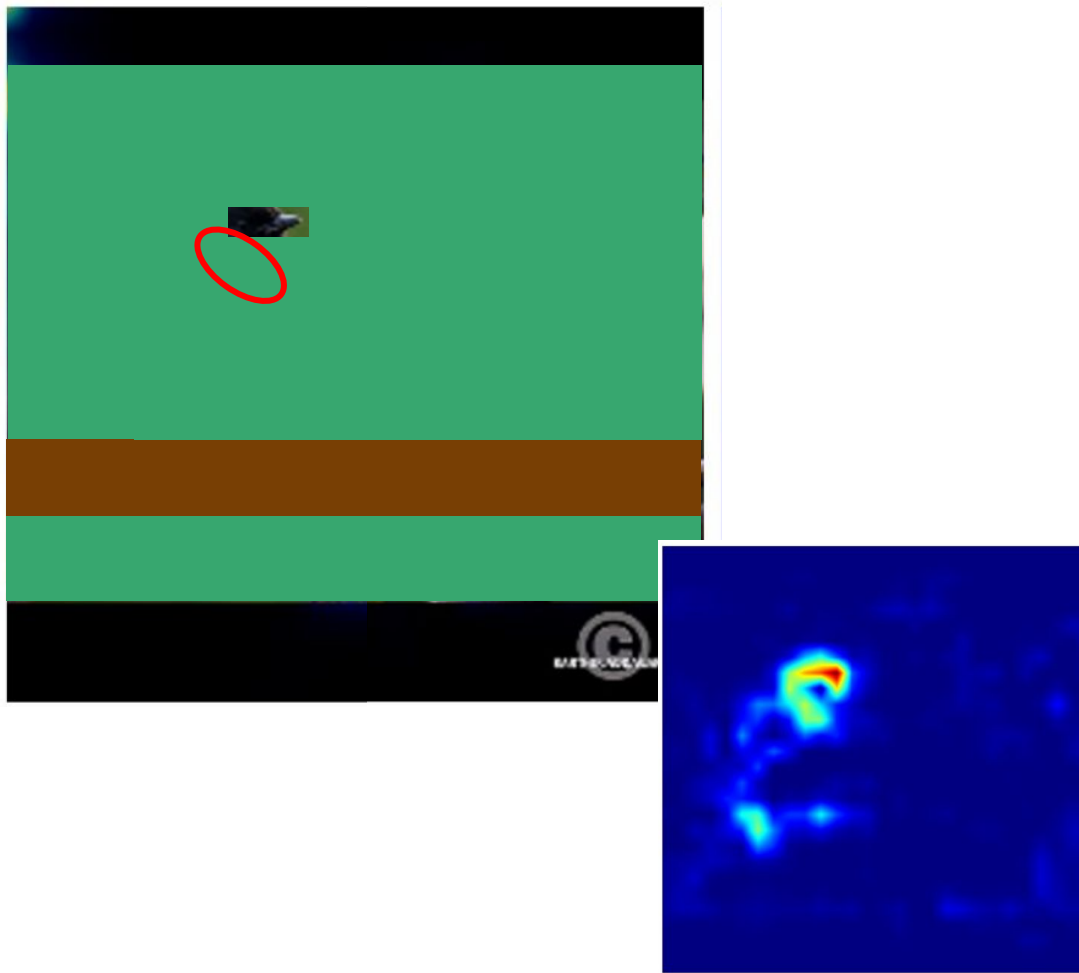
- **Insertion:** add important features and see what happens..





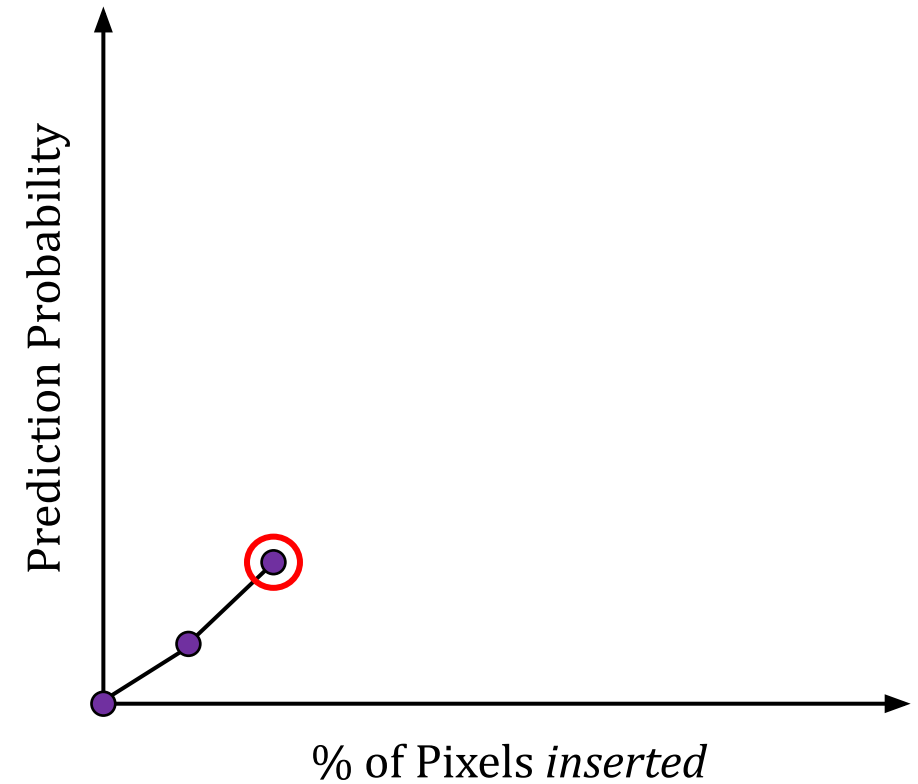
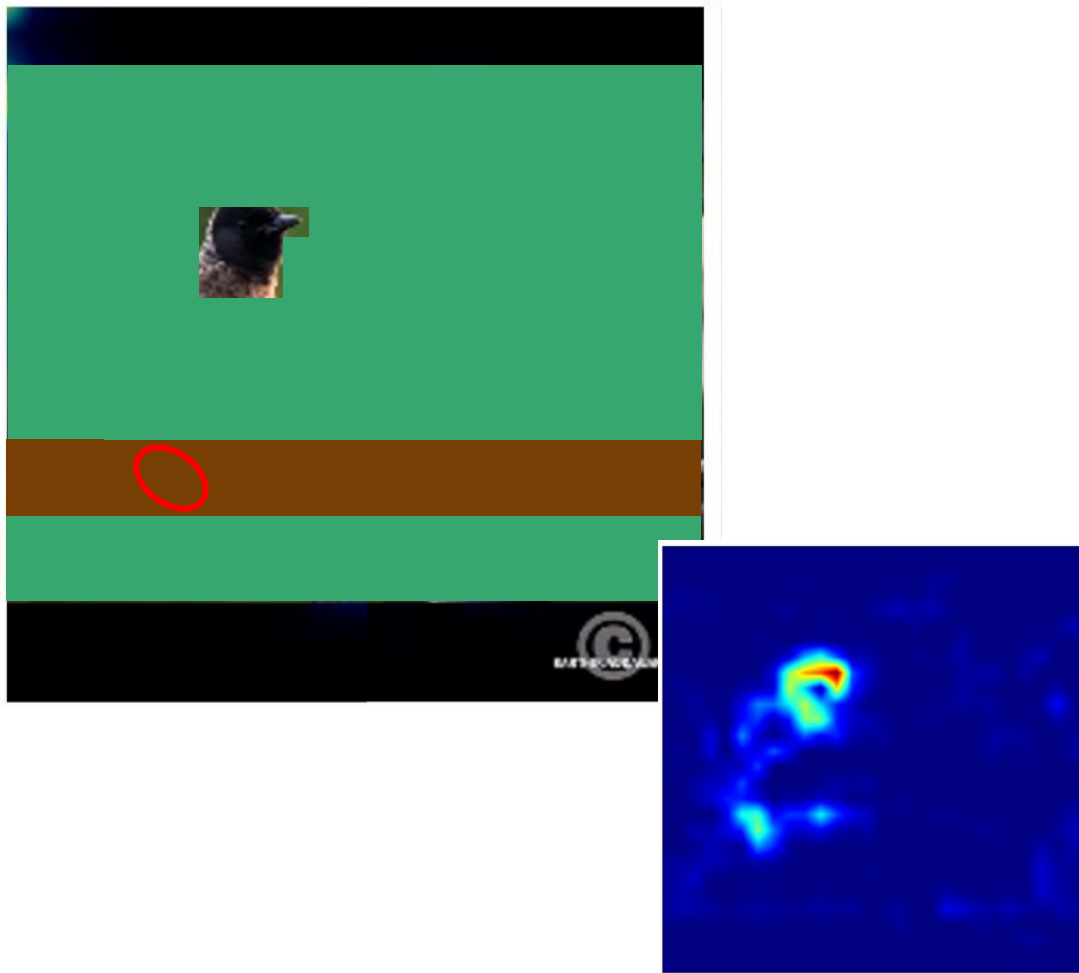
# How important are selected features?

- **Insertion:** add important features and see what happens..



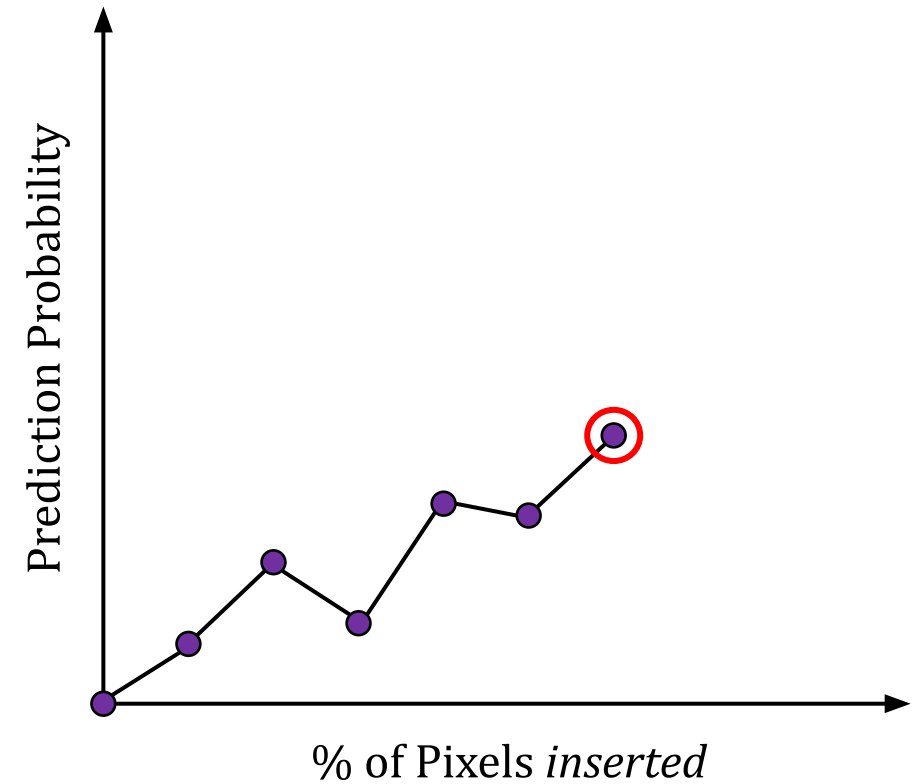
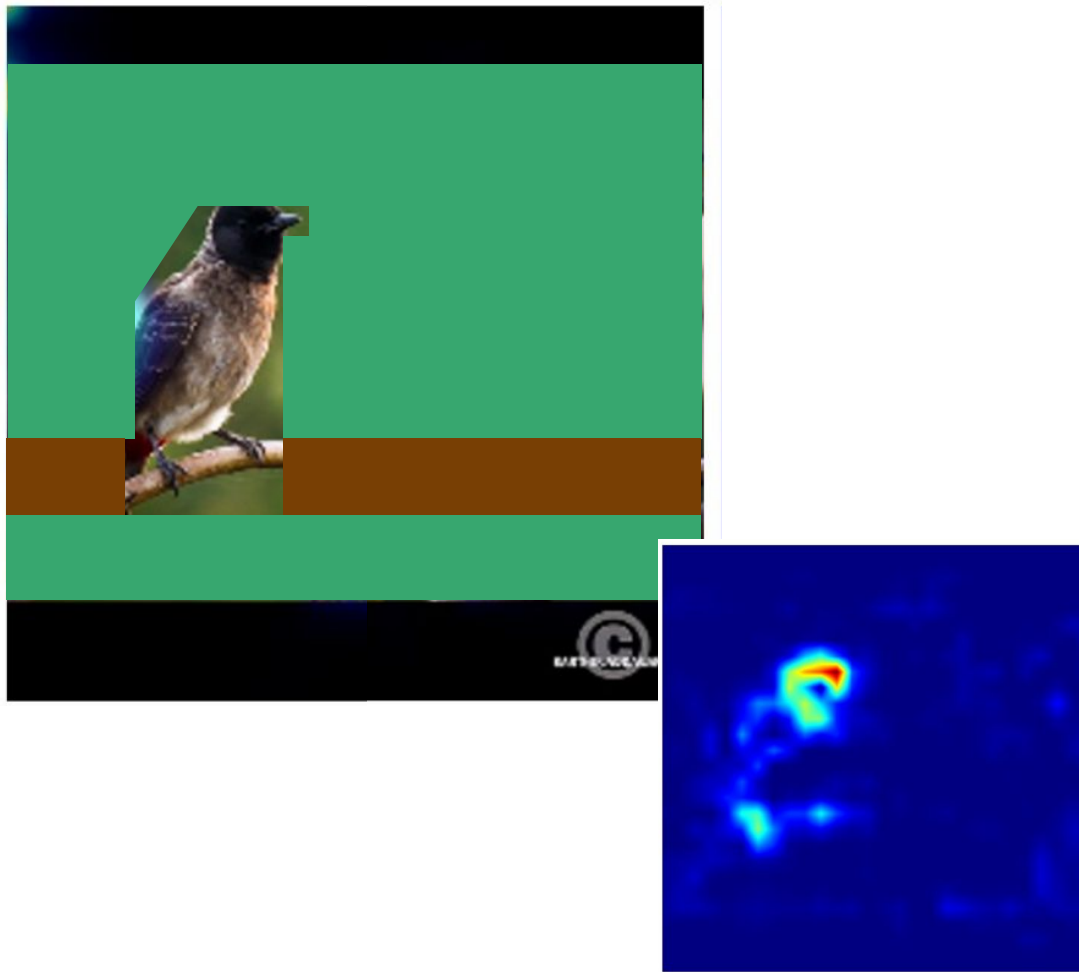
# How important are selected features?

- **Insertion:** add important features and see what happens..



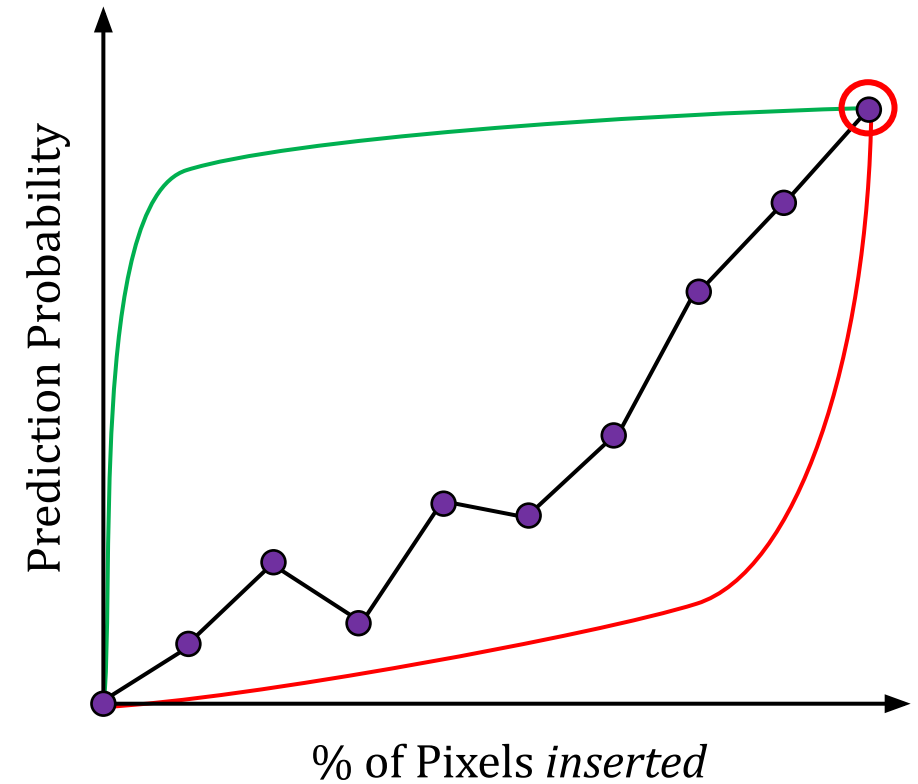
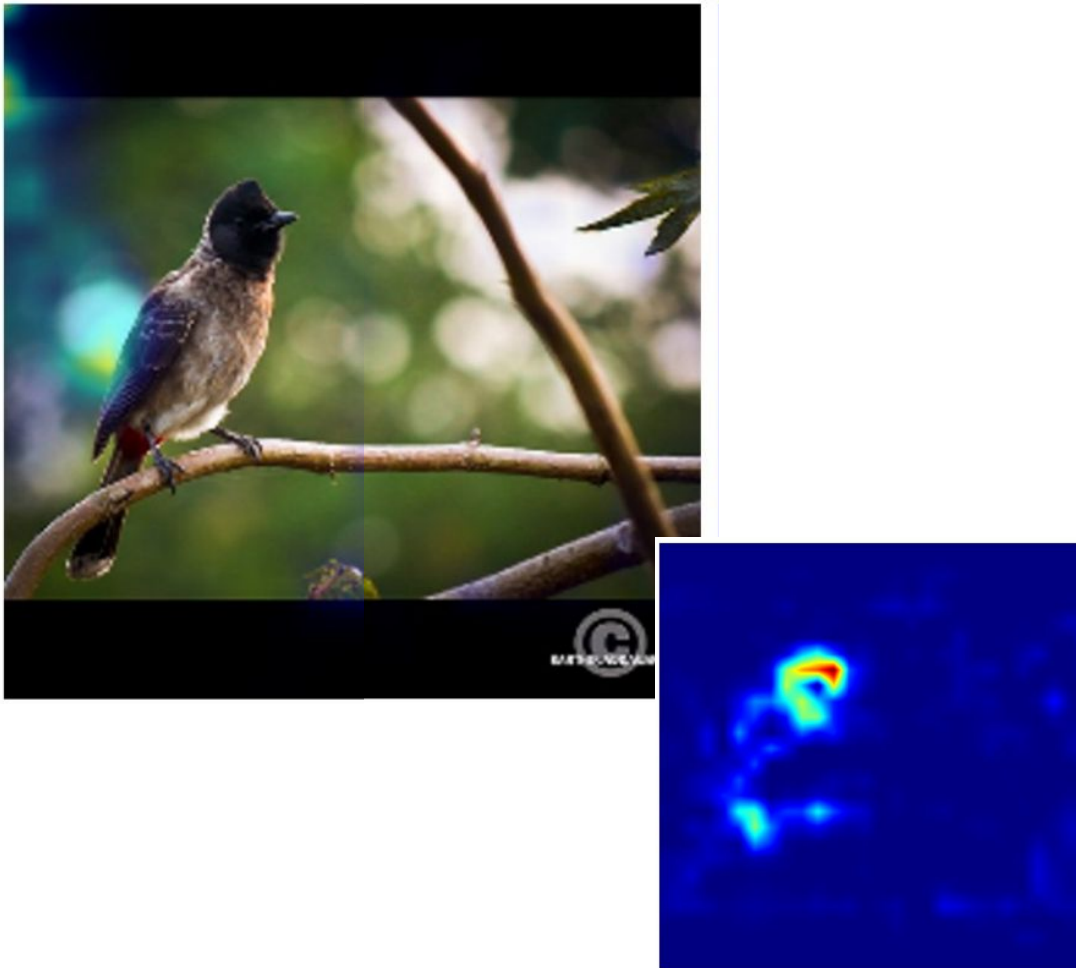
# How important are selected features?

- **Insertion:** add important features and see what happens..



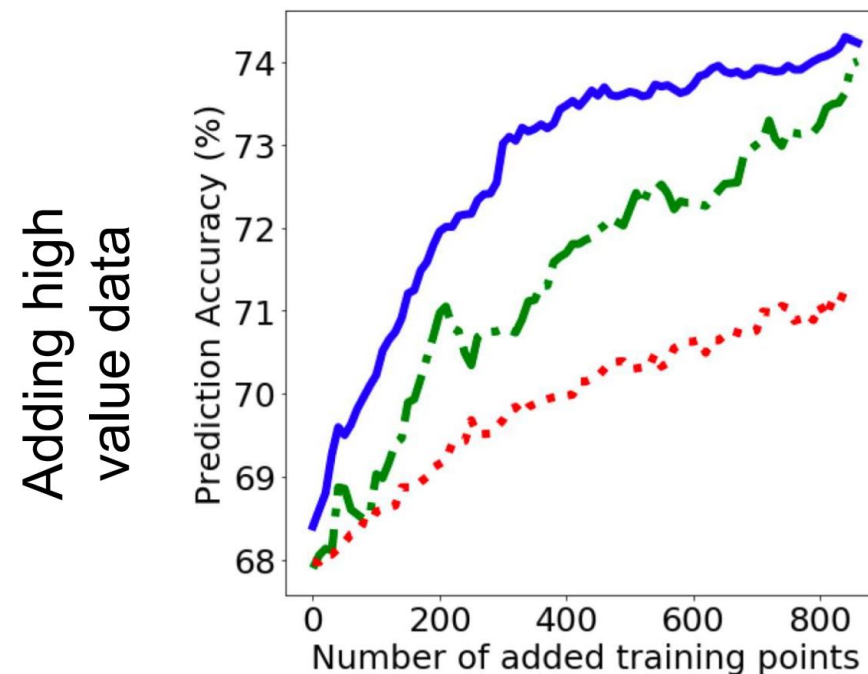
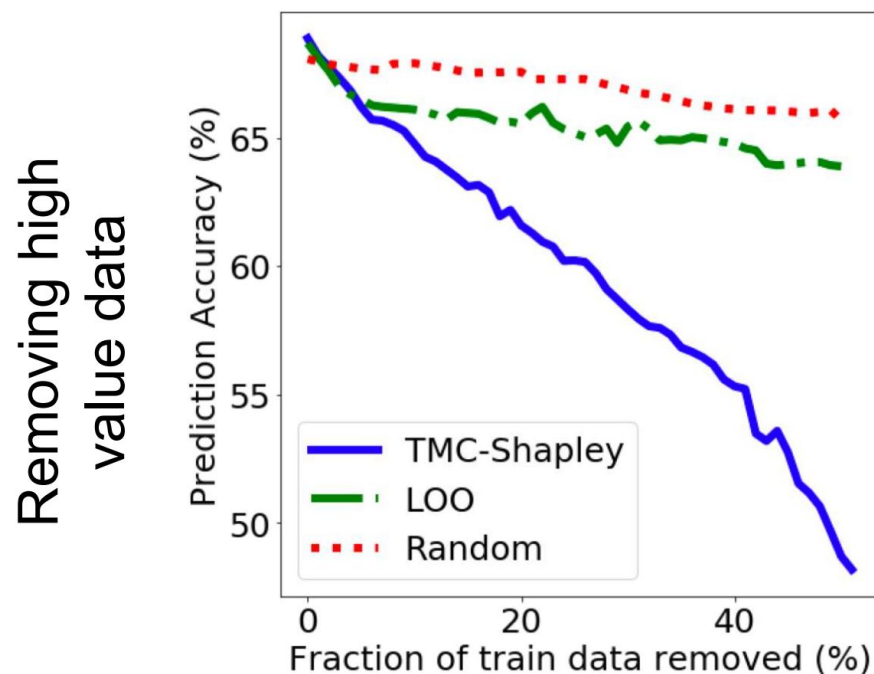
# How important are selected features?

- **Insertion:** add important features and see what happens..

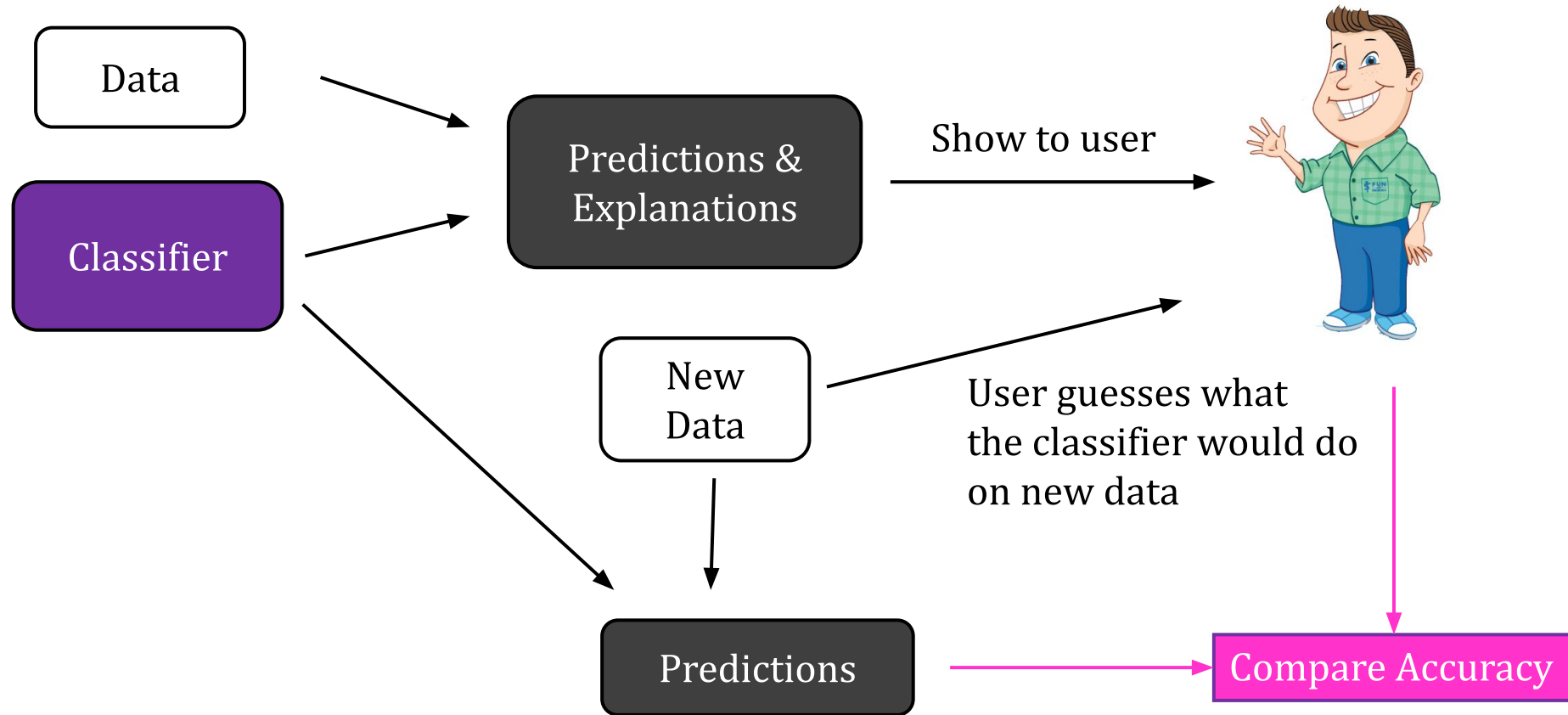


# Same Idea: For *Training* Data

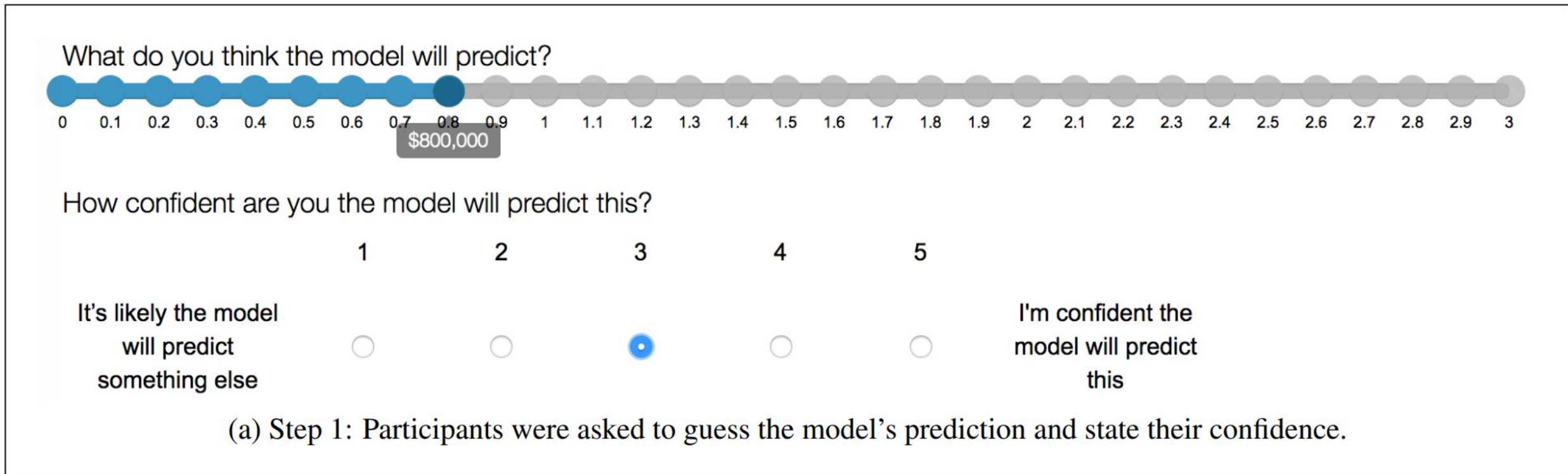
Add/remove **influential** training data, see what happens



# Predicting Behavior (“Simulation”)



# Predicting Behavior (“Simulation”)





# Evaluating Post hoc Explanations

Understand the Behavior

Help make decisions

Useful for Debugging



# 1. Detecting Problems in Classifiers



## Question 1

Would you trust this model?

Did they say no?

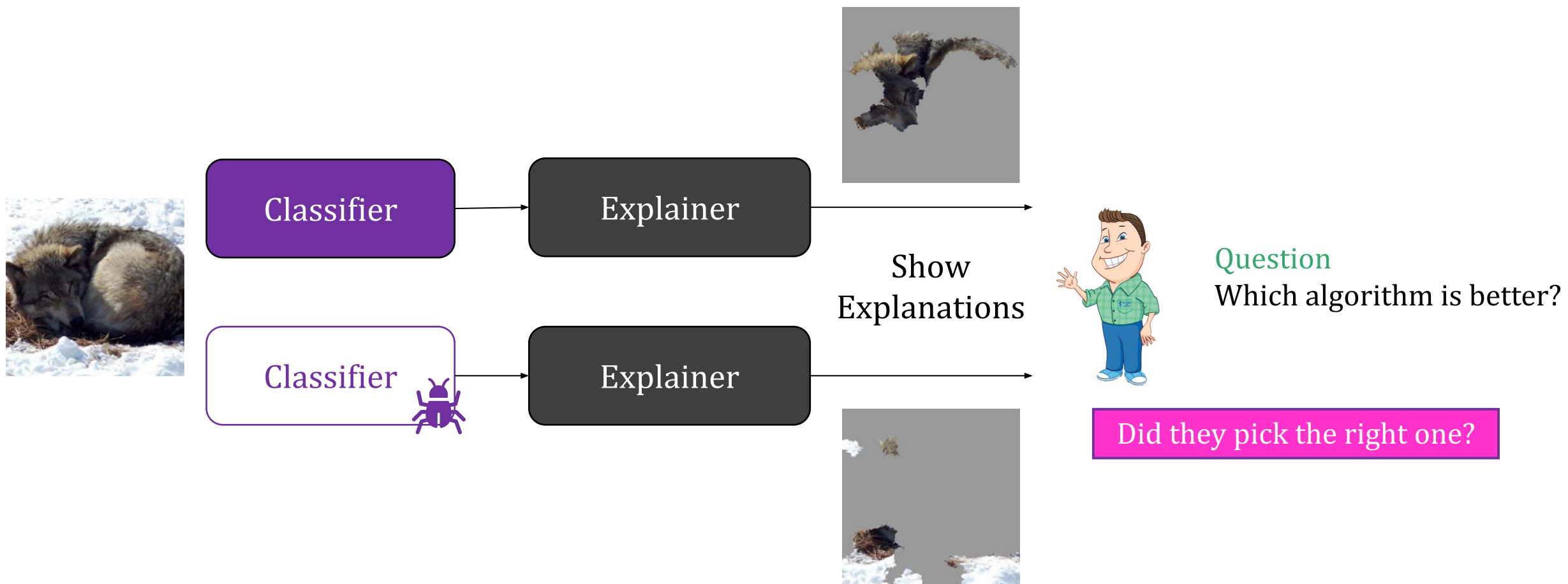
## Question 2

What is the classifier doing?

Did they get it right?

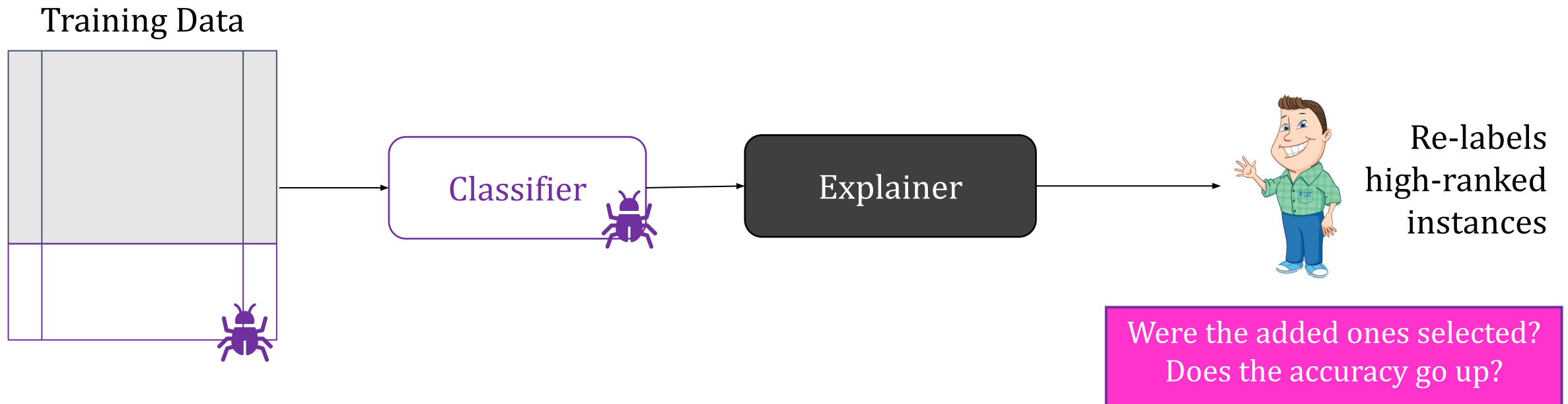


## 2. Comparing Classifiers



# 3. Finding Errors in Training Data

- **Prototypical Explanations:** important instances from training data





# Evaluating Posthoc Explanations

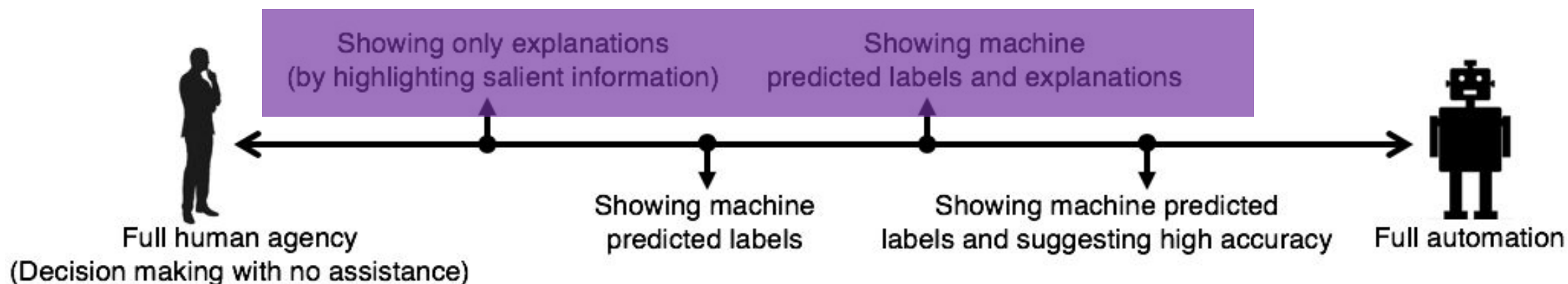
Understand the Behavior

Help make decisions

Useful for Debugging

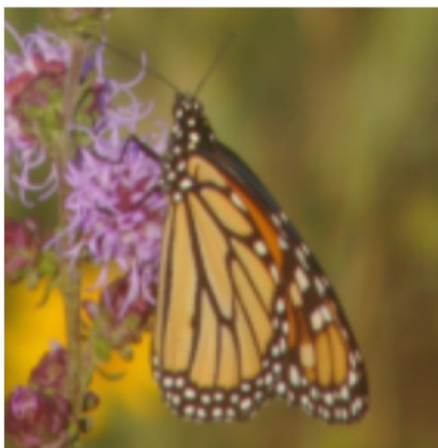
# Human-AI Collaboration

- Are Explanations Useful for Making Decisions?
  - For tasks where the algorithms are not reliable by themselves

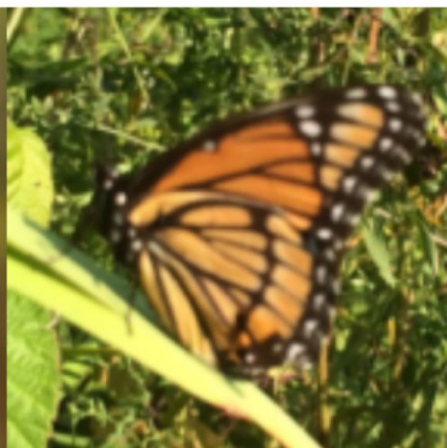


# Machine Teaching

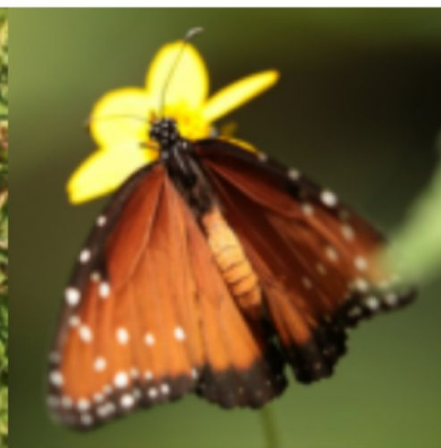
Monarch



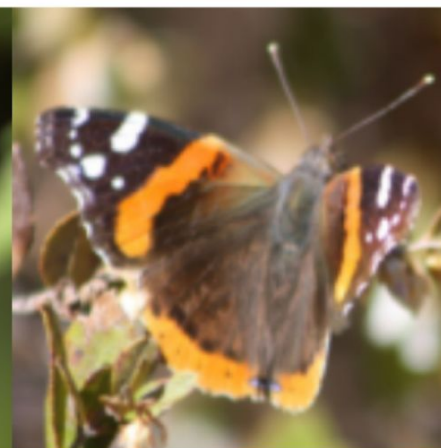
Viceroy



Queen



Red Admiral



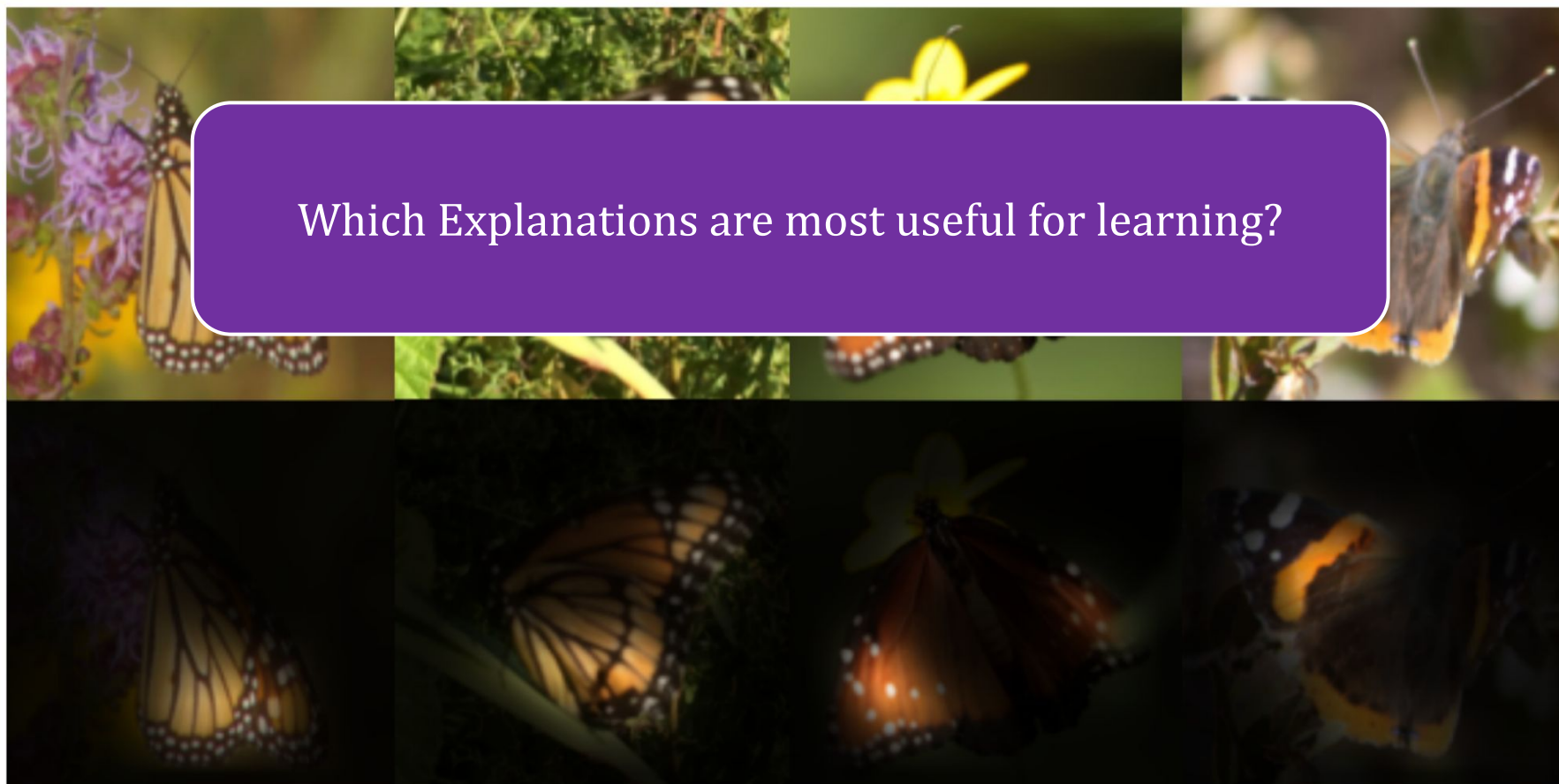
# Machine Teaching

Monarch

Viceroy

Queen

Red Admiral





# Evaluating Posthoc Explanations

Understand the Behavior

Help make decisions

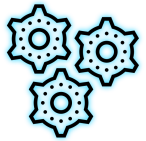
Useful for Debugging



# Limitations of Evaluating Explanations

- Evaluation setup is often **very easy/simple** (or **unrealistic**)
  - E.g. “bugs” are obvious artifacts, classifiers are different from each other
  - Instances/perturbations create out-of-domain points
- Sometimes **flawed**
  - E.g. is model explanation same as human explanation?
- Automated **metrics can be *optimized***
- User studies are **not consistent**
  - Affected by choice of: UI, phrasing, visualization, population, incentives, ...
  - ML researchers are not trained for this 😞
- **Conclusions are difficult to generalize**

# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations

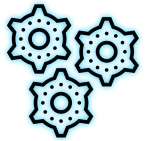


**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability

# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations



**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability

# Limits of Post hoc Explanations



# Limitations

- **Faithfulness/Fidelity**
  - Some explanation methods do not '*reflect*' the underlying model.
- **Fragility**
  - Post-hoc explanations can be easily manipulated.
- **Stability**
  - Slight changes to inputs can cause large changes in explanations.
- **Useful in practice?**
  - Unclear if a data scientist (ML engineer)/end-user can use explanations to isolate errors, improve 'trust' or simulate the model.

# Limitations

- **Faithfulness/Fidelity**
  - Some explanation methods do not '*reflect*' the underlying model.

# Do Explanations Capture Model-based Discriminative Signals?

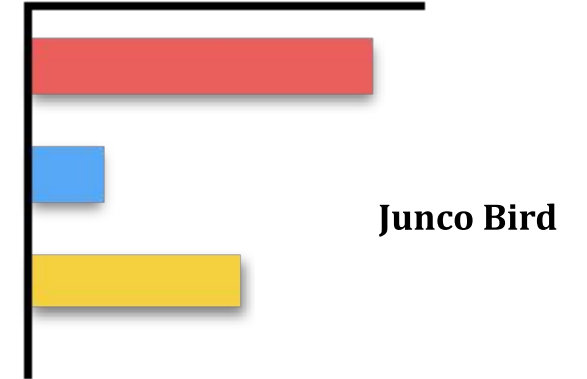
Input



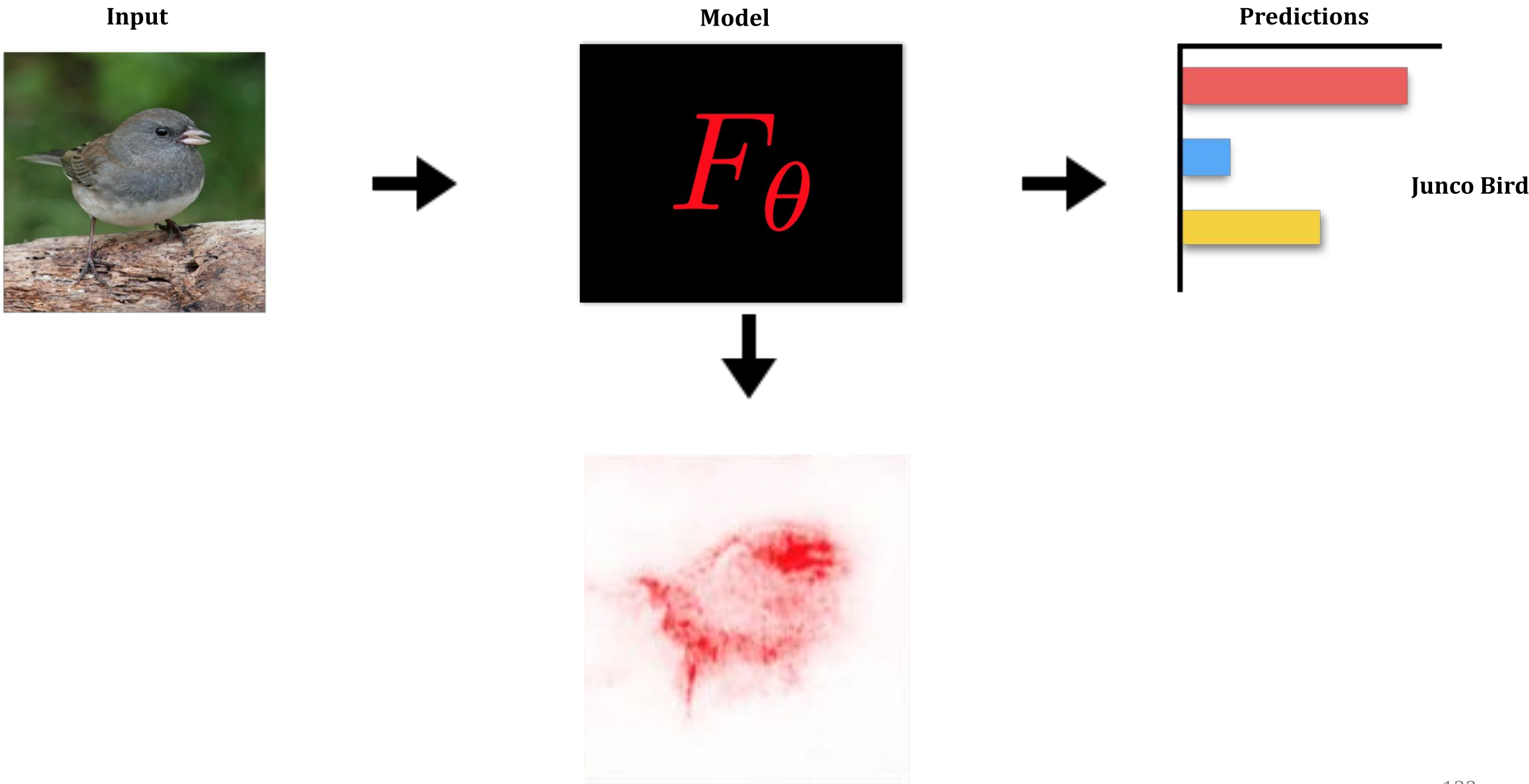
Model



Predictions

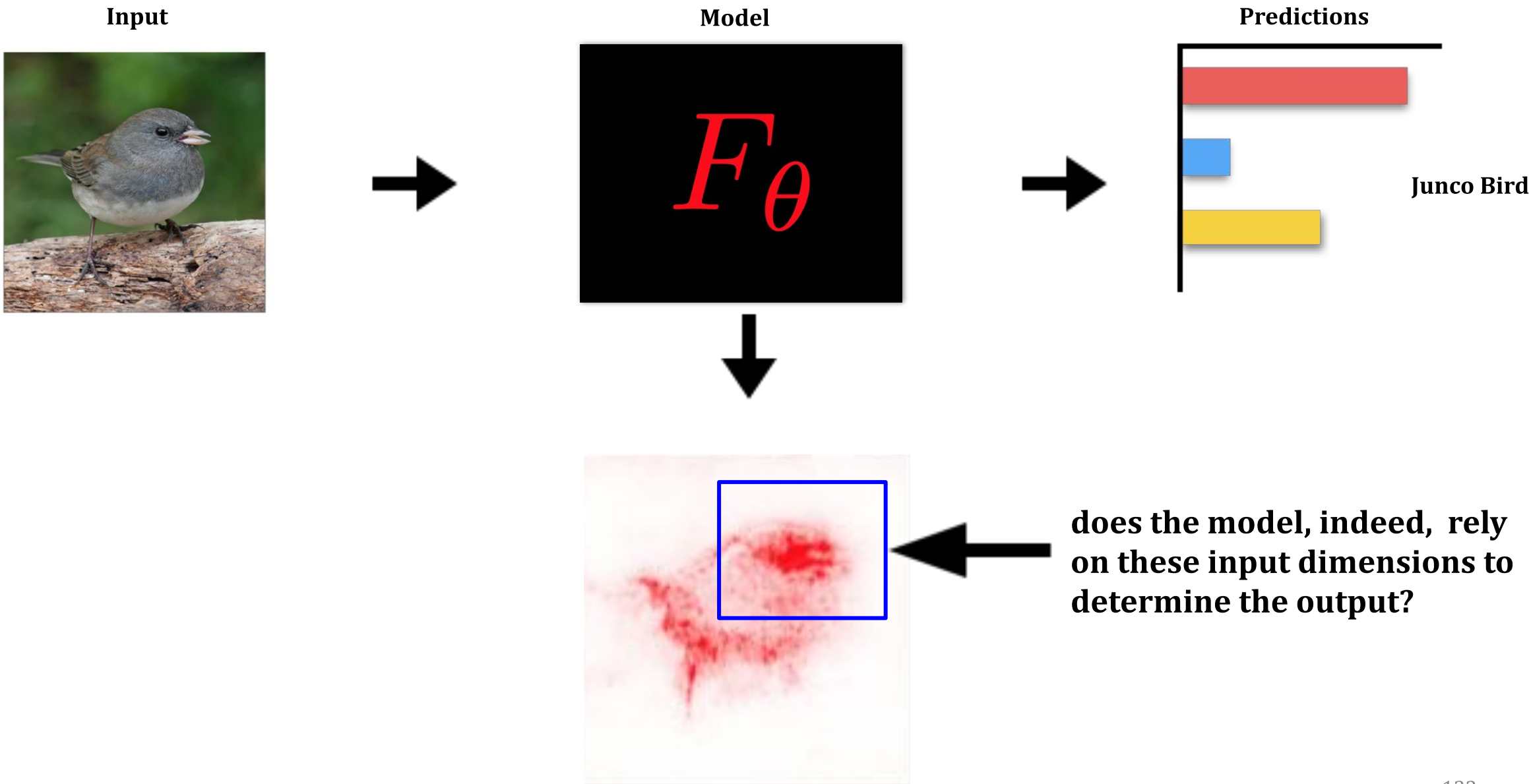


# Do Explanations Capture Model-based Discriminative Signals?





# Do Explanations Capture Model-based Discriminative Signals?



# Faithfulness/Fidelity

Does the output of an explanation method reflect the underlying *'computation or behavior'* of the black-box model?

# Sanity Check for Faithfulness/Fidelity

- **Sensitivity to Model Parameters:** if the parameter settings change, the explanations should change.

# Sanity Check for Faithfulness/Fidelity

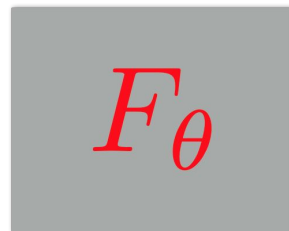
- **Sensitivity to Model Parameters:** if the parameter settings change, the explanations should change.



Parameter Setting 1

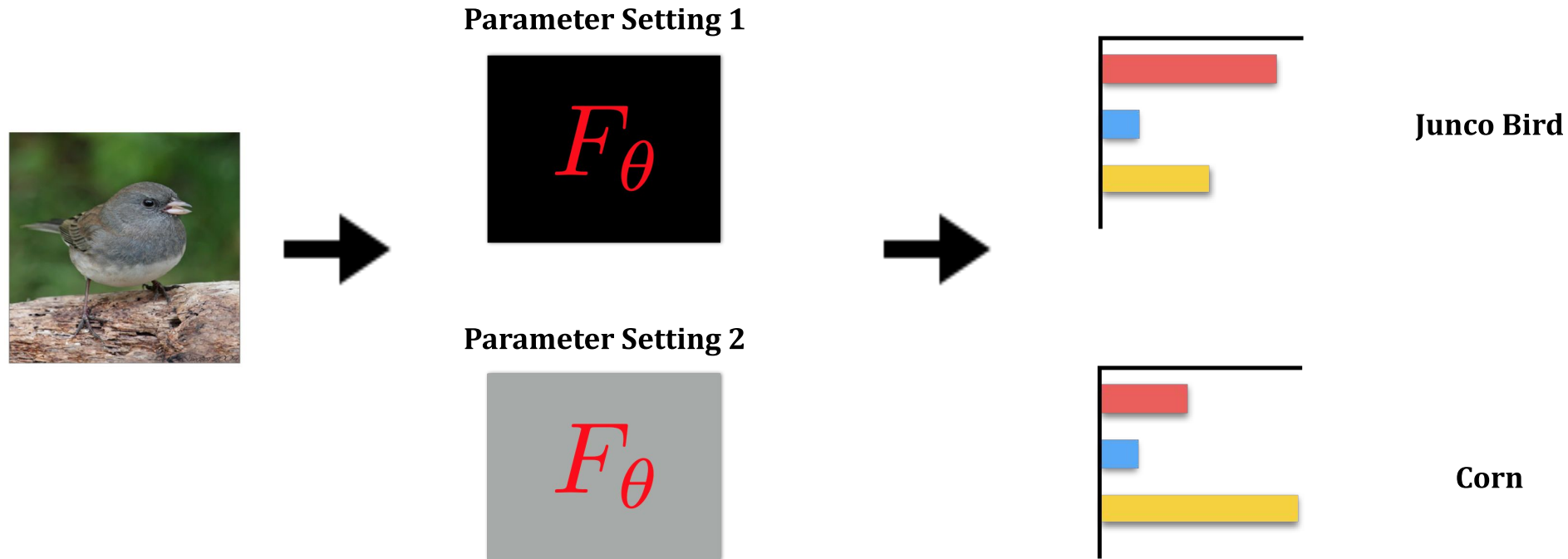


Parameter Setting 2



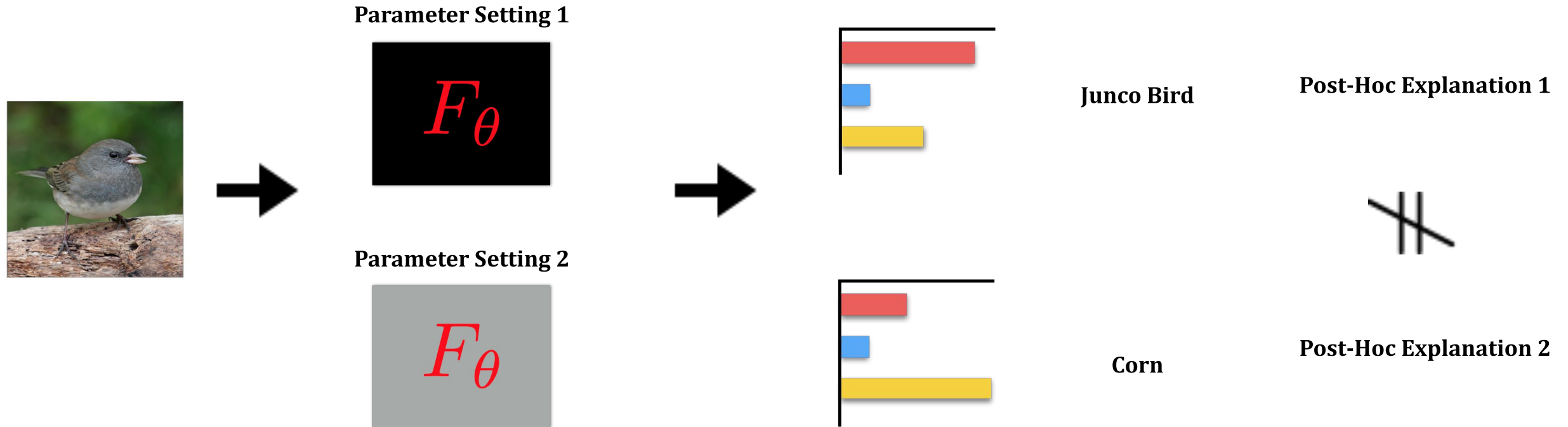
# Sanity Check for Faithfulness/Fidelity

- **Sensitivity to Model Parameters:** if the parameter settings change, the explanations should change.



# Sanity Check for Faithfulness/Fidelity

- **Sensitivity to Model Parameters:** if the parameter settings change, the explanations should change.



# Cascading Randomization Inception-V3

- **Randomize (re-initialize)** model parameters starting from top layer all the way to the input.



**Guided BackProp Explanation Inception-V3 ImageNet**

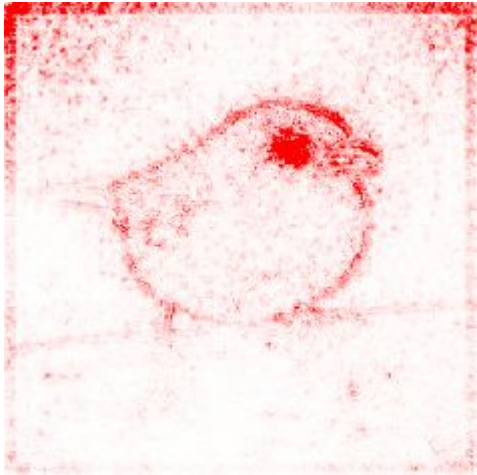
-----

# Cascading Randomization Inception-V3

- **Randomize (re-initialize)** model parameters starting from top layer all the way to the input.



**Normal Model  
Explanation**



**Guided BackProp Explanation Inception-V3 ImageNet**

-----



# Cascading Randomization Inception-V3

- **Randomize (re-initialize)** model parameters starting from top layer all the way to the input.

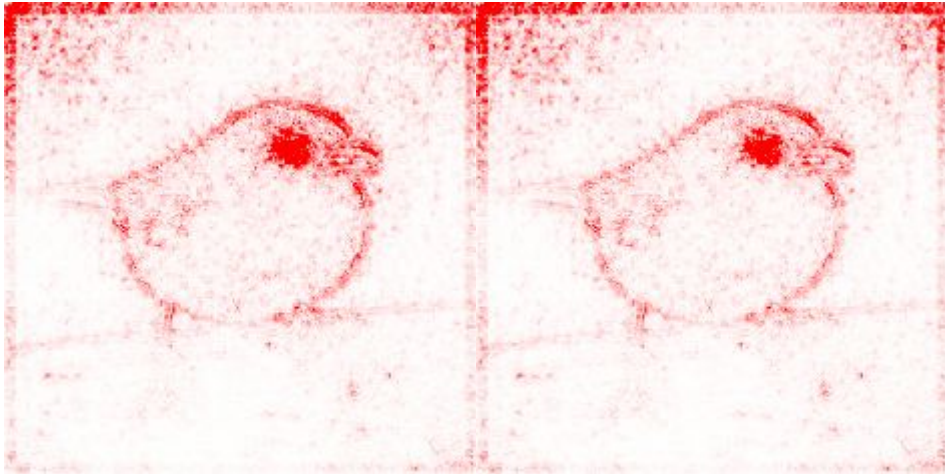


**Guided BackProp Explanation Inception-V3 ImageNet**

-----

**Normal Model  
Explanation**

**Top Layer  
Randomized**



# Cascading Randomization Inception-V3

- **Randomize (re-initialize)** model parameters starting from top layer all the way to the input.

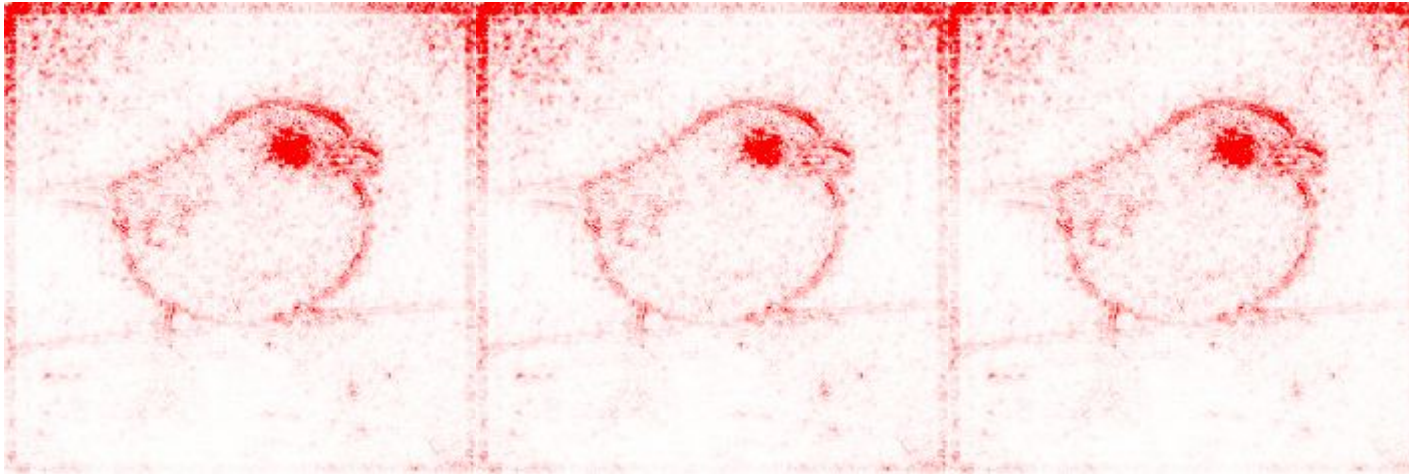


**Guided BackProp Explanation Inception-V3 ImageNet**

-----

**Normal Model  
Explanation**

**Top Layer  
Randomized**



# Cascading Randomization Inception-V3

- **Randomize (re-initialize)** model parameters starting from top layer all the way to the input.

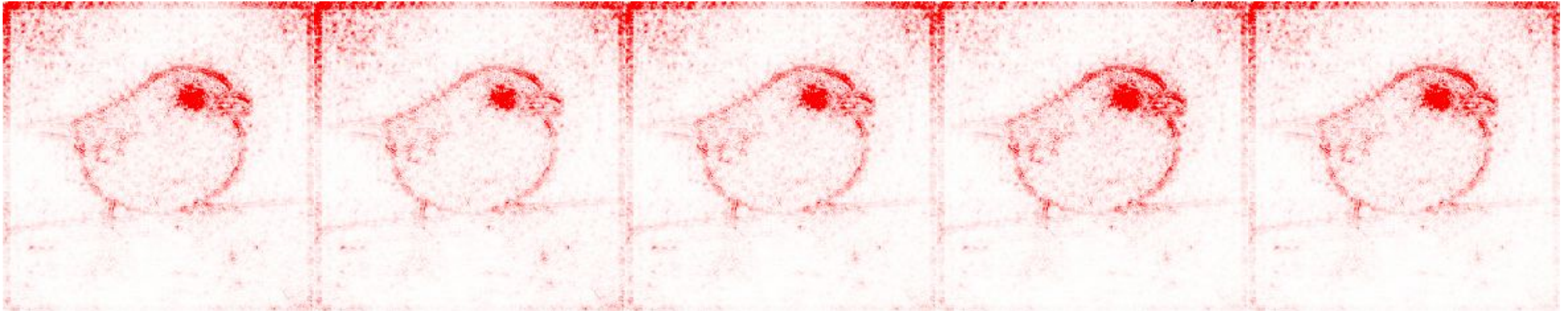


**Guided BackProp Explanation Inception-V3 ImageNet**

-----

**Normal Model  
Explanation**

**Top Layer  
Randomized**



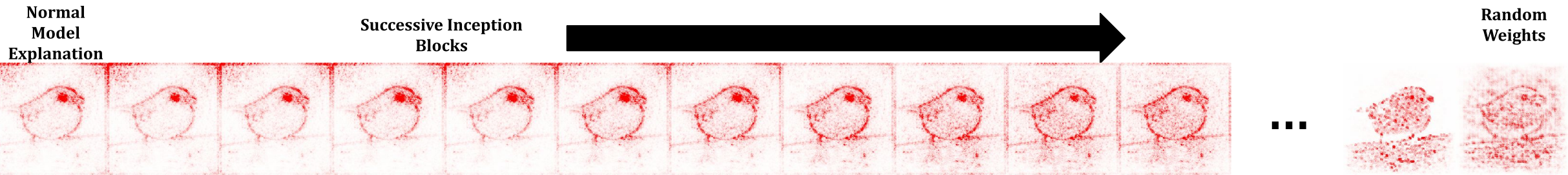
# Cascading Randomization Inception-V3

- **Randomize (re-initialize)** model parameters starting from top layer all the way to the input.



**Guided BackProp Explanation Inception-V3 ImageNet**

-----

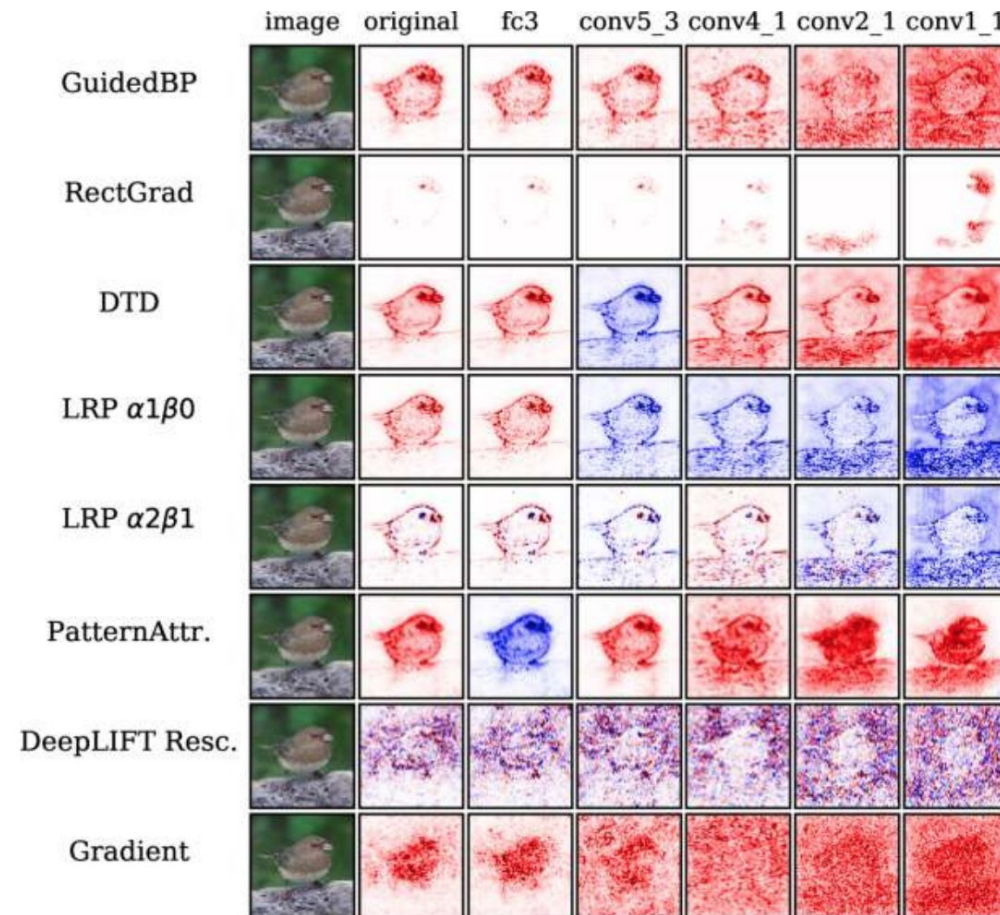


**Guided BackProp is invariant to the higher level weights.**



# ‘Modified backprop approaches’ are invariant

Method that compute relevance via modified backpropagation and performance positive aggregation along the way are invariant to higher layers.



# Source of Invariance

- **Guided BackProp and DeConvNet seek to approximately reconstruct the input** ([Nie et. al. 2018](#)).
- **These modified backprop methods converge to a rank-1 matrix!**  
This is because the product of a sequence of non-negative matrices (non-orthogonal columns, along with other assumptions) converges to a rank-1 matrix ([Theorem 1 in Sixt et. al. 2020](#)).

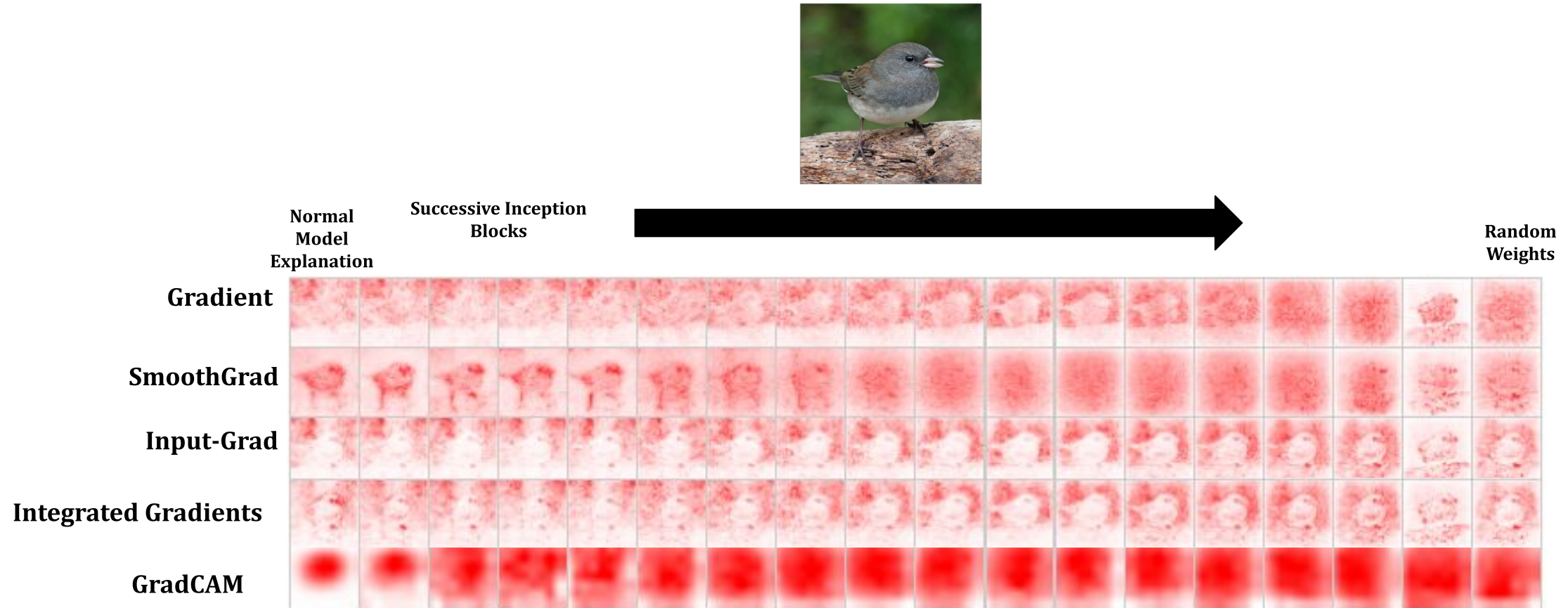
# Source of Invariance

- **Guided BackProp and DeConvNet seek to approximately reconstruct the input** ([Nie et. al. 2018](#)).
- **These modified backprop methods converge to a rank-1 matrix!**  
This is because the product of a sequence of non-negative matrices (non-orthogonal columns, along with other assumptions) converges to a rank-1 matrix ([Theorem 1 in Sixt et. al. 2020](#)).

-----

- |                   |   |
|-------------------|---|
| ● DeConvNet       | ● Deep Taylor Decomposition                         |
| ● Guided BackProp | ● Pattern Net and Pattern Attribution (empirically) |
| ● Guided GradCAM  | ● RectGrad  |

# Cascading Randomization Inception-V3





# Limitations

## ~~● Faithfulness/Fidelity~~

- ~~■ Some explanation methods do not *'reflect'* the underlying model.~~

## ● Fragility

- Post-hoc explanations can be easily manipulated.

# Post-hoc Explanations are Fragile

Post-hoc explanations can be easily manipulated.

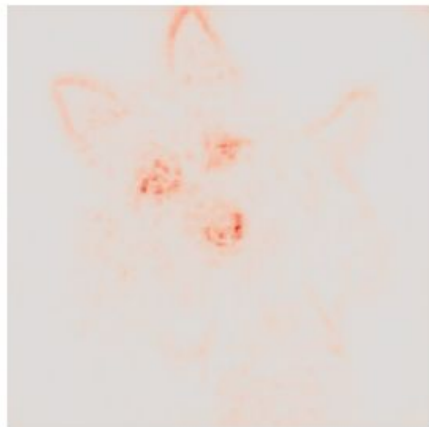
Original Image



# Post-hoc Explanations are Fragile

Post-hoc explanations can be easily manipulated.

Original Image



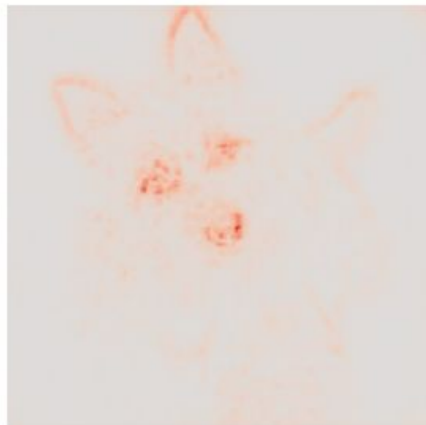
# Post-hoc Explanations are Fragile

Post-hoc explanations can be easily manipulated.

Original Image



Manipulated Image



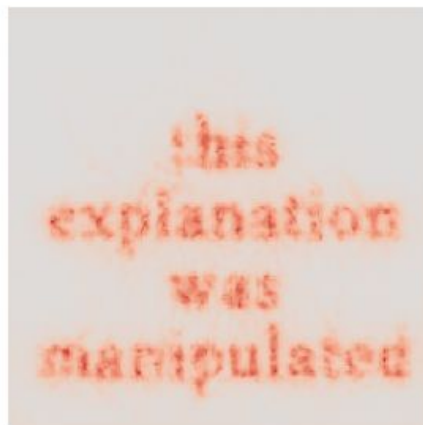
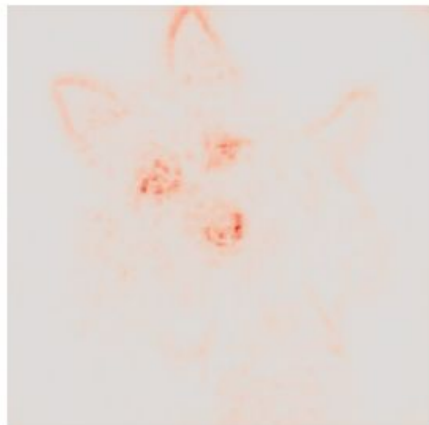
# Post-hoc Explanations are Fragile

Post-hoc explanations can be easily manipulated.

Original Image

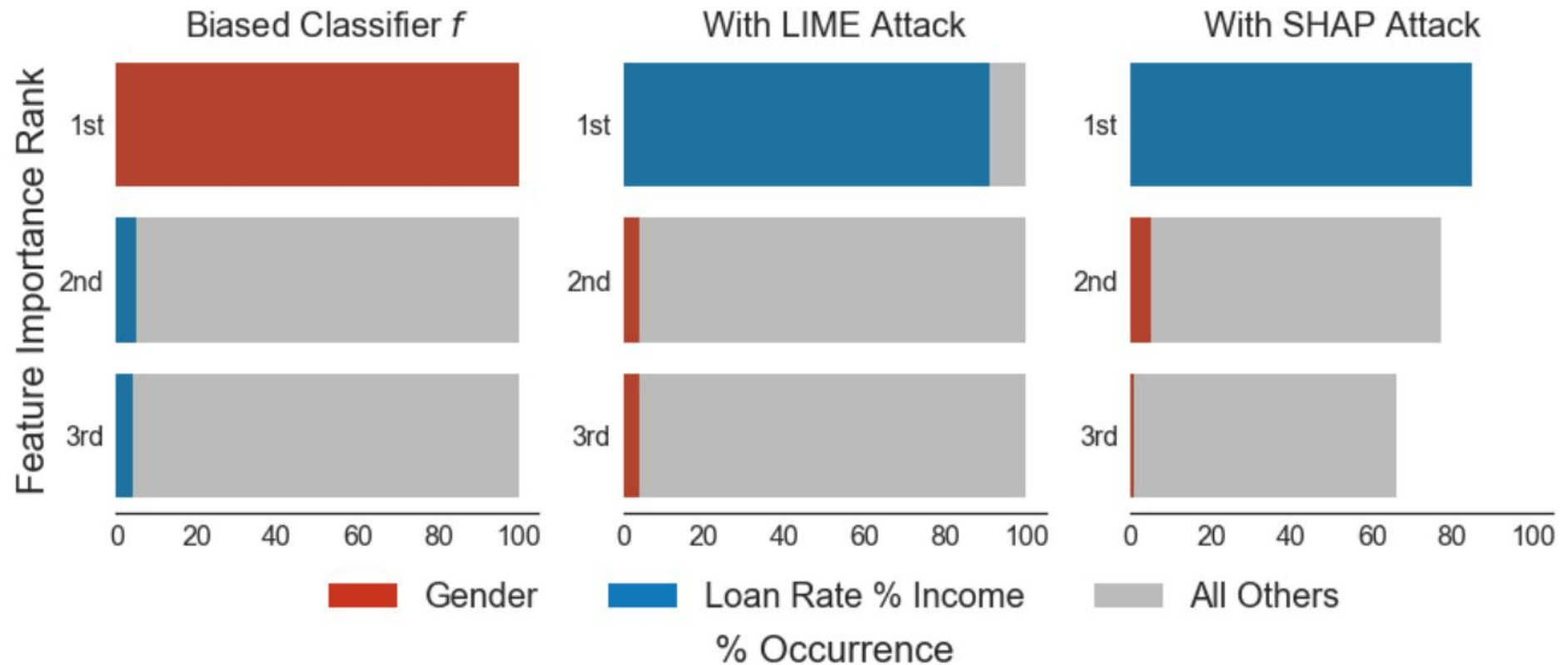


Manipulated Image



# Scaffolding Attack on LIME & SHAP

Scaffolding attack used to **hide classifier dependence on gender**.



# Adversarial Attack on Explanations

Minimally modify the input with a **small perturbation without changing the model prediction.**

$$\arg \max_{\delta} \mathcal{D} (I(\mathbf{x}_t; \mathcal{N}), I(\mathbf{x}_t + \delta; \mathcal{N}))$$

# Adversarial Attack on Explanations

Minimally modify the input with a **small perturbation without changing the model prediction.**

$$\arg \max_{\boldsymbol{\delta}} \mathcal{D}(\mathbf{I}(\mathbf{x}_t; \mathcal{N}), \mathbf{I}(\mathbf{x}_t + \boldsymbol{\delta}; \mathcal{N}))$$

subject to:  $\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon,$



# Adversarial Attack on Explanations

Minimally modify the input with a **small perturbation without changing the model prediction.**

$$\arg \max_{\boldsymbol{\delta}} \mathcal{D}(\mathbf{I}(\mathbf{x}_t; \mathcal{N}), \mathbf{I}(\mathbf{x}_t + \boldsymbol{\delta}; \mathcal{N}))$$

$$\text{subject to: } \|\boldsymbol{\delta}\|_{\infty} \leq \epsilon,$$

$$\text{Prediction}(\mathbf{x}_t + \boldsymbol{\delta}; \mathcal{N}) = \text{Prediction}(\mathbf{x}_t; \mathcal{N})$$

# Other Attacks

- Shift attack by [Kindermans & Hooker et. al. \(2017\)](#).
- Augmented loss function attack by [Dombrowski et. al. \(2019\)](#).
- Passive and Active fooling loss augmentation attack by [Heo et. al. \(2019\)](#).

# Other Attacks

- Shift attack by [Kindermans & Hooker et. al. \(2017\)](#).
- Augmented loss function attack by [Dombrowski et. al. \(2019\)](#).
- Passive and Active fooling loss augmentation attack by [Heo et. al. \(2019\)](#).

## Methods Affected

-----

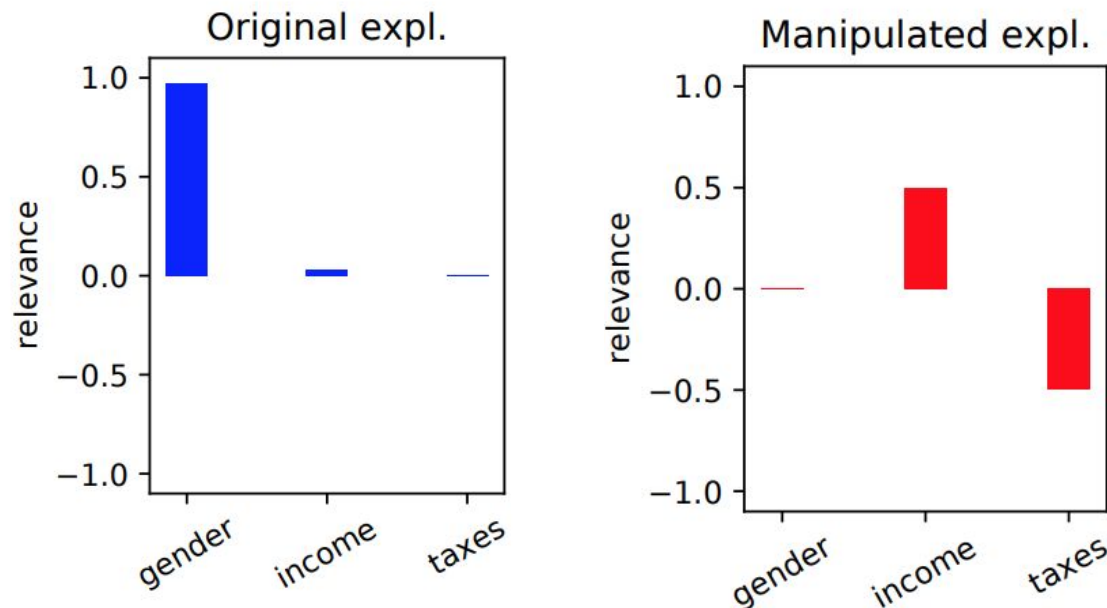
- |                   |                             |
|-------------------|-----------------------------|
| ● LIME            | ● SHAP                      |
| ● Gradient        | ● Integrated Gradients      |
| ● Input-Gradient  | ● LRP                       |
| ● DeConvNet       | ● Deep Taylor Decomposition |
| ● Guided BackProp | ● Pattern Attribution       |
| ● GradCAM         | ● Training Point Ranking    |

# Defense Against Manipulation

Anders et. al. (2020) propose: 1) Hyperplane method & 2) Autoencoder to defend explanations against manipulation.

## Credit Scoring Example

-----

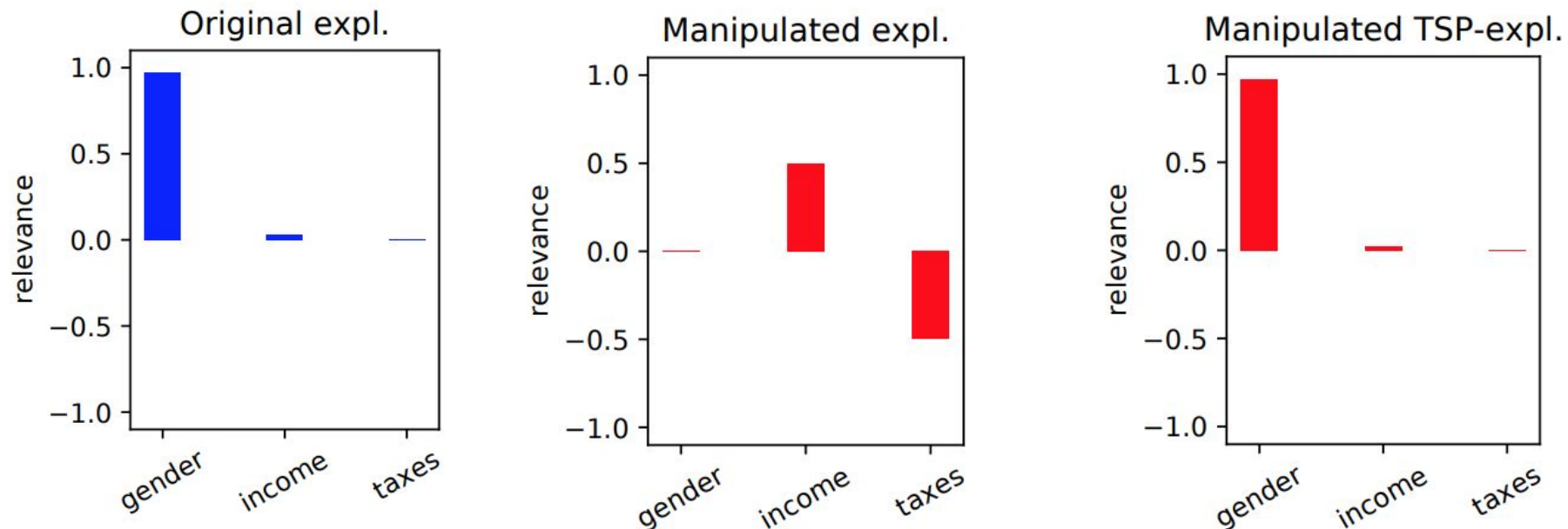


# Defense Against Manipulation

Anders et. al. (2020) propose: 1) Hyperplane method & 2) Autoencoder to defend explanations against manipulation.

## Credit Scoring Example

-----



# Limitations

## ~~● Faithfulness/Fidelity~~

- ~~■ Some explanations do not reflect the underlying model.~~

## ~~● Fragility~~

- ~~■ Post-hoc explanations can be easily manipulated.~~

## ● Stability

- Slight changes to inputs can cause large changes in explanations.

# Limitations: Stability

Post-hoc explanations can be unstable to small, **non-adversarial**, perturbations to the input.

# Limitations: Stability

Post-hoc explanations can be unstable to small, **non-adversarial**, perturbations to the input.

## ‘Local Lipschitz Constant’

$$\hat{L}(x_i) = \operatorname{argmax}_{x_j \in B_\epsilon(x_i)} \frac{\|f(x_i) - f(x_j)\|_2}{\|x_i - x_j\|_2}$$

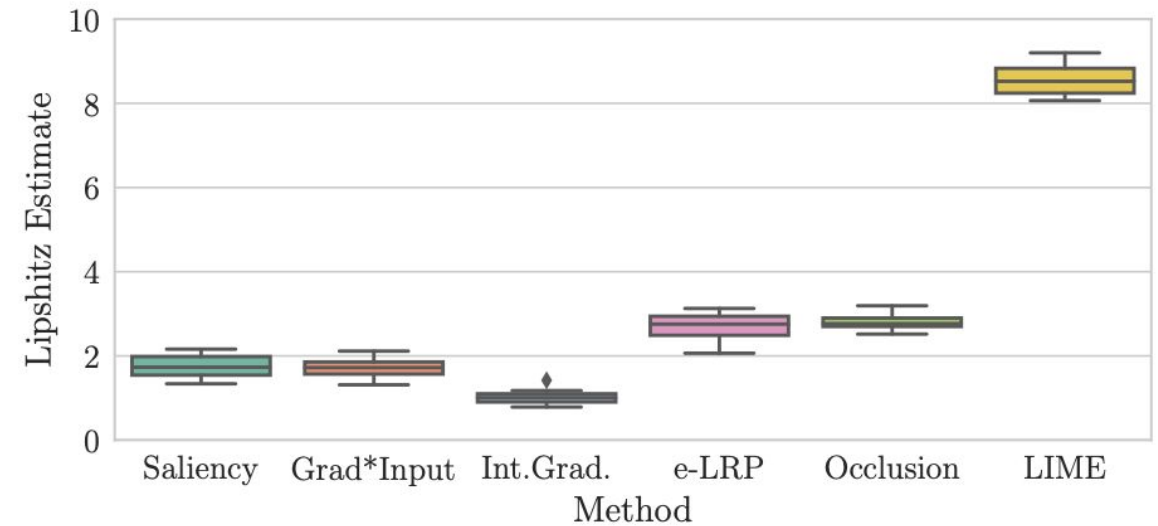
Diagram illustrating the calculation of the Local Lipschitz Constant ( $\hat{L}(x_i)$ ):

- An upward arrow labeled "Input" points to  $x_i$  in the denominator of the formula.
- A downward arrow labeled "Explanation function: LIME, SHAP, Gradient...etc." points from the function  $f$  in the numerator.



# Limitations: Stability

- Perturbation approaches like LIME can be unstable.
- [Yeh et. al. \(2019\)](#) analytically derive bounds on explanations sensitive for certain popular methods and propose stable variants.

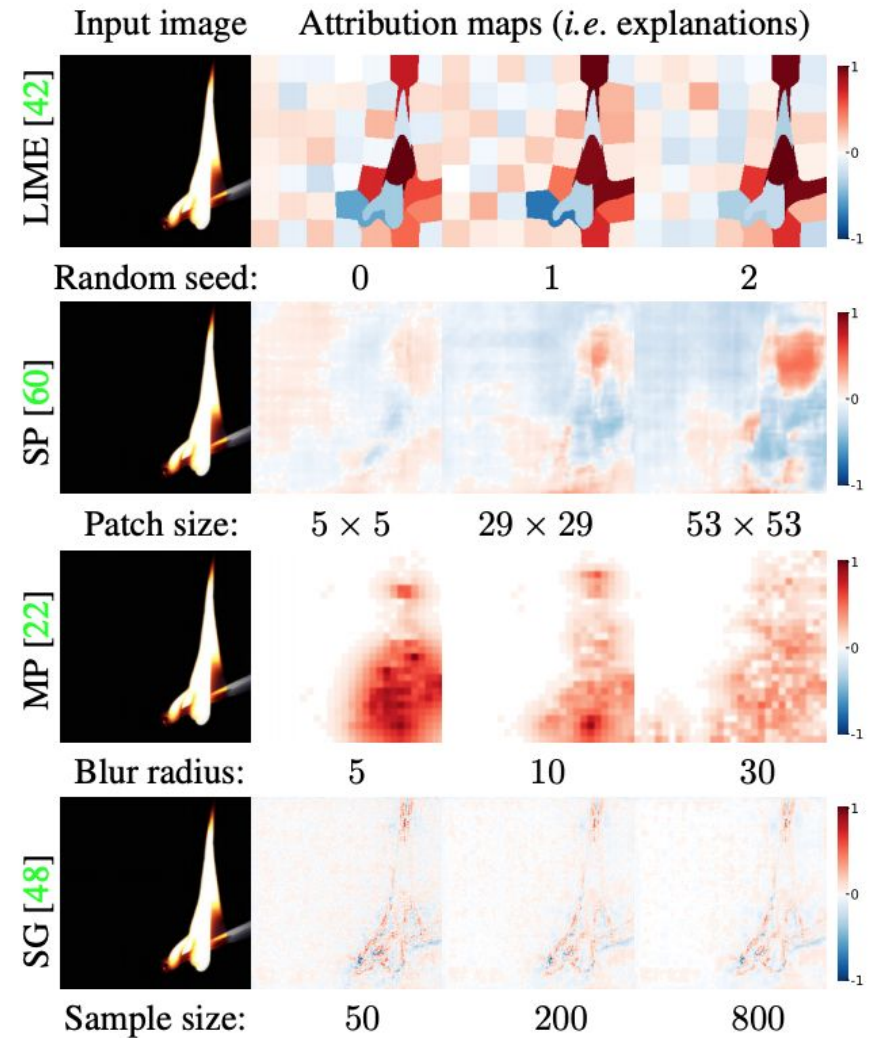


Estimate for 100 tests for an MNIST Model.

[Alvarez et. al. 2018.](#)

# Sensitivity to Hyperparameters

Explanations can be highly sensitive to hyperparameters such as **random seed**, number of perturbations, patch size, etc.



# Limitations

## ● ~~Faithfulness/Fidelity~~

- ~~Some explanations do not reflect the underlying model.~~

## ● ~~Fragility~~

- ~~Post-hoc explanations can be easily manipulated.~~

## ● ~~Stability~~

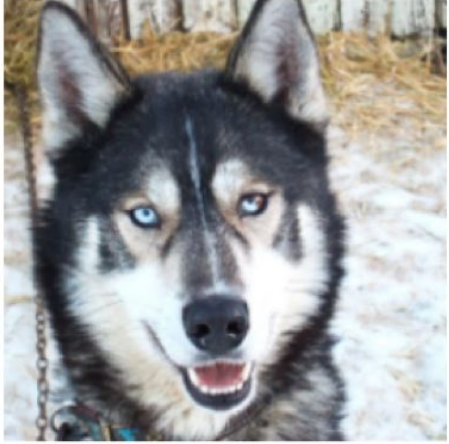
- ~~Slight changes to inputs can cause large changes in explanations.~~

## ● Useful in practice?

- Unclear if a data scientist (ML engineer)/lay person use explanations to isolate errors, improve 'trust', and 'simulatability' in practice?

# Model Debugging: Spurious Signals

**True Label:** Siberian Husky



**Model**



**LIME**



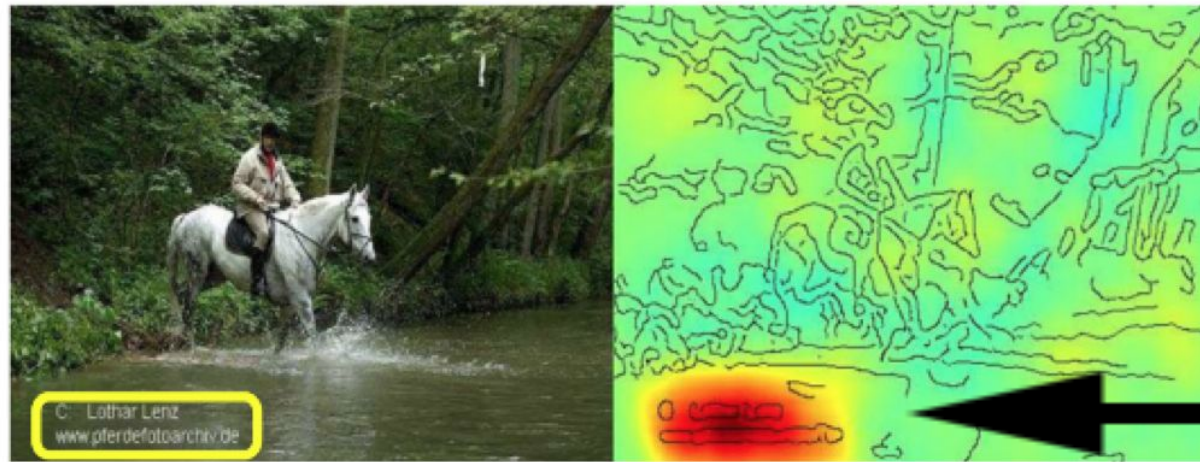
**Predictions**



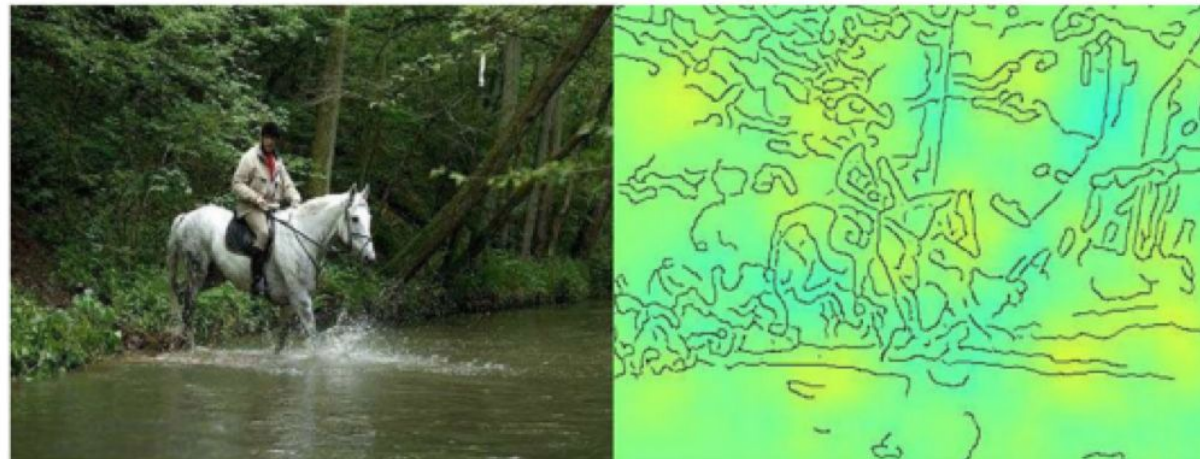
**Relying on snow background**

# Model Debugging: Spurious Signals

Horse-picture from Pascal VOC data set



**Relying on Image Captions to find horses.**



# Explanations as Priors & Model ‘Simulatability’

- Regularizing explanations during training:
  - reduces reliance on **spurious training signals** ([Ross et. al., 2017](#); [Reiger et. al., 2020](#); & [Erion et. al. 2020](#));
  - improves **robustness to adversarial examples** ([Ross et. al., 2018](#)).



# Explanations as Priors & Model ‘Simulatability’

- Regularizing explanations during training:
  - reduces reliance on **spurious training signals** ([Ross et. al., 2017](#); [Reiger et. al., 2020](#); & [Erion et. al. 2020](#));
  - improves **robustness to adversarial examples** ([Ross et. al., 2018](#)).
- Explanations help improve ability of **end-users to simulate the model**:
  - tabular LIME improves forward and counterfactual simulatability ([Hase et. al. 2020](#));
  - prototype explanation improves counterfactual simulatability ([Hase et. al. 2020](#)).

# Explanations with perfect fidelity can still mislead

In a bail adjudication task, **misleading** high-fidelity explanations improve end-user (domain experts) trust.

## True Classifier relies on race

If **Race**  $\neq$  African American:

If **Prior-Felony** = Yes and **Crime-Status** = Active, then **Risky**

If **Prior-Convictions** = 0, then **Not Risky**

If **Race** = African American:

If **Pays-rent** = No and **Gender** = Male, then **Risky**

If **Lives-with-Partner** = No and **College** = No, then **Risky**

If **Age**  $\geq 35$  and **Has-Kids** = Yes, then **Not Risky**

If **Wages**  $\geq 70K$ , then **Not Risky**

Default: **Not Risky**



# Explanations with perfect fidelity can still mislead

In a bail adjudication task, **misleading** high-fidelity explanations improve end-user (domain experts) trust.

## True Classifier relies on race

If **Race**  $\neq$  **African American**:  
If **Prior-Felony** = **Yes** and **Crime-Status** = **Active**, then **Risky**  
If **Prior-Convictions** = **0**, then **Not Risky**

If **Race** = **African American**:  
If **Pays-rent** = **No** and **Gender** = **Male**, then **Risky**  
If **Lives-with-Partner** = **No** and **College** = **No**, then **Risky**  
If **Age**  $\geq 35$  and **Has-Kids** = **Yes**, then **Not Risky**  
If **Wages**  $\geq 70K$ , then **Not Risky**

Default: **Not Risky**

## High fidelity 'misleading' explanation

If **Current-Offense** = **Felony**:  
If **Prior-FTA** = **Yes** and **Prior-Arrests**  $\geq 1$ , then **Risky**  
If **Crime-Status** = **Active** and **Owns-House** = **No** and **Has-Kids** = **No**, then **Risky**  
If **Prior-Convictions** = **0** and **College** = **Yes** and **Owns-House** = **Yes**, then **Not Risky**

If **Current-Offense** = **Misdemeanor** and **Prior-Arrests**  $> 1$ :  
If **Prior-Jail-Incarcerations** = **Yes**, then **Risky**  
If **Has-Kids** = **Yes** and **Married** = **Yes** and **Owns-House** = **Yes**, then **Not Risky**  
If **Lives-with-Partner** = **Yes** and **College** = **Yes** and **Pays-Rent** = **Yes**, then **Not Risky**

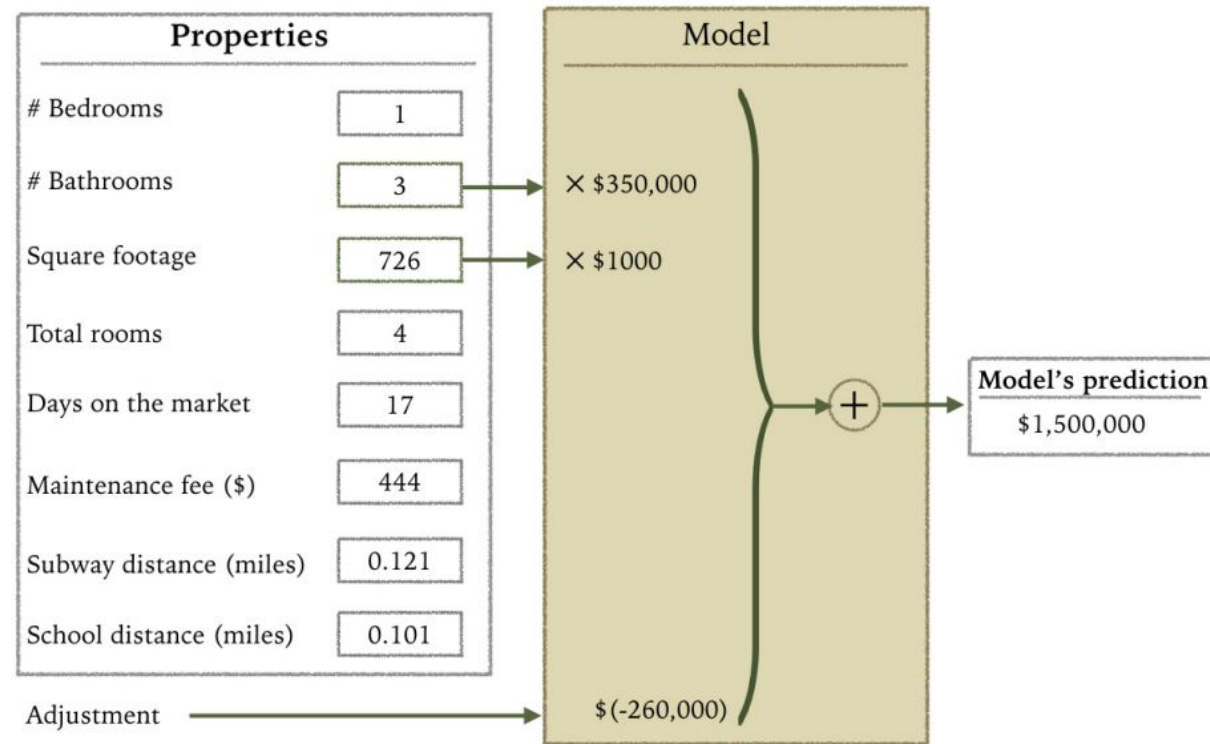
If **Current-Offense** = **Misdemeanor** and **Prior-Arrests**  $\leq 1$ :  
If **Has-Kids** = **No** and **Owns-House** = **No** and **Prior-Jail-Incarcerations** = **Yes**, then **Risky**  
If **Age**  $\geq 50$  and **Has-Kids** = **Yes** and **Prior-FTA** = **No**, then **Not Risky**

Default: **Not Risky**

# Difficulty using explanations for debugging

In a housing price prediction task, Amazon mechanical turkers are unable to use linear model coefficients to diagnose model mistakes.

Attention: This apartment has an unusual combination of # Bedrooms and # Bathrooms.



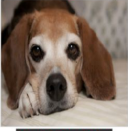
Please take the unusual configuration of this apartment into consideration when making predictions.

# Difficulty using explanations for debugging

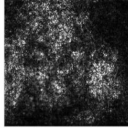
In a dog breeds classification task, users familiar with machine learning **rely on labels, instead of saliency maps**, for diagnosing model errors.

Using the output and explanation of the dog classification model below, do you think this specific model is ready to be sold to customers?


Algorithm Prediction Image



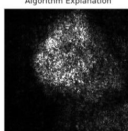
Algorithm Explanation



Algorithm Prediction Image



Algorithm Explanation



DEFINITELY NOT

PROBABLY NOT

UNSURE/MAYBE

PROBABLY

DEFINITELY

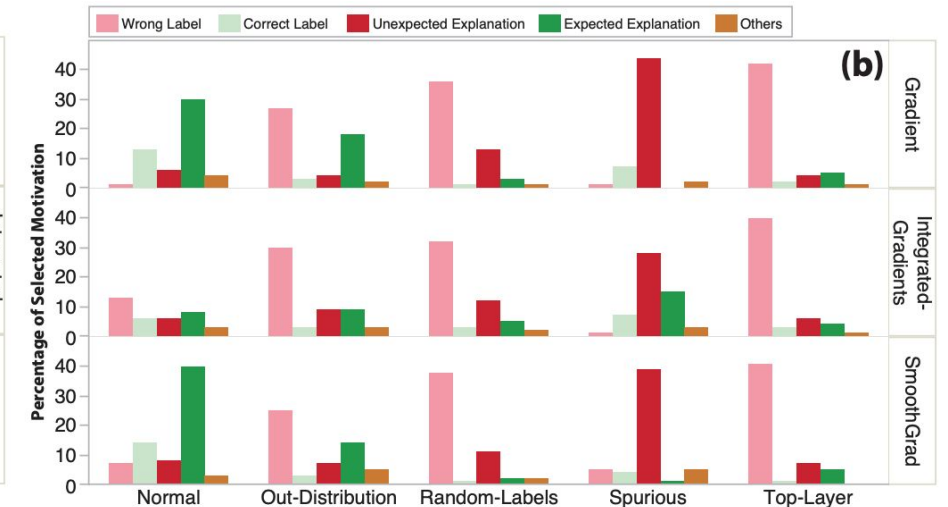
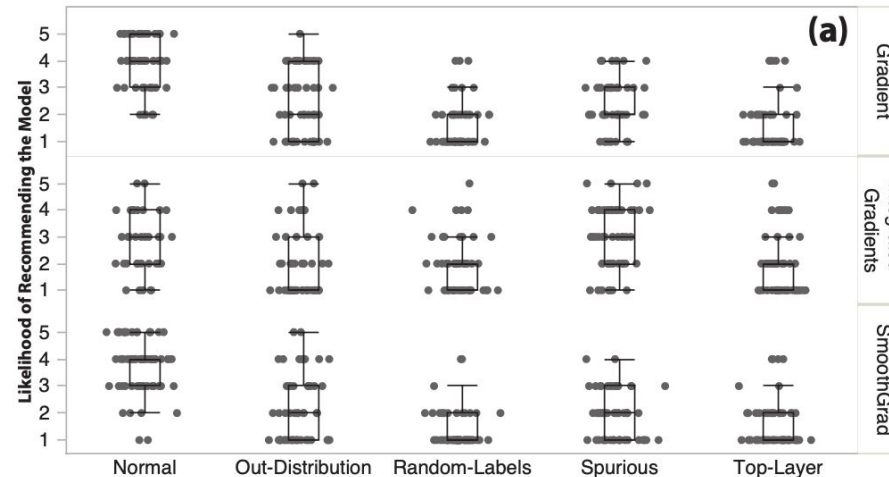
What were your motivation for your response above?

☐ On some or all of the images, the dog breed was wrong.

☐ The dog breeds were correct.

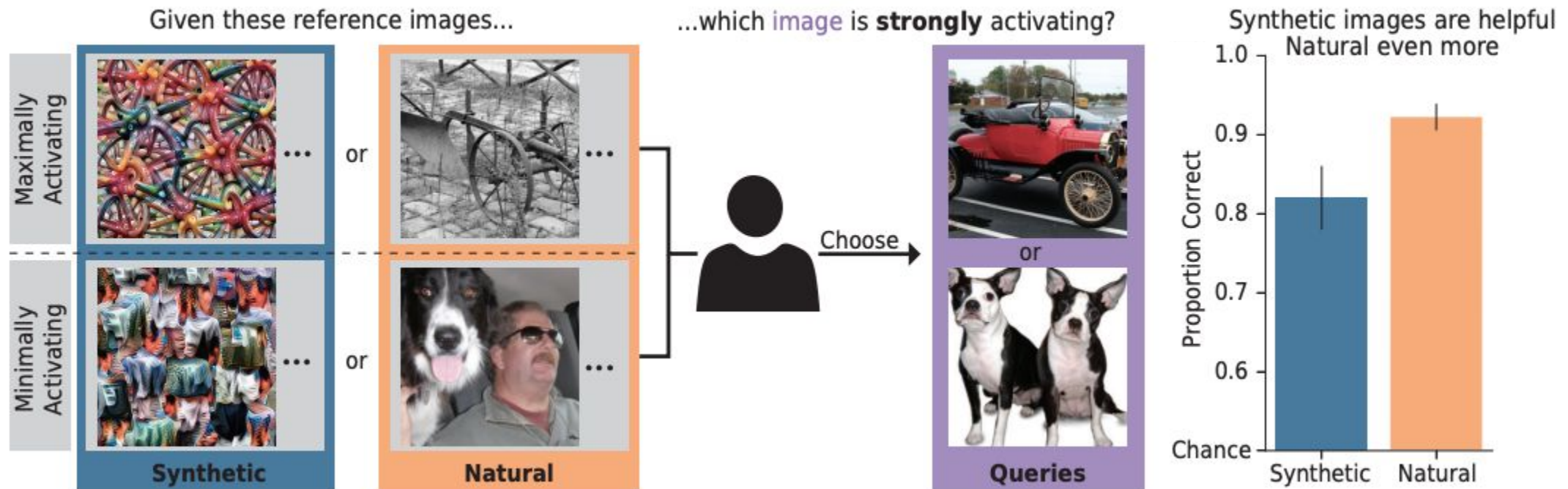
☐ The explanation did not highlight the part of the image that I expected it to focus on.

☐ Other, please specify



# Natural images more helpful than feature visualization

Users found natural images more helpful than feature visualization in deciding whether an image strongly activated a neuron.



# Conflicting Evidence on Utility of Explanations

- **Mixed evidence:**
  - simulation and benchmark studies show that explanations are useful for debugging;
  - however, recent user studies show limited utility in practice.

# Conflicting Evidence on Utility of Explanations

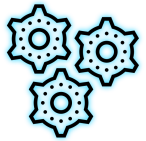
- **Mixed evidence:**
  - simulation and benchmark studies show that explanations are useful for debugging;
  - however, recent user studies show limited utility in practice.
- Rigorous **user studies** and **pilots with end-users** can continue to help provide feedback to researchers on what to address (see: [Alqaraawi et. al. 2020](#), [Bhatt et. al. 2020](#) & [Kaur et. al. 2020](#)).



# Limitations

- **Faithfulness/Fidelity**
  - Some explanation methods do not '*reflect*' the underlying model.
- **Fragility**
  - Post-hoc explanations can be easily manipulated.
- **Stability**
  - Slight changes to inputs can cause large changes in explanations.
- **Useful in practice?**
  - Unclear if a data scientist (ML engineer)/end-user can use explanations to isolate errors, improve 'trust' or simulate the model.

# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations



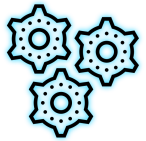
**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability



# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations



**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability

# Future of Post hoc Explainability

Emerging Topics in Explainability Research



# Future of Post hoc Explainability

## Towards Better Post hoc Explanations

Methods for More Reliable  
Post hoc Explanations

Theoretical Analysis of  
Post hoc Explanation Methods

Rigorous Evaluation of the Utility of  
Post hoc Explanations

## Other Emerging Directions

Post hoc Explainability  
Beyond Classification

Intersections with Differential Privacy

Intersections with Fairness

# Future of Post hoc Explainability

## Towards Better Post hoc Explanations



Methods for More Reliable  
Post hoc Explanations

Theoretical Analysis of  
Post hoc Explanation Methods

Rigorous Evaluation of the Utility of  
Post hoc Explanations

## Other Emerging Directions

Post hoc Explainability  
Beyond Classification

Intersections with Differential Privacy

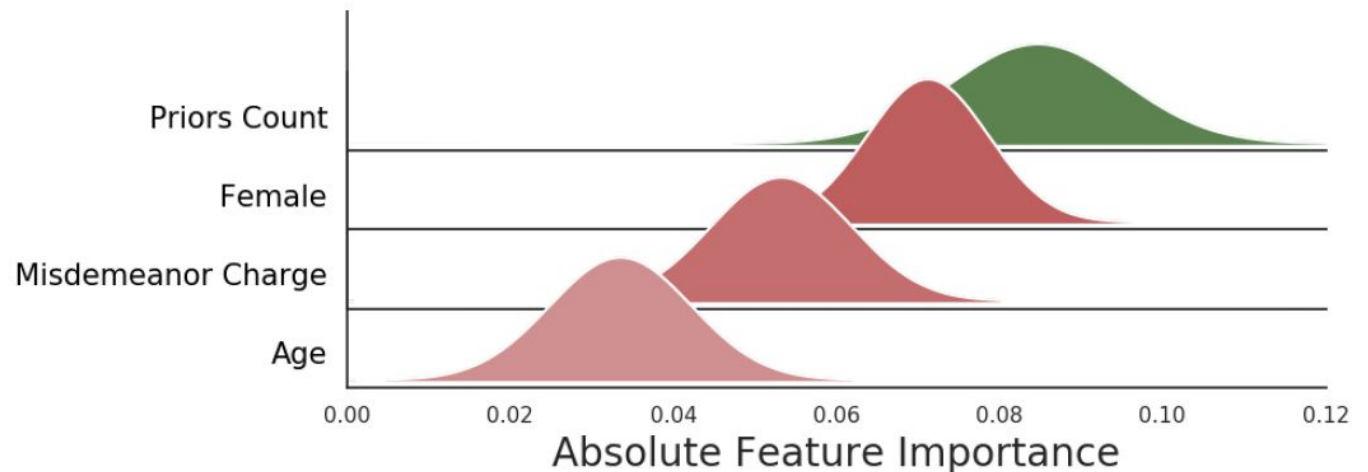
Intersections with Fairness

# Methods for More Reliable Post hoc Explanations

Post hoc explanations have several limitations:  
not faithful to the underlying model, unstable, fragile

- Modeling **uncertainty** in post hoc explanations [Guo et. al. 2018, Slack et. al. 2020]

*Bayesian versions of LIME/SHAP  
with closed form solutions*



# Methods for More Reliable Post hoc Explanations

Post hoc explanations have several limitations:  
not faithful to the underlying model, unstable, fragile

- Generating post hoc explanations that are **stable** as well as **robust to distribution shifts** [Chalasani et. al., 2020, Lakkaraju et. al. 2020]
  - Use adversarial training i.e., minimize the worst case mismatch between explanation and (black box) model predictions.

# Methods for More Reliable Post hoc Explanations

Post hoc explanations have several limitations:  
not faithful to the underlying model, unstable, fragile

- Identifying vulnerabilities in existing post hoc explanation methods and proposing approaches to address these vulnerabilities is a critical research direction going forward!

# Future of Post hoc Explainability

## Towards Better Post hoc Explanations

Methods for More Reliable  
Post hoc Explanations



Theoretical Analysis of  
Post hoc Explanation Methods

Rigorous Evaluation of the Utility of  
Post hoc Explanations

## Other Emerging Directions

Post hoc Explainability  
Beyond Classification

Intersections with Differential Privacy

Intersections with Fairness



# Theoretical Analysis of Post hoc Explanation Methods

- Theoretical analysis of LIME

Theoretical analysis shedding light on the fidelity, stability, and fragility of post hoc explanation methods can be extremely valuable to the progress of the field!

- The coefficients obtained are proportional to the gradient of the function to be explained
- Local error of surrogate model is bounded away from zero with high probability

# Future of Post hoc Explainability

## Towards Better Post hoc Explanations

Methods for More Reliable  
Post hoc Explanations

Theoretical Analysis of  
Post hoc Explanation Methods



Rigorous Evaluation of the Utility of  
Post hoc Explanations

## Other Emerging Directions

Post hoc Explainability  
Beyond Classification

Intersections with Differential Privacy

Intersections with Fairness

# Rigorous Evaluation of the Utility of Post hoc Explanations

- Domain experts and end users seem to be over trusting explanations & the underlying models based on explanations
  - Law school students trusted underlying model 9.8 times more when shown a misleading explanation which “white-washes” the model
  - Data scientists over trusted explanations without even comprehending them -- *“Participants trusted the tools because of their visualizations and their public availability”*

# Responses from Data Scientists Using Explainability Tools (GAM and SHAP)

“I didn’t fully grasp what SHAP values were. This is a pretty popular tool and I get the log-odds concept in general. I figure they were showing SHAP values for a reason. Maybe it’s easier to judge relationships using log-odds instead of predicted value. Anyway, so it made sense I suppose.” (P6, SHAP)

“[The tool] assigns a value that is important to know, but it’s showing that in a way that makes you misinterpret that value. Now I want to go back and check all my answers”... [later] “Okay, so, it’s not showing me a whole lot more than what I can infer on my own. Now I’m thinking... is this an ‘interpretability tool’?” (P4, SHAP)

“Age 38 seems to have the highest positive influence on income based on the plot. Not sure why, but the explanation clearly shows it... makes sense.” (P9, GAMs)

“[The tool] shows visualizations of ML models, which is not something anything else I have worked with has done. It’s very transparent, and that makes me trust it more” (P9, GAMs).

# Are Explanations Helping Humans in Real World Tasks?

- Evaluating the effect of explanations on human-AI collaboration

- Rigorous user studies and evaluations to ascertain the utility of different post hoc explanation methods in various contexts is extremely critical for the progress of the field!

# Future of Post hoc Explainability

## Towards Better Post hoc Explanations

Methods for More Reliable  
Post hoc Explanations

Theoretical Analysis of  
Post hoc Explanation Methods

Rigorous Evaluation of the Utility of  
Post hoc Explanations



## Other Emerging Directions

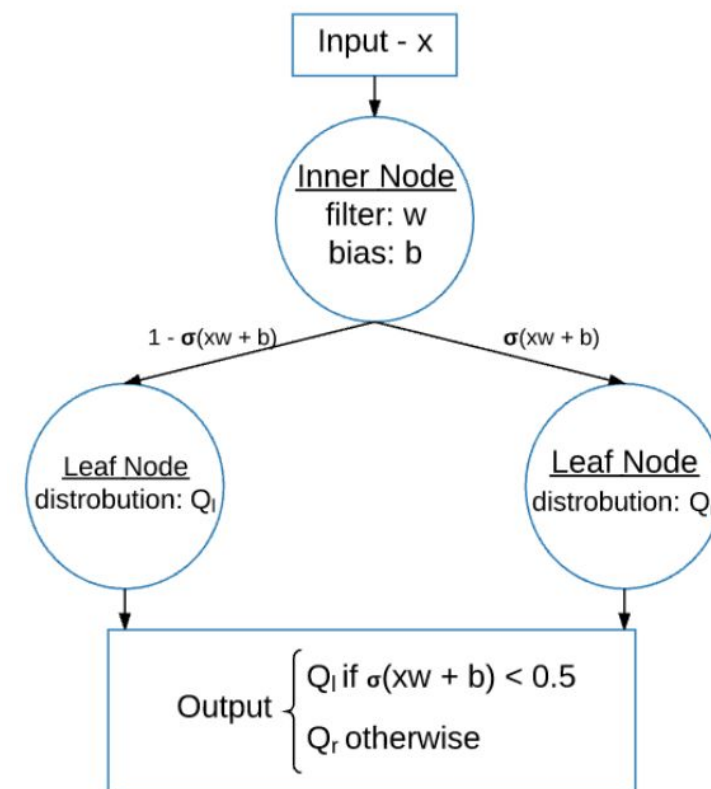
Post hoc Explainability  
Beyond Classification

Intersections with Differential Privacy

Intersections with Fairness

# Beyond Classification: Explainability for RL

- Model distillation using soft decision trees to understand RL policies
  - Map states to actions
- Summarize agent behavior by identifying important states in a policy
  - A state is important if different actions lead to substantially different outcomes



# Beyond Classification: Explainability for RL

- Causal explanations of the behavior of model free RL agents
- Generate explanations of agent behaviour based on counterfactual analysis of the causal model

## Explaining the actions of a StarCraft II agent

*Question* Why not *build\_barracks* ( $A_b$ )?

*Explanation* Because it is more desirable to do action *build\_supply\_depot* ( $A_s$ ) to have more Supply Depots ( $S$ ) as the goal is to have more Destroyed Units ( $D_u$ ) and Destroyed buildings ( $D_b$ ).



# Beyond Classification: Explainability for GNNs

Takes a trained GNN and its predictions and returns an explanation in the form of a set of nodes and edges in the input graph.

Lots of real world applications call for models/algorithms that go beyond classification. Exciting opportunities to explore explainability in these settings!



# Future of Post hoc Explainability

## Towards Better Post hoc Explanations

Methods for More Reliable  
Post hoc Explanations

Theoretical Analysis of  
Post hoc Explanation Methods

Rigorous Evaluation of the Utility of  
Post hoc Explanations



## Other Emerging Directions

Post hoc Explainability  
Beyond Classification

Intersections with Differential Privacy

Intersections with Fairness

# Intersections with Differential Privacy

- 
- 

Need for more theoretical, methodological, and empirical research exploring this intersection!

learning them

# Future of Post hoc Explainability

## Towards Better Post hoc Explanations

Methods for More Reliable  
Post hoc Explanations

Theoretical Analysis of  
Post hoc Explanation Methods

Rigorous Evaluation of the Utility of  
Post hoc Explanations

## Other Emerging Directions

Post hoc Explainability  
Beyond Classification

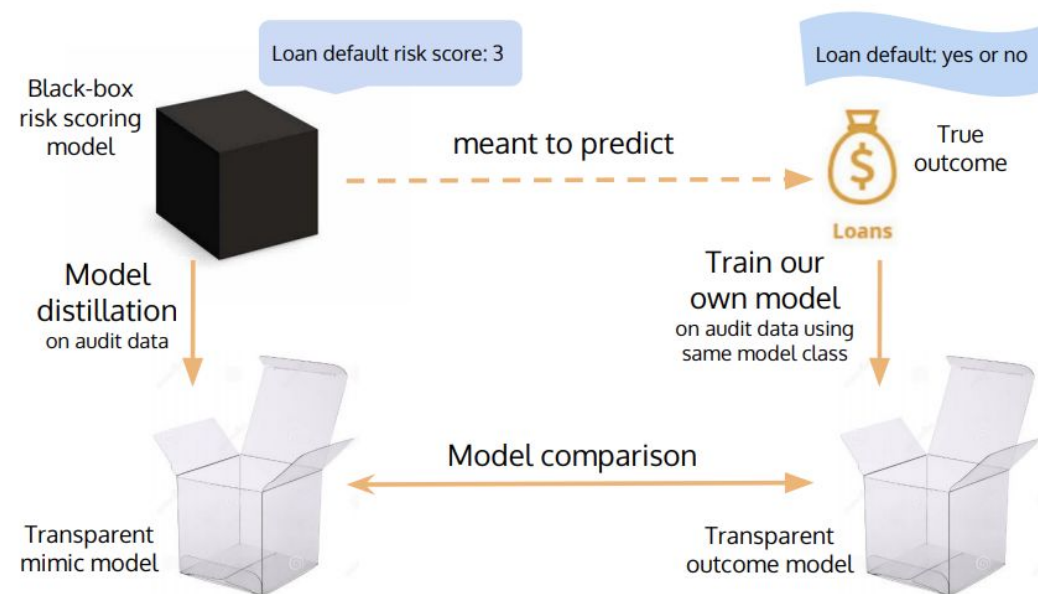
Intersections with Differential Privacy

Intersections with Fairness



# Intersections with Fairness

**Distill and Compare:** Compare the transparent/distilled down versions of risk scoring model and true outcome model to detect biases in risk scoring models.



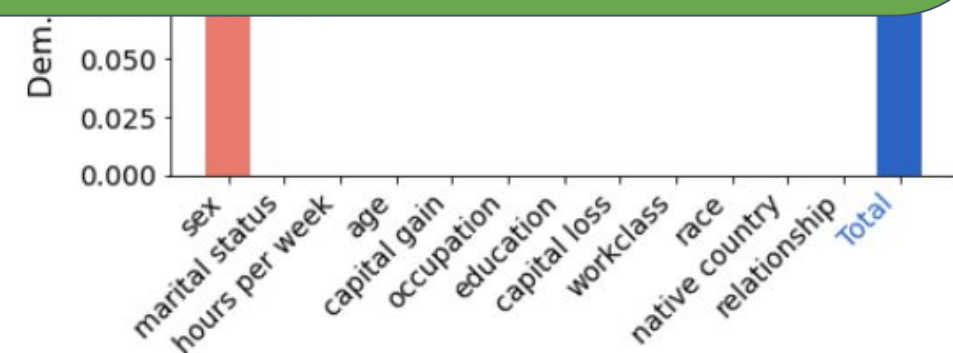
# Intersections with Fairness

- It is **commonly hypothesized** that post hoc explanations can help with **detecting model biases**.
  - Need for more **rigorous theoretical and empirical studies** to quantitatively evaluate this hypothesis
- Can post hoc explanations **help detect unfairness**?
  - How do they complement existing statistical notions of unfairness?

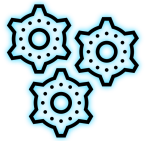
# Intersections with Fairness

The connections between explainability and fairness need to be explored more thoroughly both through rigorous analysis and user studies.

functions which 'explain' the unfairness



# Tutorial on Post hoc Explanations



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations



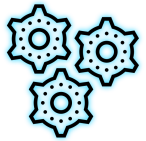
**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability



# Summary of Tutorial



**Approaches** for Post hoc Explainability



**Evaluation** of Explanations



**Limits** of Post hoc Explainability



**Future** of Post hoc Explainability

# Parting Thoughts...

When introducing a new explanation method:

- Who are the **target end users** that the method will help?
- A clear statement about **what capability and/or insight the method aims to provide** to its end users
- **Careful analysis and exposition of the limitations and vulnerabilities** of the proposed method
- **Rigorous user studies** (preferably with actual end users) to evaluate if the method is achieving the desired effect
- Use **quantitative metrics (and not anecdotal evidence)** to make claims about explainability

# Thank You!



**Julius Adebayo**  
MIT



**Hima Lakkaraju**  
Harvard University



**Sameer Singh**  
UC Irvine

**Slides and Video:** [explainml-tutorial.github.io](https://github.com/explainml-tutorial/explainml-tutorial)