



# Interpretability: the myth, questions, and some answers.

Been Kim

Presenting a subset of work  
with a lot of awesome people inside and outside of Google:

Martin Wattenberg, Finale Doshi-Velez, Julius Adebayo, Heinrich Jiang, Maya Gupta, Ike Lage, Andrew Ross, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, Rory Sayres, Ian Goodfellow, Mortiz Hardt, Sam Gershman, Menaka Narayanan, Emily Chen, Jeffrey He

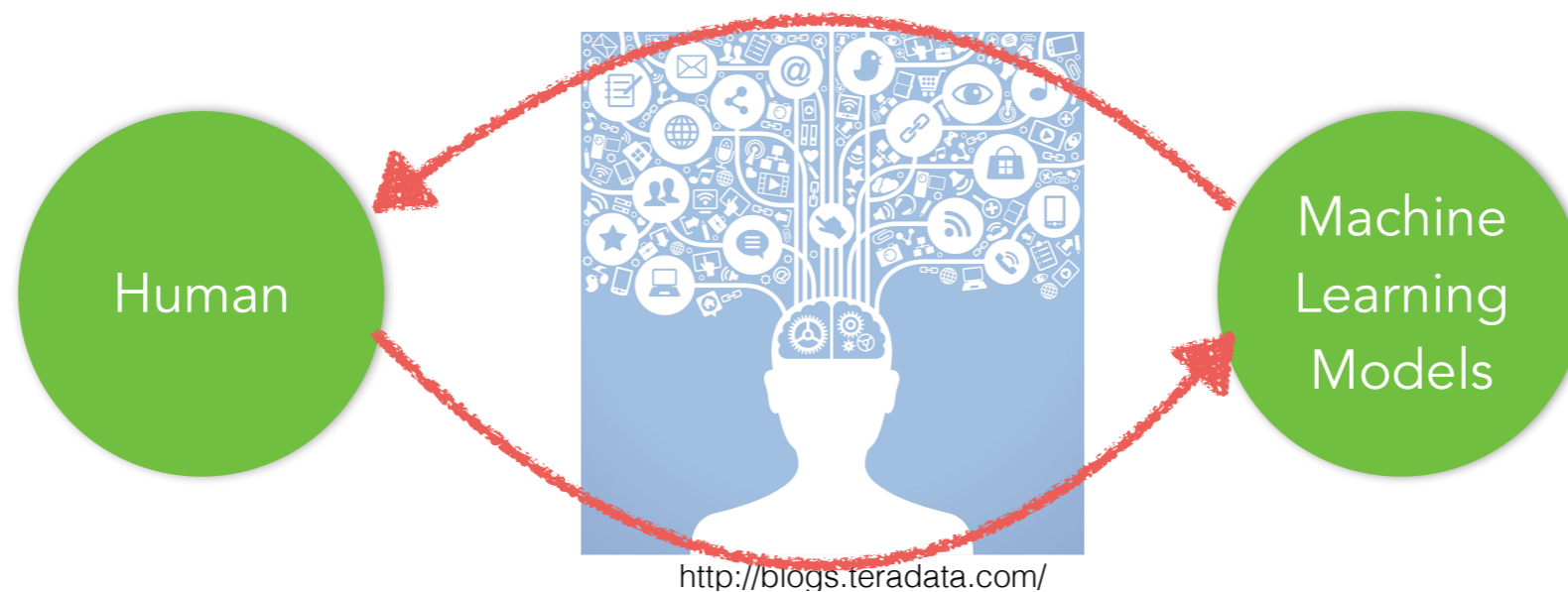
My goal

# interpretability

To use machine learning **responsibly**

we need to ensure that

1. our **values** are aligned
2. our **knowledge** is reflected  
**for everyone.**

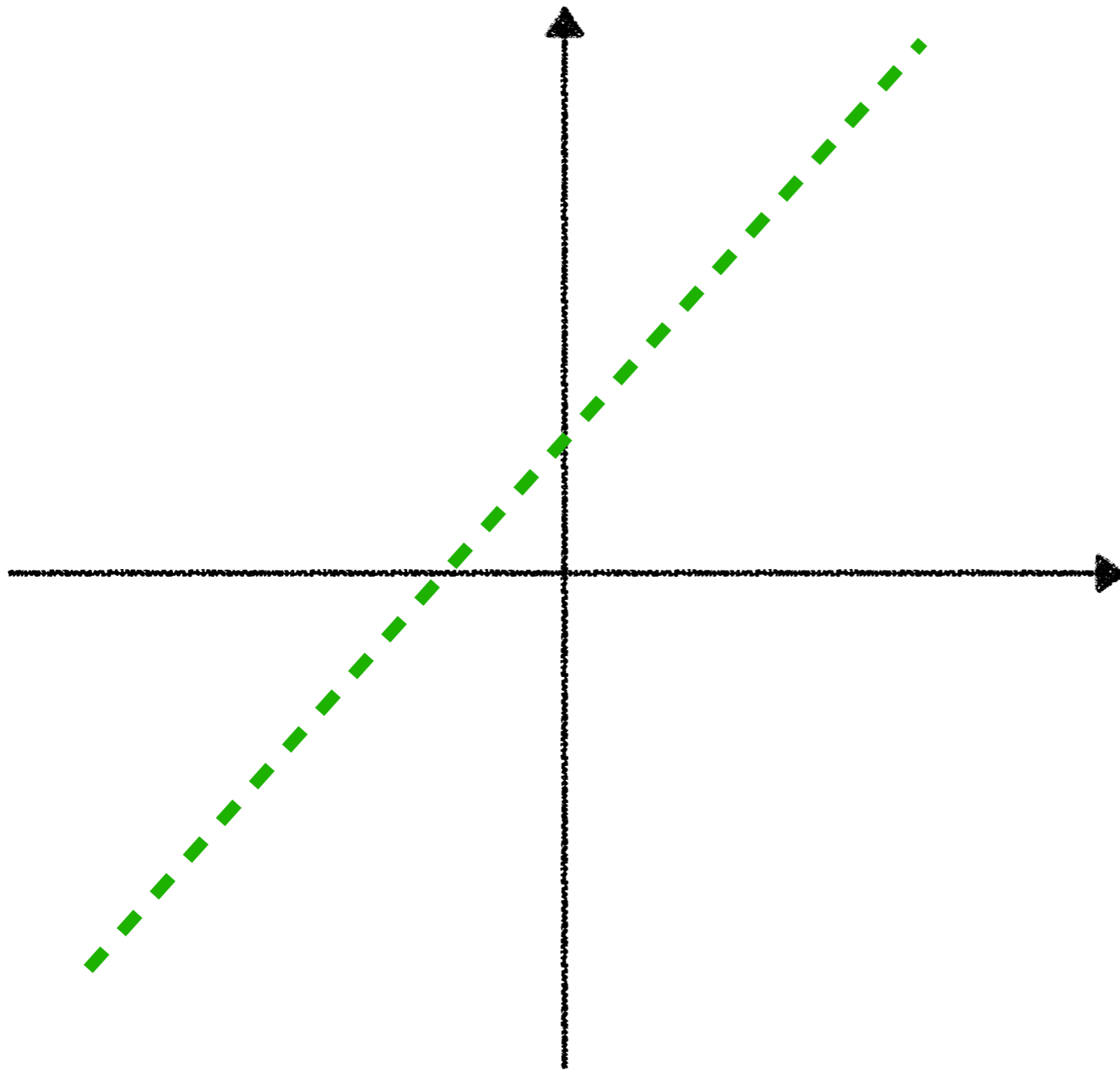


ingredients for interpretability methods.

$$\operatorname{argmax}_E Q(E|?)$$

Some quality function

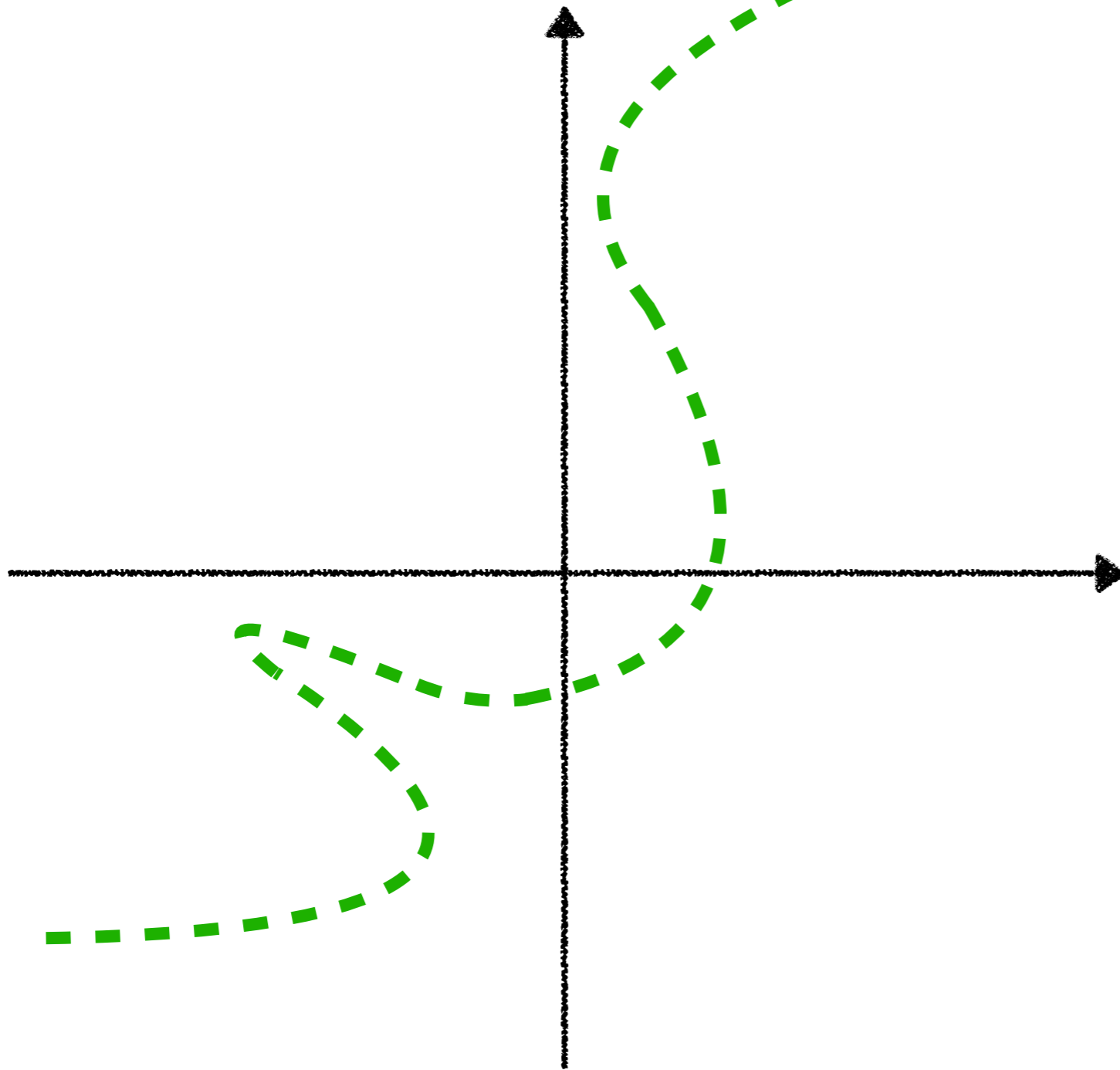
$$\operatorname{argmax}_E Q(E|M, \text{score})$$



**Model**

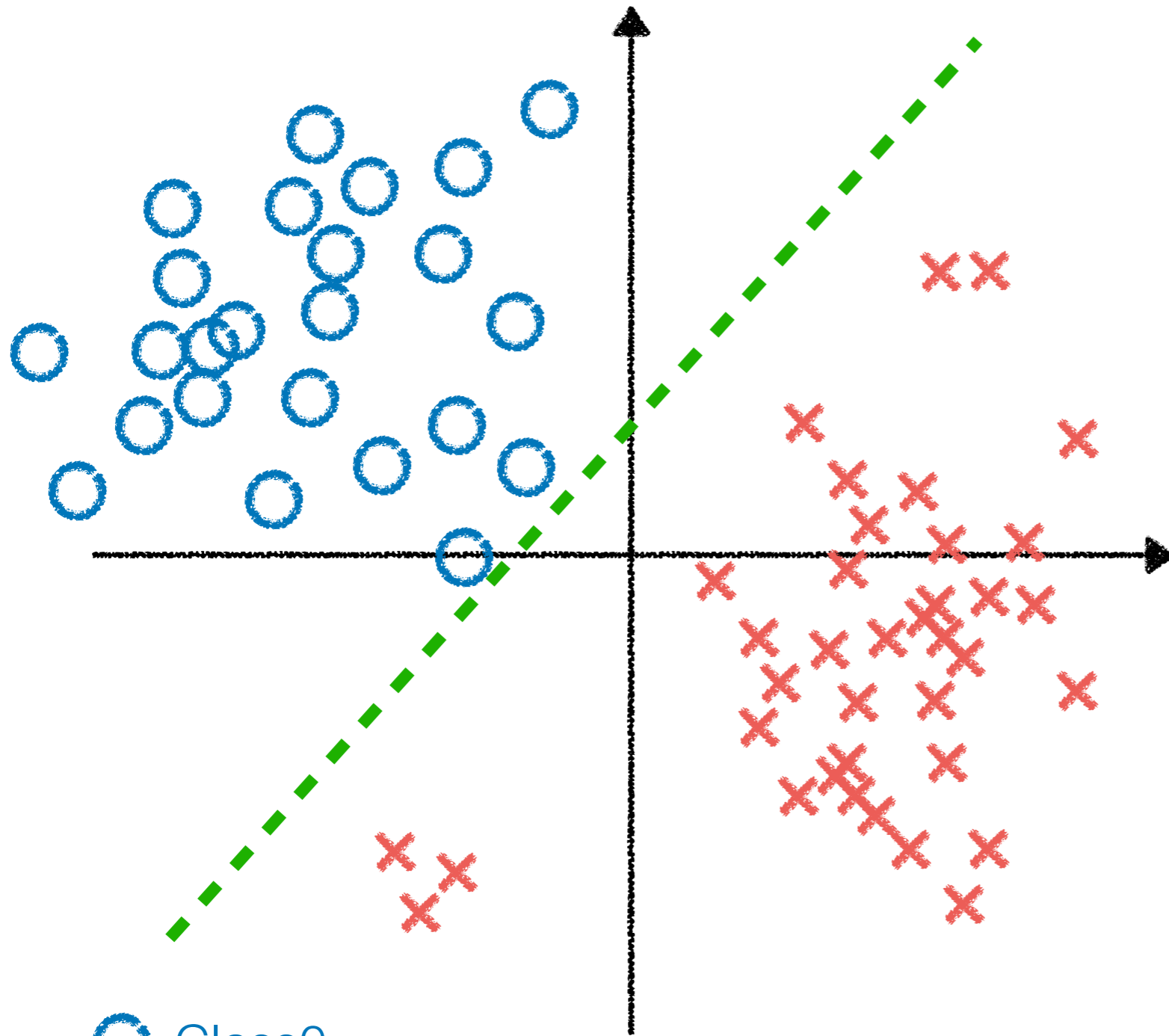


$$\operatorname{argmax}_E Q(E|M, \text{score})$$



Model

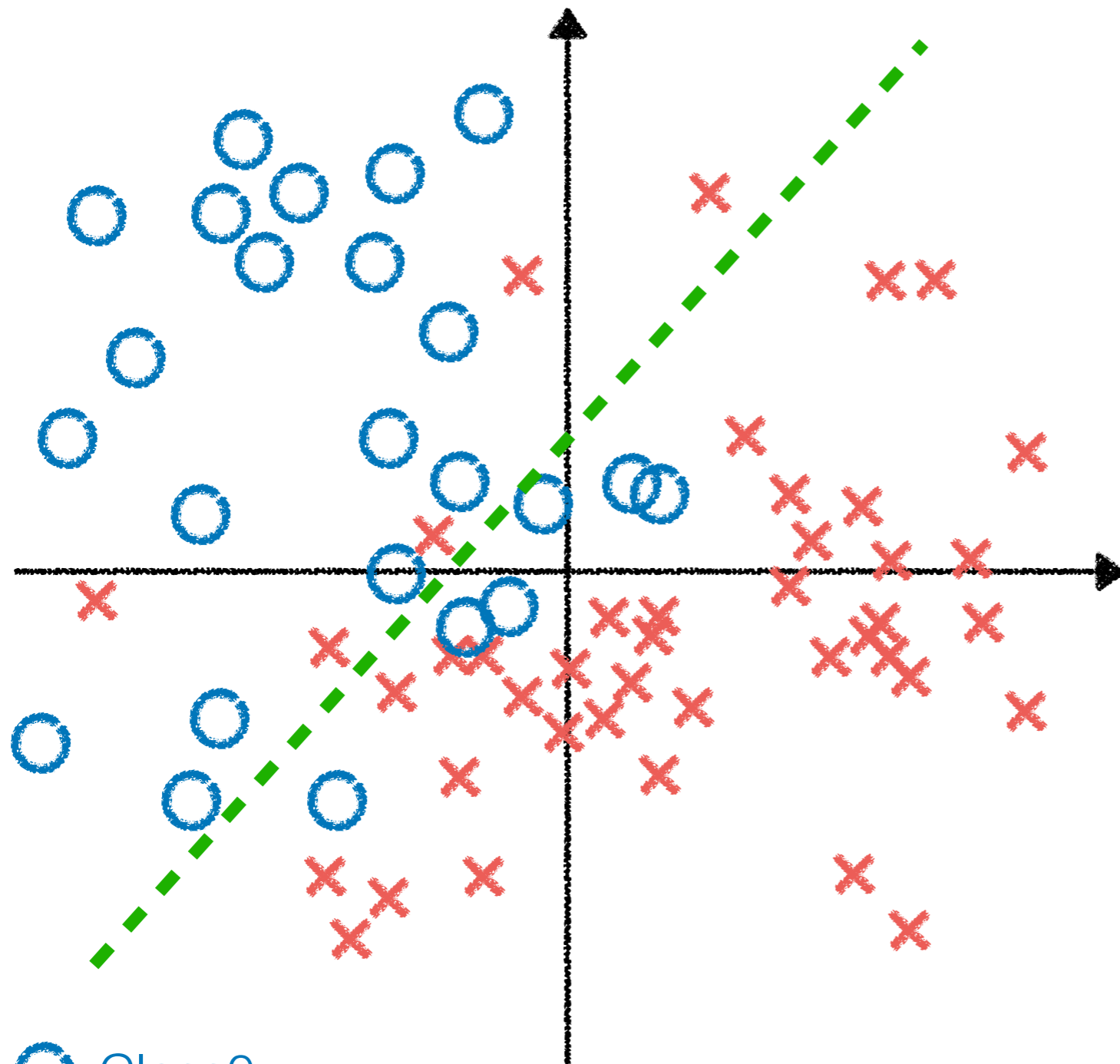
$$\operatorname{argmax}_E Q(E|M, \underline{D}, )$$



○ Class0

× Class1

$$\operatorname{argmax}_E Q(E|M, \underline{D}, )$$

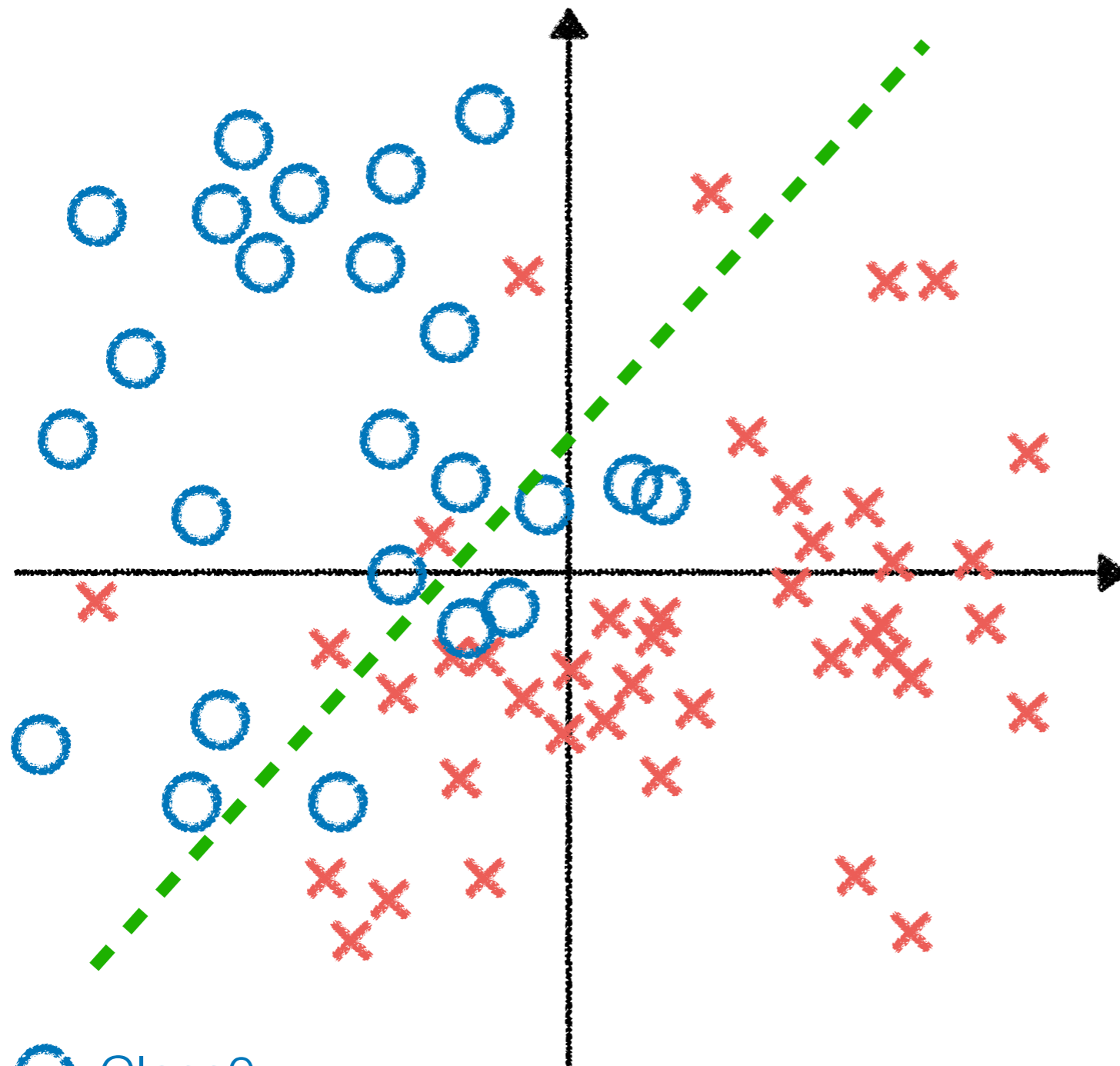


Data

○ Class0

× Class1

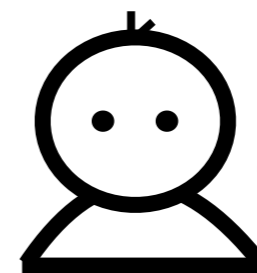
$$\operatorname{argmax}_E Q(E|M, H, D, \quad)$$



○ Class0

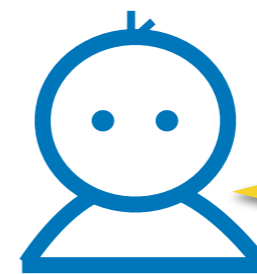
× Class1

## Human



newbie

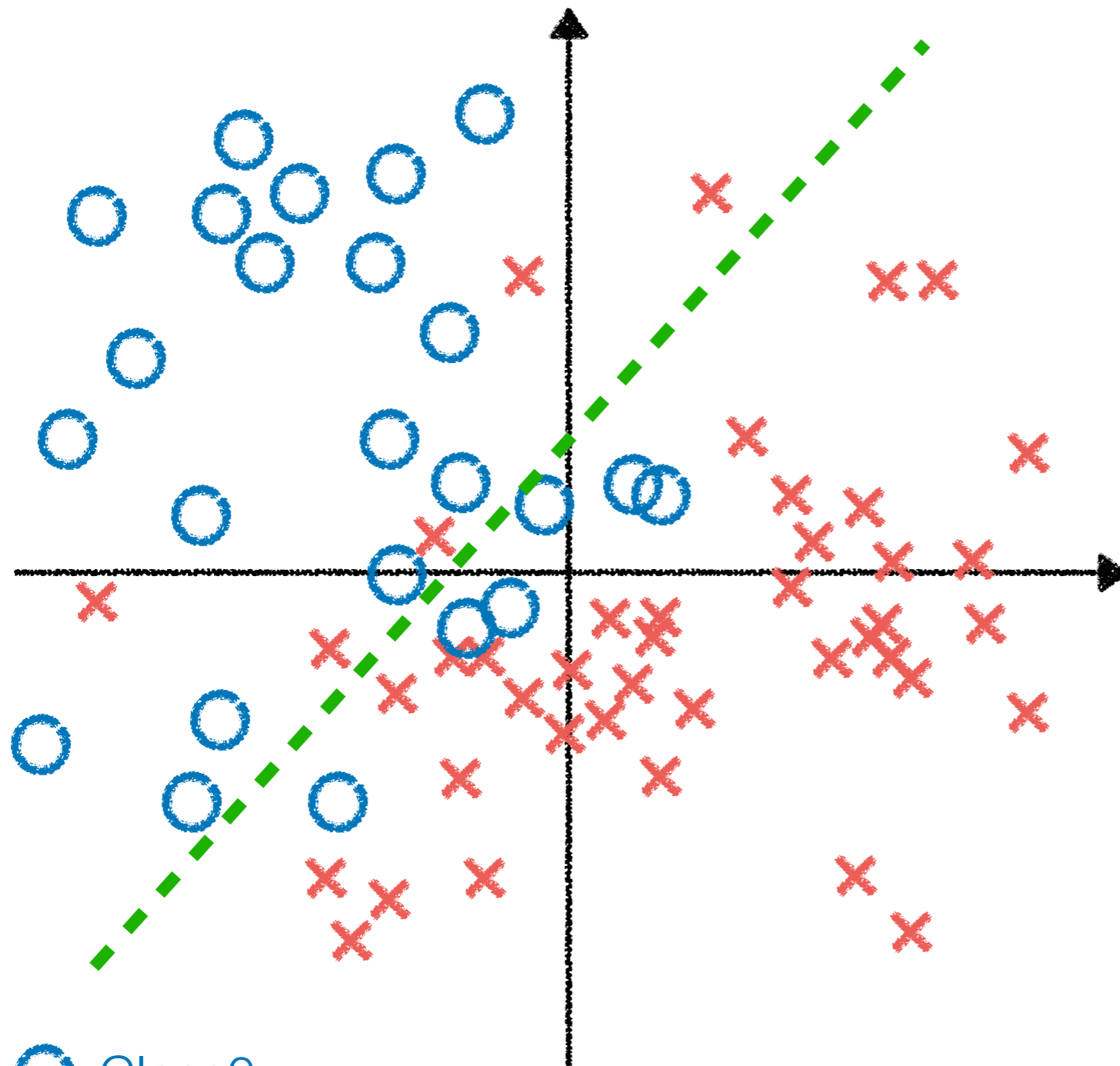
What's ML?



expert (you)

If I were you, I would train a neural network.

$$\operatorname{argmax}_E Q(E|M, H, D, \underline{T})$$



○ Class0

× Class1

**Task**

- Local vs. global
- Simple explanations vs. more complex but more accurate explanations
- Low or high stake domains

# Agenda

Post-training explanations

$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

Building inherently interpretable models

$$\operatorname{argmax}_{E, M} Q(\mathbf{Explanation}, \mathbf{Model} | \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

# Agenda

$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$

# Agenda

1. Revisit some existing methods:  
Sanity check questions

2. Make explanations  
that work for lay people.

$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

3. Understand how humans  
understand explanations

4. Make explanations to detect  
trustworthy predictions.



# Agenda

1. Revisit some existing methods:  
Sanity check questions

2. Make explanations  
that work for lay people.

$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

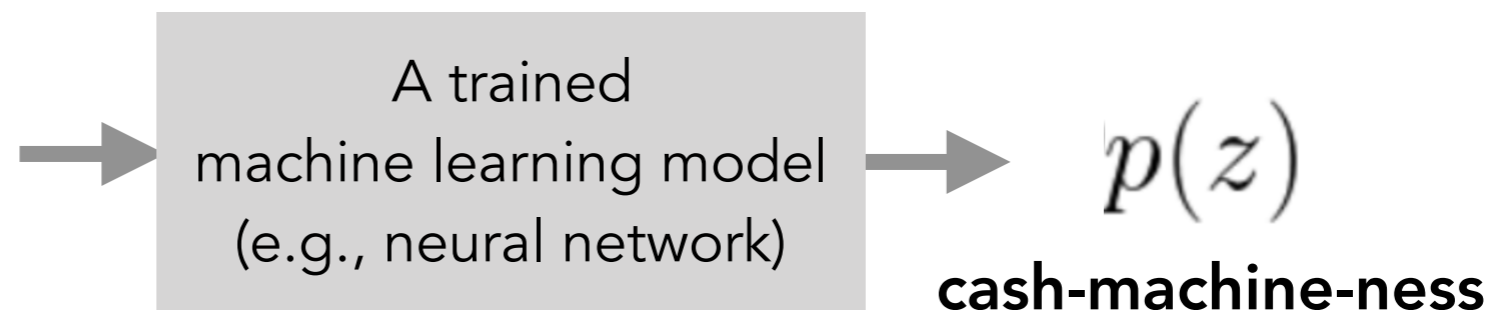
3. Understand how humans  
understand explanations

4. Make explanations to detect  
trustworthy predictions.

# Problem:

## Post-training explanation

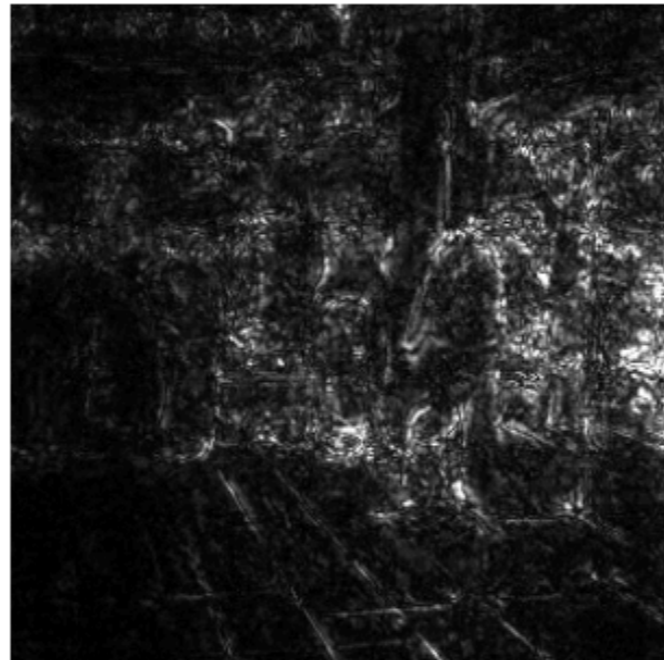
$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$



Why was this a cash machine?

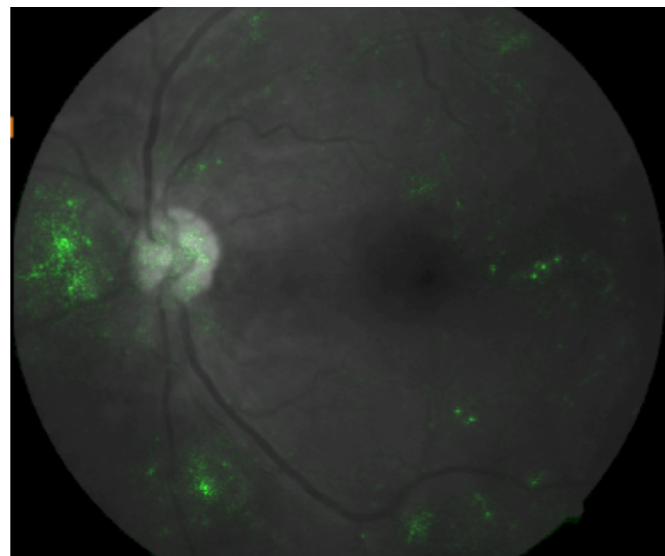
One of the most popular interpretability methods for images:

# Saliency maps

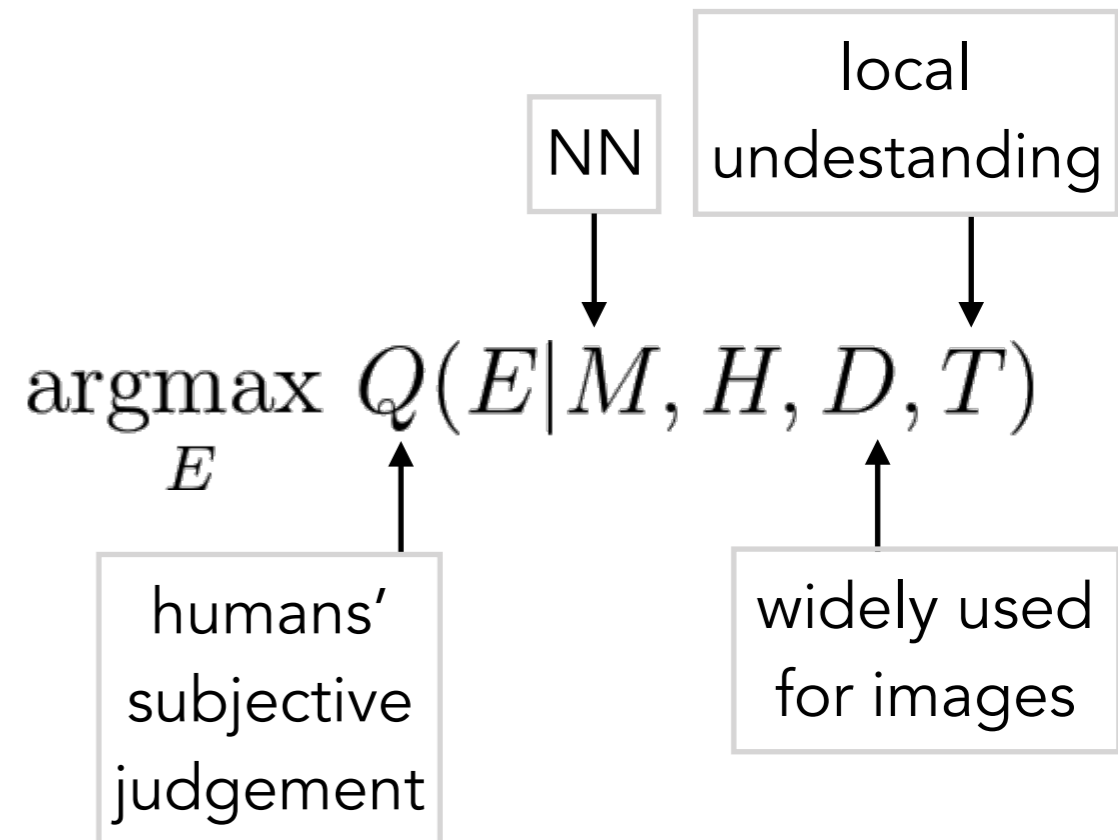


Used for image classification and medical applications.

$$\begin{aligned} \text{a logit} &\rightarrow \frac{\partial p(z)}{\partial z} \\ \text{pixel } i,j &\rightarrow \frac{\partial x_{i,j}}{\partial z} \end{aligned}$$



picture credit: @sayres



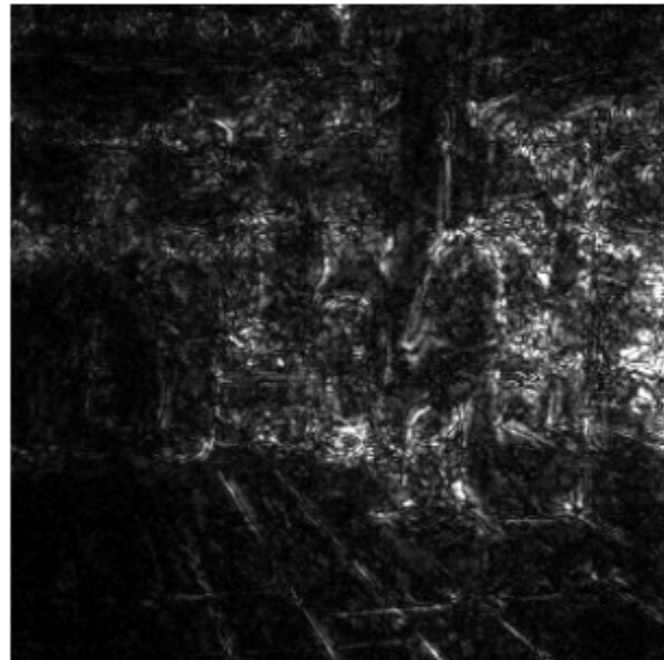
SmoothGrad [Smilkov, Thorat, K., Viégas, Wattenberg '17]

Integrated gradient [Sundararajan, Taly, Yan '17]



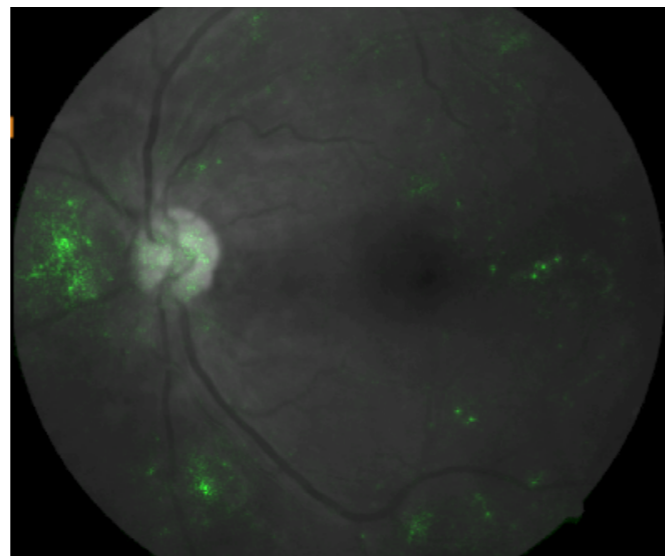
One of the most popular interpretability methods for images:

# Saliency maps



Used for image classification and medical applications.

$$\begin{aligned} \text{a logit} &\rightarrow \frac{\partial p(z)}{\partial z} \\ \text{pixel } i,j &\rightarrow \frac{\partial x_{i,j}}{\partial z} \end{aligned}$$



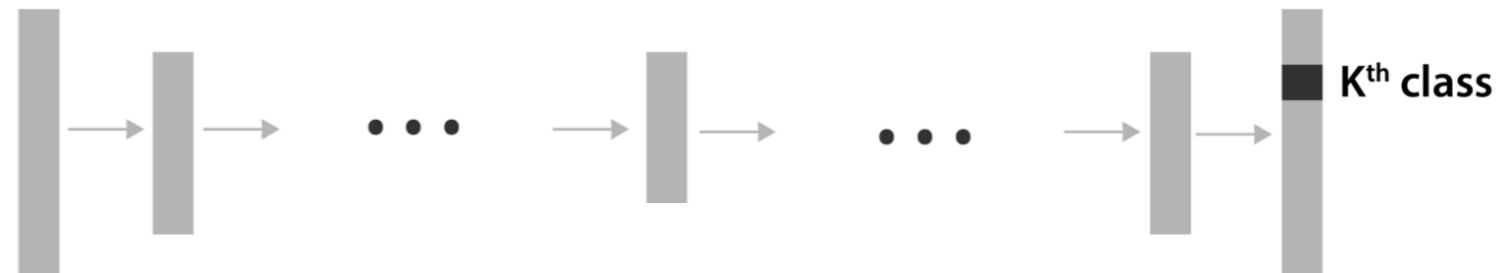
$$\operatorname{argmax}_E Q(E|M, H, D, T)$$

Sanity check:

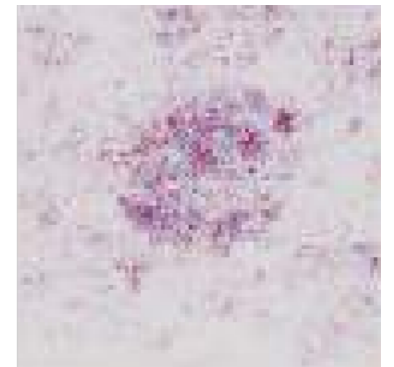
If I change M a lot, will human perceive that E has changed a lot?

# Some confusing behaviors of saliency maps.

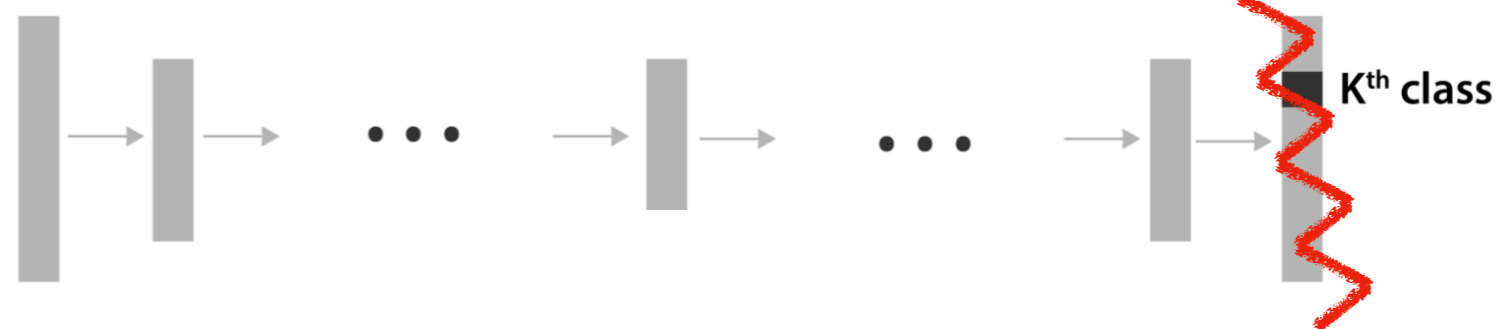
Original Image



Saliency map

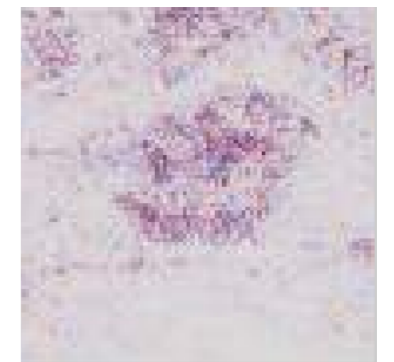


Original Image



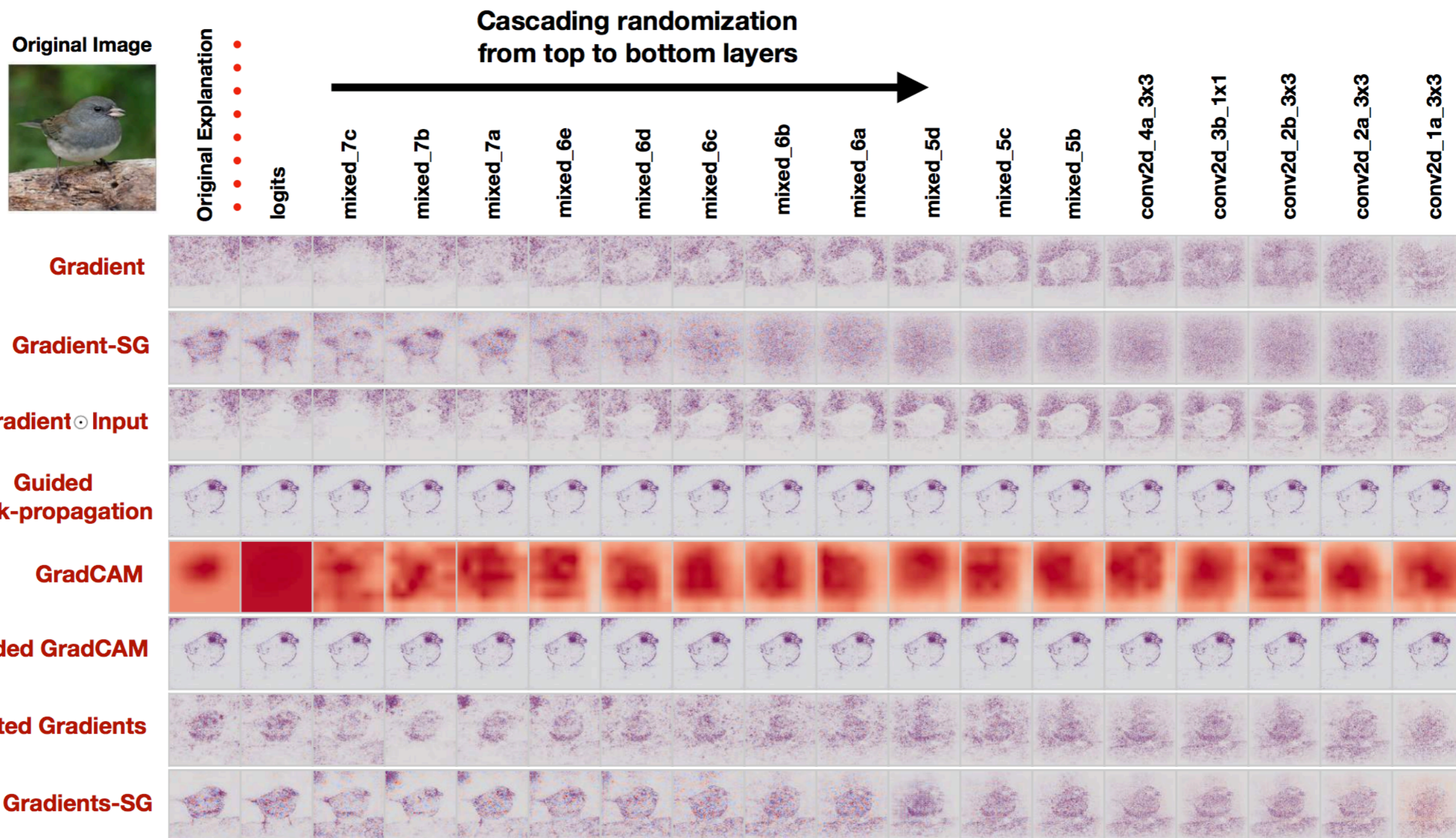
Randomized weights!  
Network now makes garbage prediction.

!!!!!!????!?





# Some saliency maps look similar when we randomize the network.

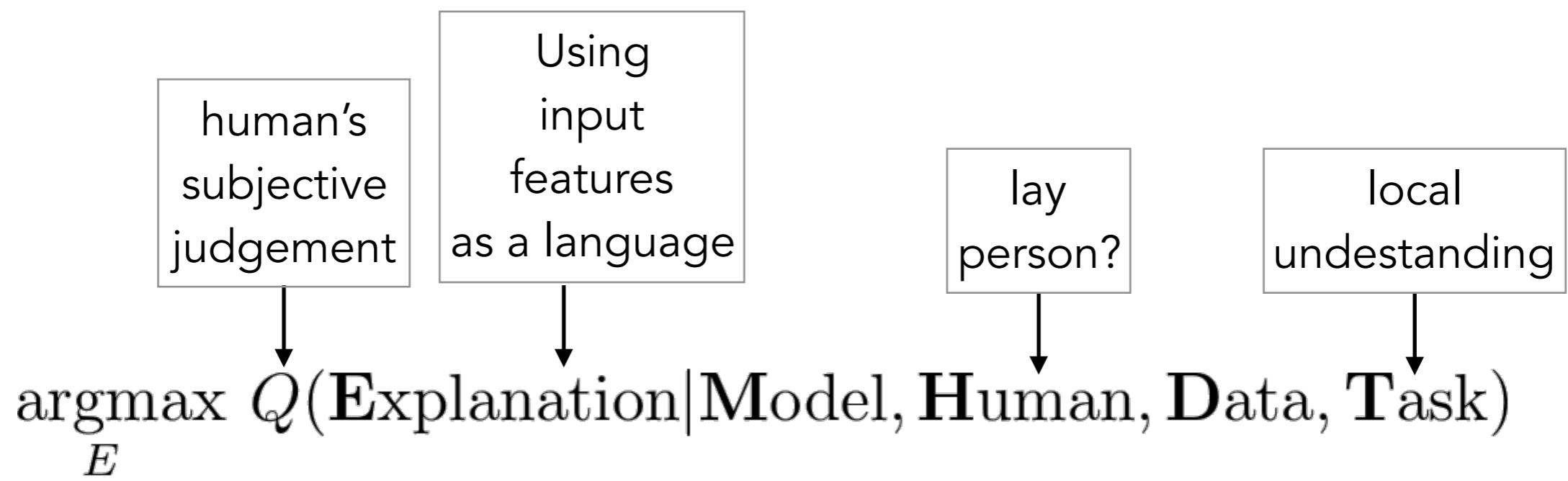


# What can we learn from this?

- Potential human confirmation bias: Just because it “makes sense” to humans, doesn’t mean they reflect evidence for the prediction.
- Our discovery is consistent with other findings [Nie, Zhang, Patel '18] [Ulyanov, Vedaldi, Lempitsky '18]
- Some of these methods have been shown to be useful in practice. Explaining predictions or features? More studies needed.



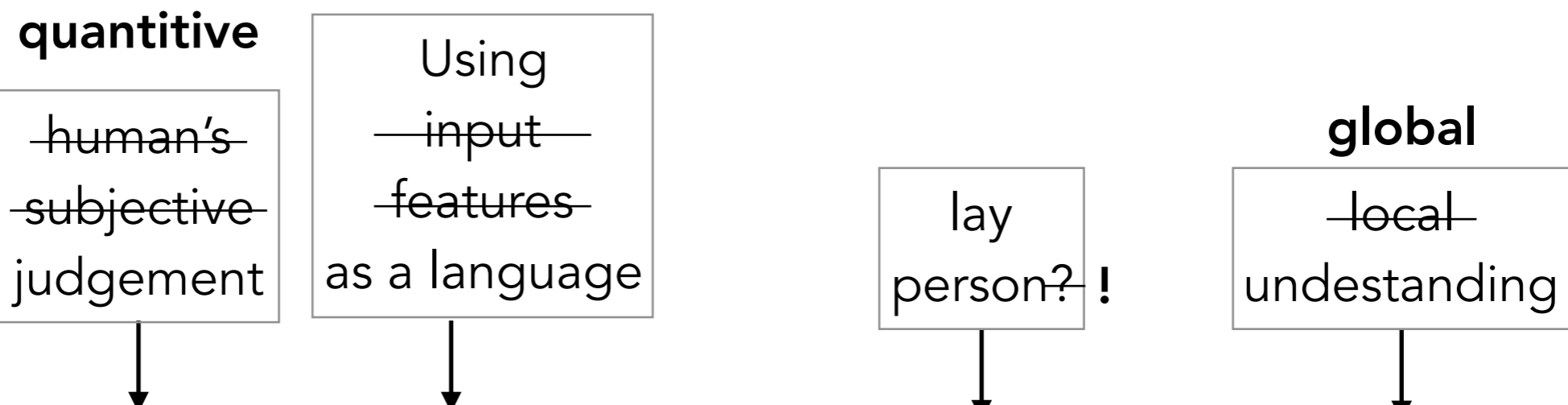
# What can we do better? Creating a wishlist.





# What can we do better? Creating a wishlist.

Something more human-friendly?



$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

# Agenda

1. Revisit some existing methods:  
Sanity check questions

2. Make explanations  
that work for lay people.

$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

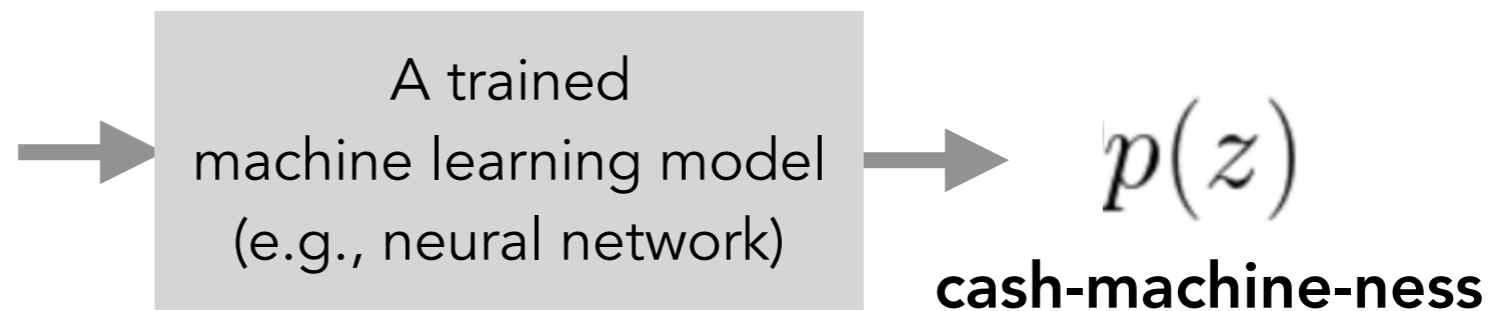
3. Understand how humans  
understand explanations

4. Make explanations to detect  
trustworthy predictions.

# Problem:

## Post-training explanation

$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$



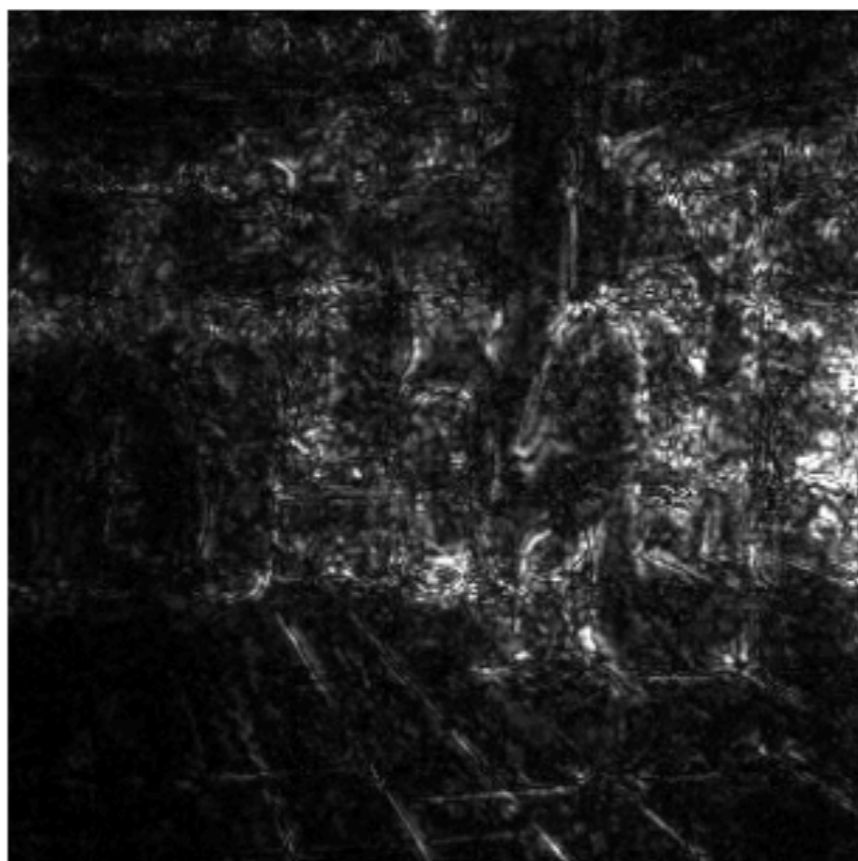
Why was this a cash machine?

# Common solution: Saliency map

prediction:  
Cash machine



Let's use this to help us think about what we really want to ask.

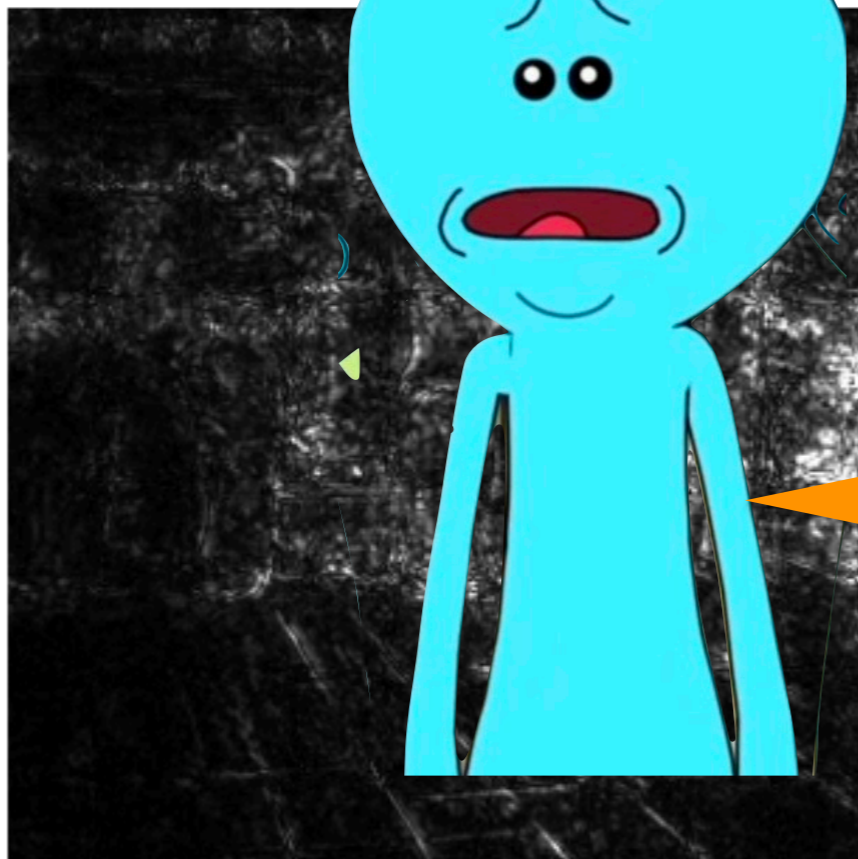


<https://pair-code.github.io/saliency/>



# What we really want to ask...

prediction:  
Cash machine



Were there more pixels on the cash machine than on the person?

Did the 'human' concept matter?  
Did the 'wheels' concept matter?

Which concept mattered more?

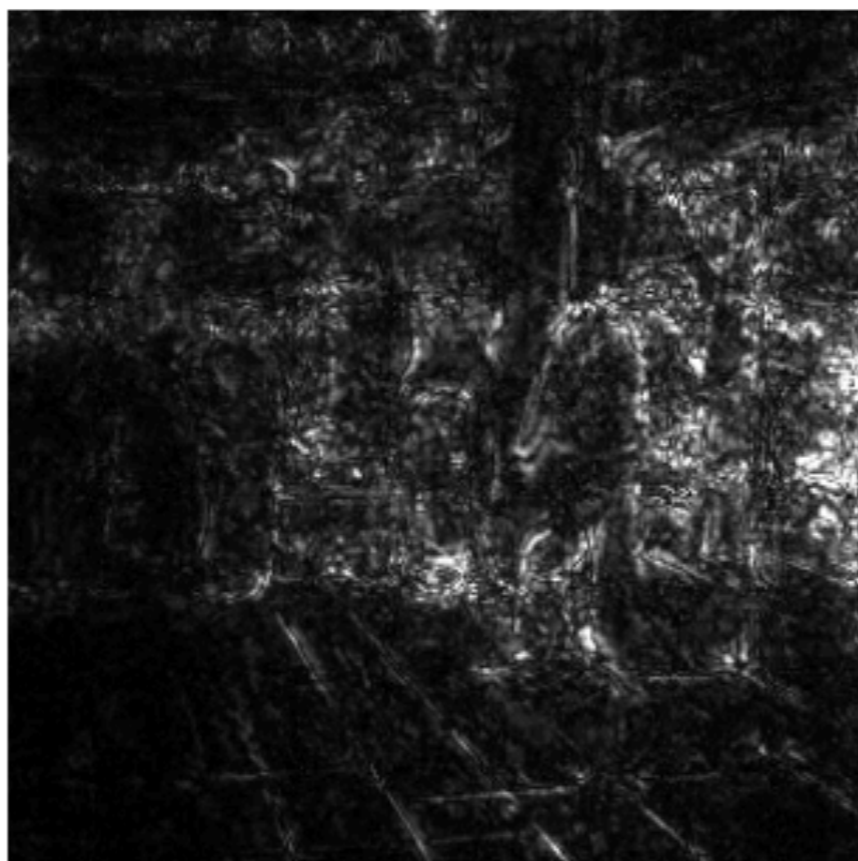
Is this true for all other cash machine predictions?

Oh no! I can't express these concepts as pixels!!  
They weren't my input features either!

<https://pair-code.github.io/saliency/>

# What we really want to ask...

prediction:  
Cash machine



Were there more pixels on the cash machine than on the person?

Did the 'human' concept matter?  
Did the 'wheels' concept matter?

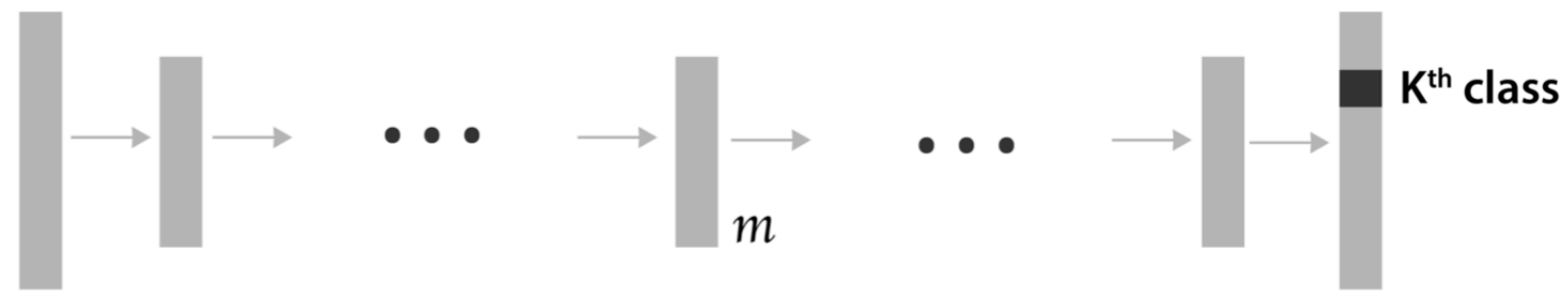
Which concept mattered more?

Is this true for all other cash machine predictions?

Wouldn't it be great if we can **quantitatively** measure how important *any* of these **user-chosen concepts** are?

<https://pair-code.github.io/saliency/>

# Goal of TCAV: Testing with Concept Activation Vectors

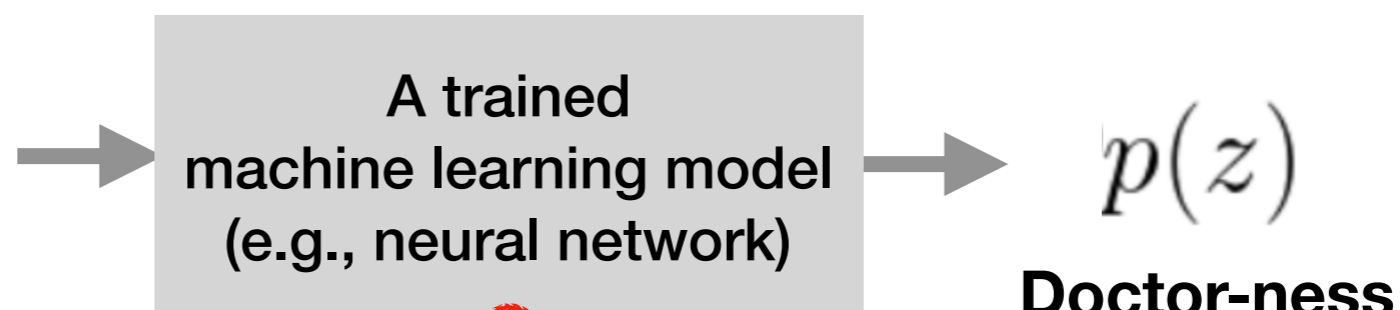


**Quantitative** explanation: how much a **concept** (e.g., gender, race) was important for a **prediction** in a trained model.

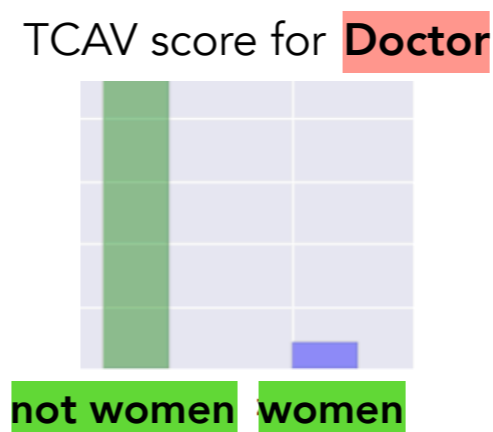
...even if the **concept** was not part of the training.



# Goal of TCAV: Testing with Concept Activation Vectors



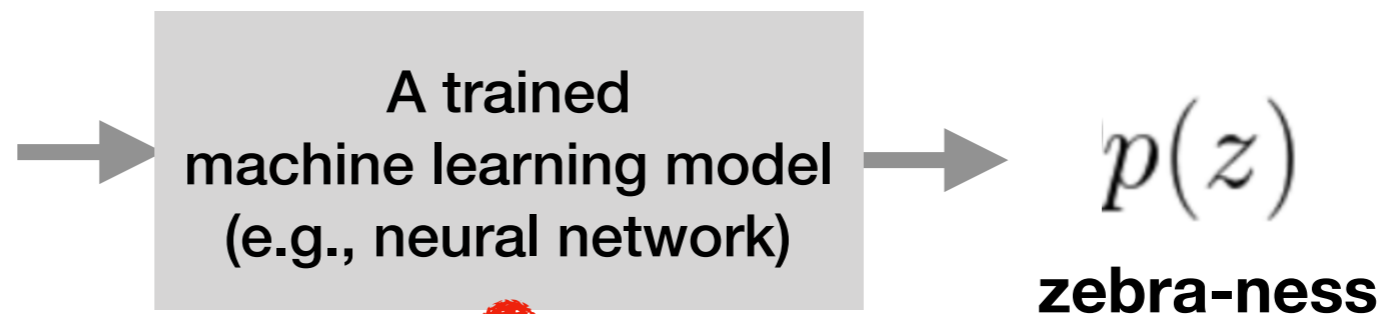
Was **gender concept** important to this **doctor** image classifier?



TCAV provides **quantitative importance** of a concept **if and only if** your network learned about it.



# Goal of TCAV: Testing with Concept Activation Vectors



Was striped concept important to this zebra image classifier?



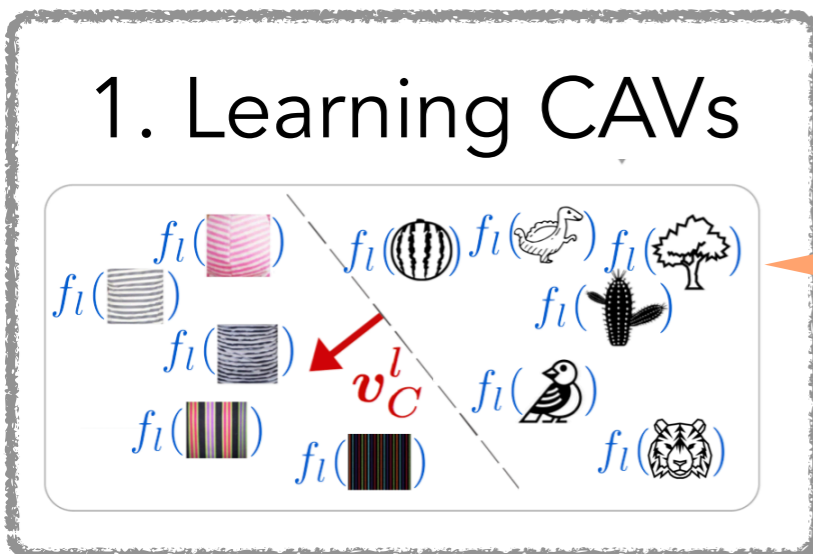
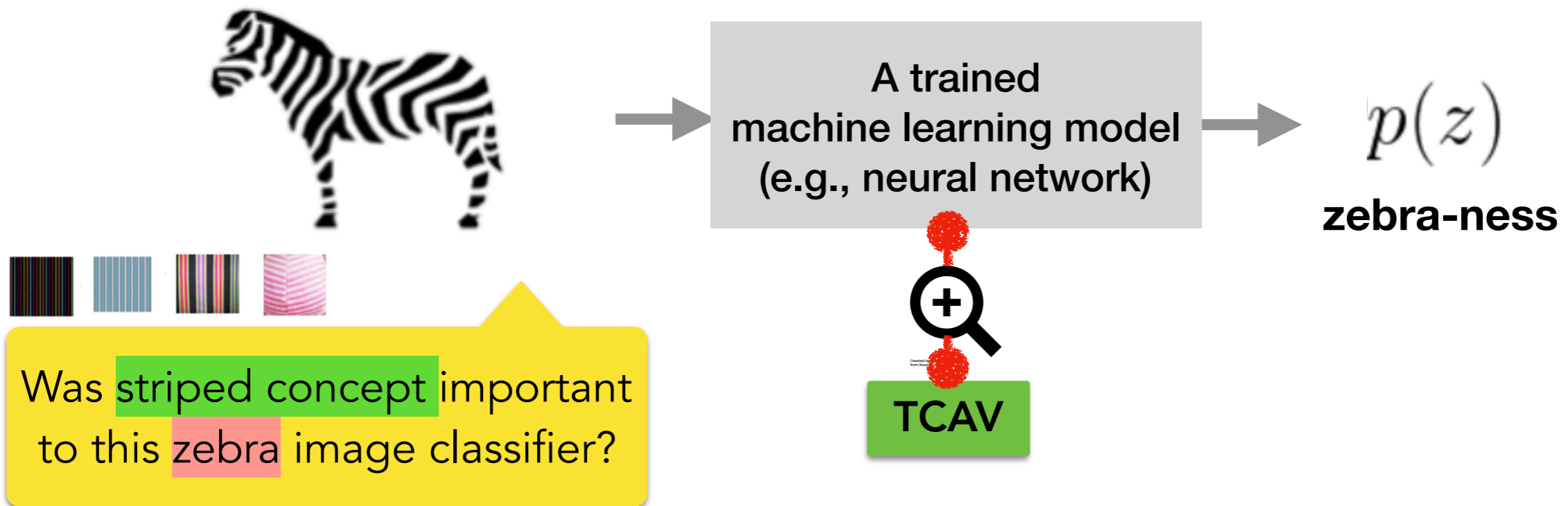
TCAV score for Zebra



TCAV provides **quantitative importance** of a concept **if and only if** your network learned about it.

# TCAV:

## Testing with Concept Activation Vectors



1. How to define concepts?

# Defining concept activation vector (CAV)

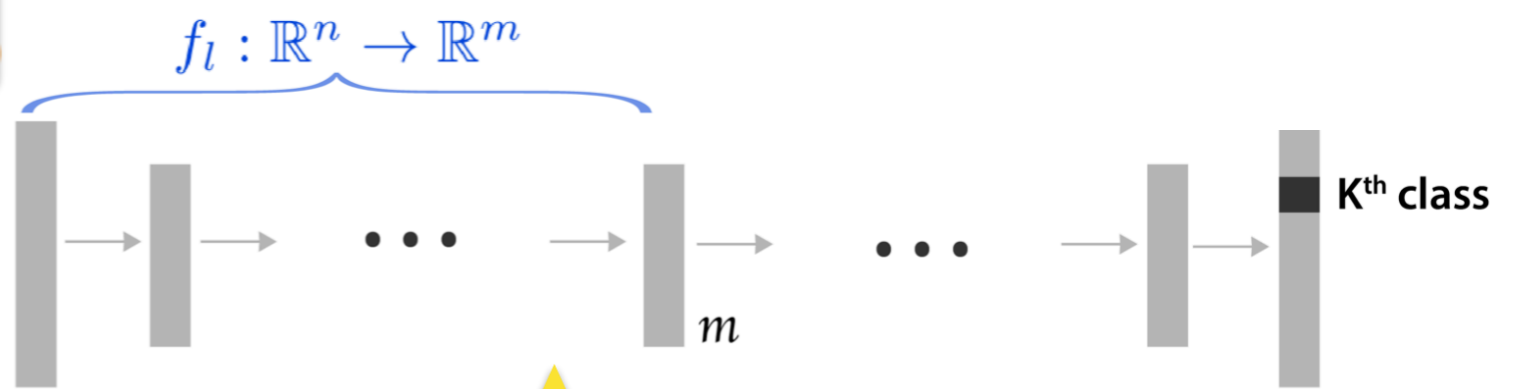
## Inputs:

a



Examples of concepts

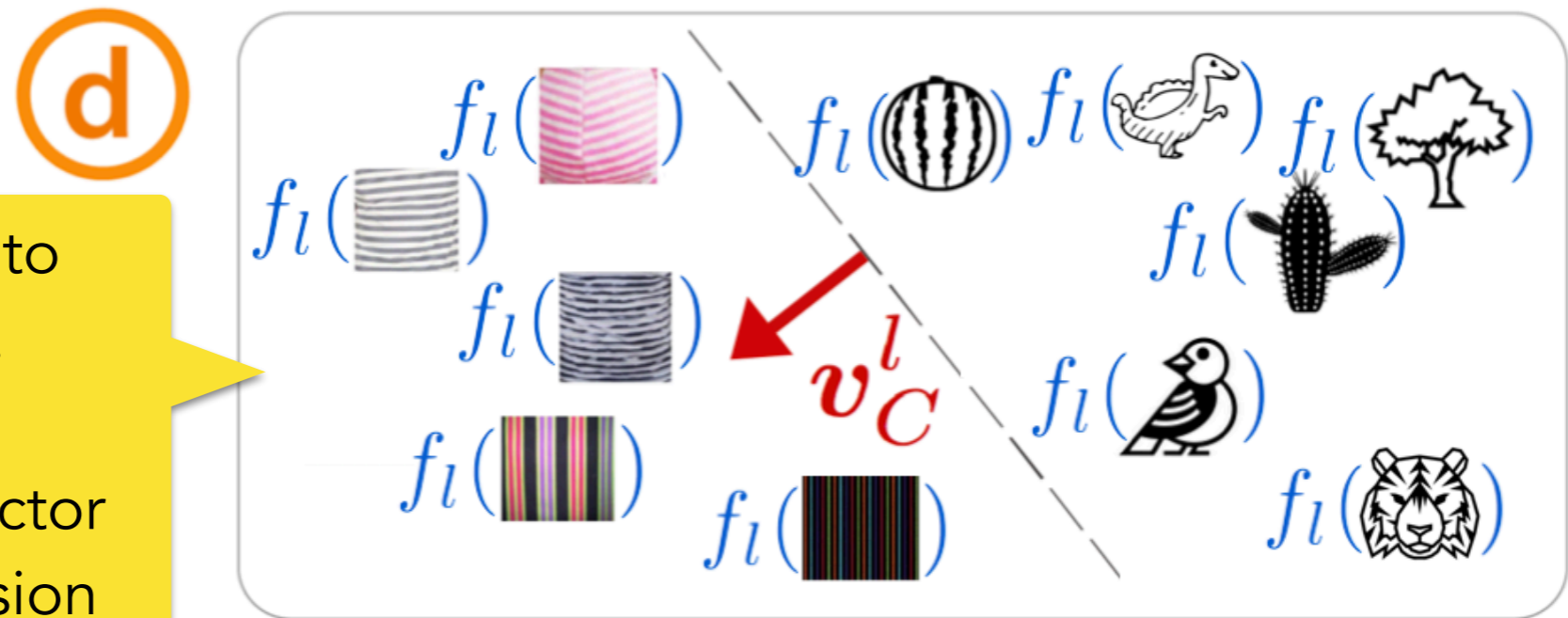
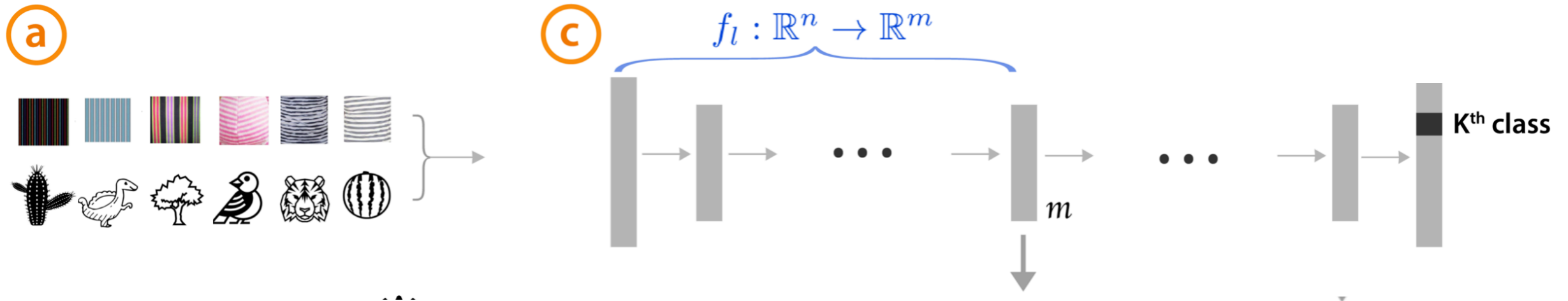
Random images



A trained network under investigation and Internal tensors

# Defining concept activation vector (CAV)

## Inputs:



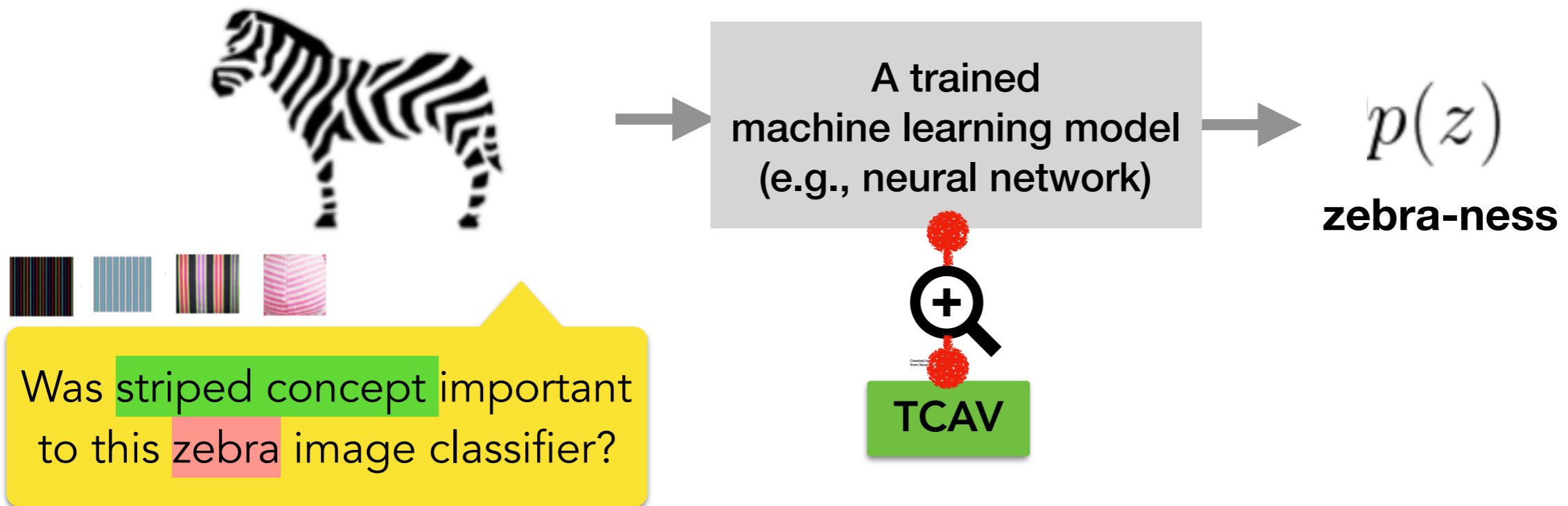
Train a linear classifier to separate activations.

CAV ( $v_C^l$ ) is the vector **orthogonal** to the decision boundary.

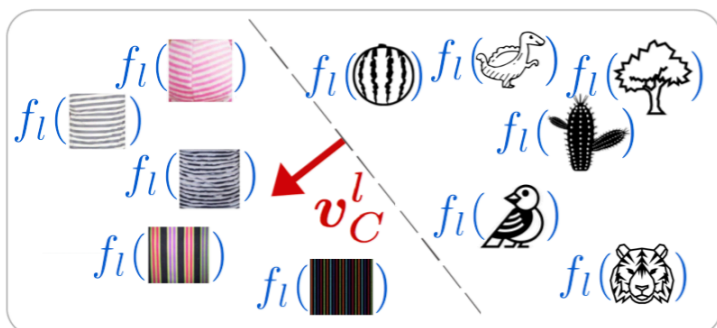
[Smilkov '17, Bolukbasi '16, Schmidt '15]

# TCAV:

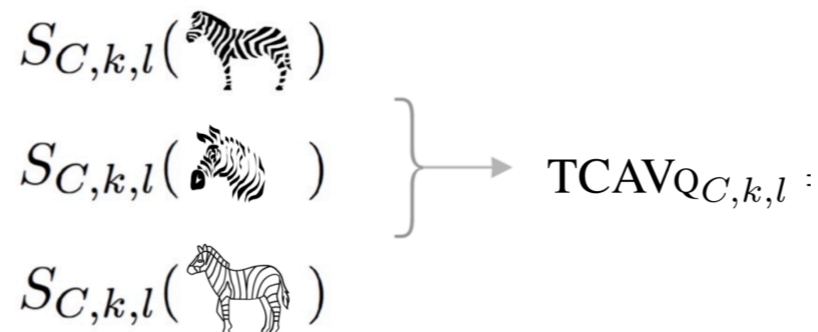
## Testing with Concept Activation Vectors



### 1. Learning CAVs



### 2. Getting TCAV score

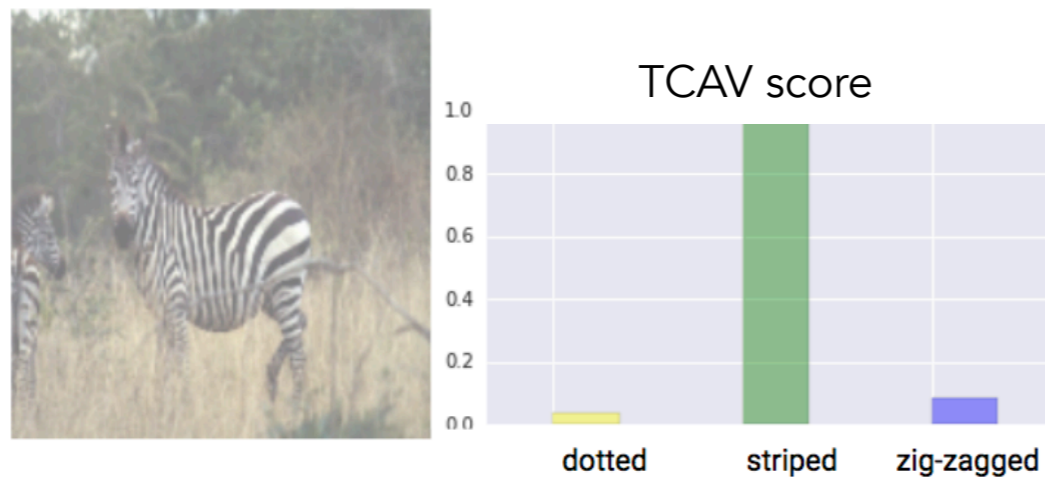


2. How are the CAVs useful to get explanations?

# TCAV core idea:

Derivative with CAV to get prediction sensitivity

## TCAV



$$\begin{aligned}
 & S_{C,k,l}(\text{zebra}) \\
 & S_{C,k,l}(\text{zebra head}) \\
 & S_{C,k,l}(\text{zebra body}) \\
 & S_{C,k,l}(\text{zebra tail})
 \end{aligned}$$

zebra-ness  $\rightarrow \frac{\partial p(z)}{\partial \mathbf{v}_C^l} = S_{C,k,l}(\mathbf{x})$

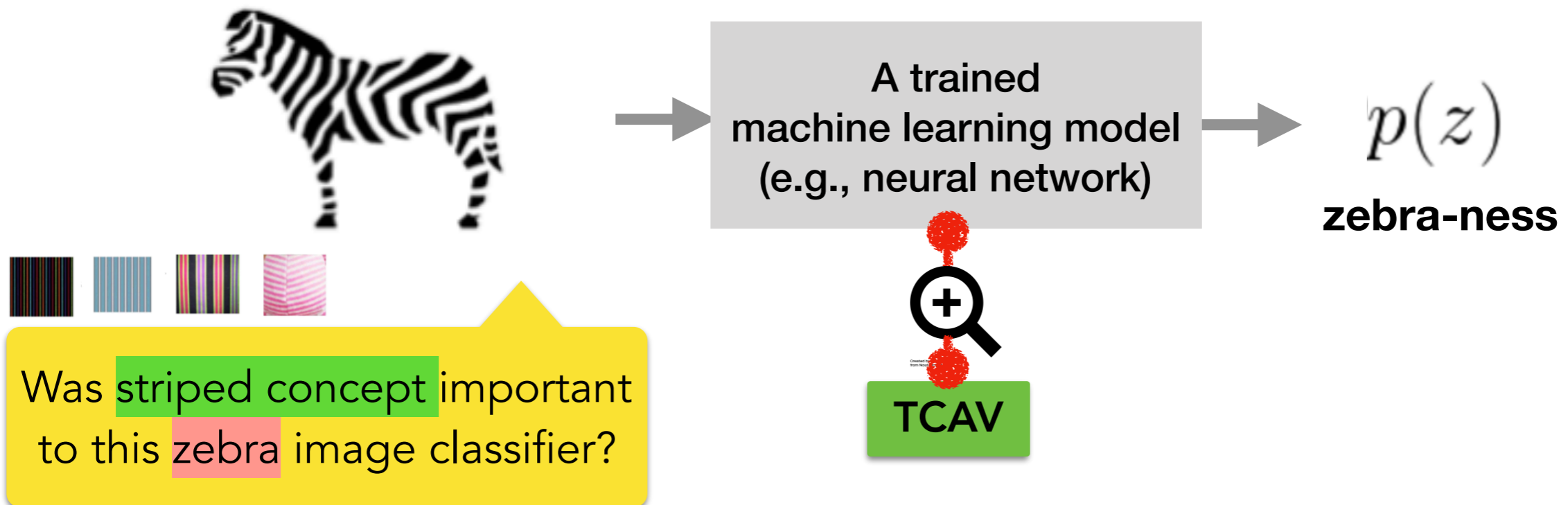
striped CAV  $\rightarrow$

$$\text{TCAV}_{Q_{C,k,l}} = \frac{|\{\mathbf{x} \in X_k : S_{C,k,l}(\mathbf{x}) > 0\}|}{|X_k|}$$

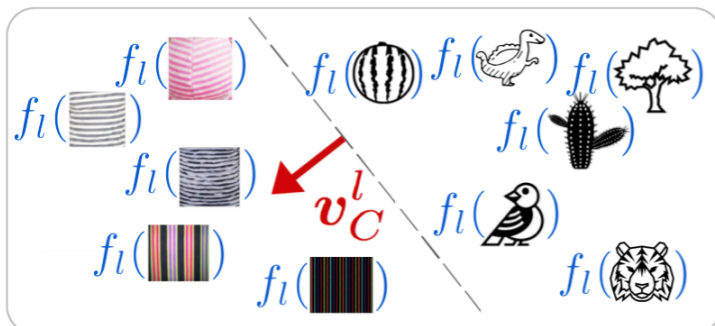
Directional derivative with CAV

# TCAV:

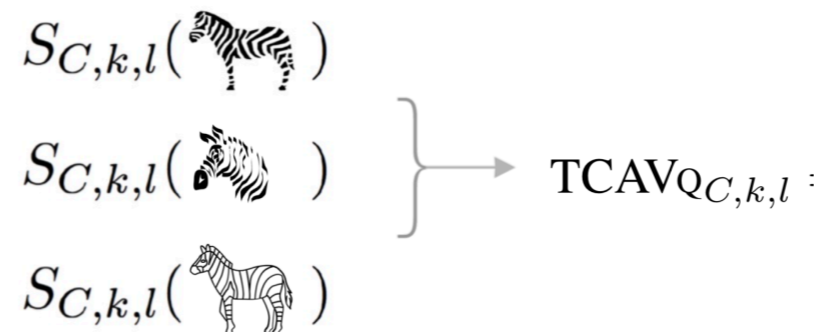
## Testing with Concept Activation Vectors



### 1. Learning CAVs



### 2. Getting TCAV score



### 3. CAV validation

Qualitative  
Quantitative



Quantitative validation:

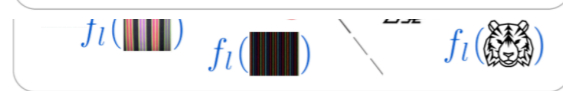
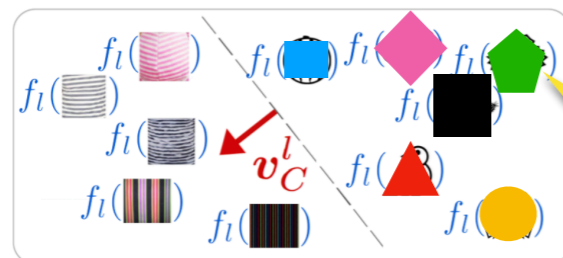
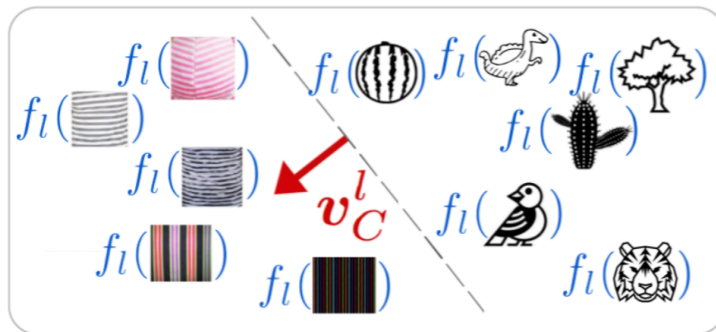
# Guarding against spurious CAV

Did my CAVs returned high sensitivity by chance?



Quantitative validation:

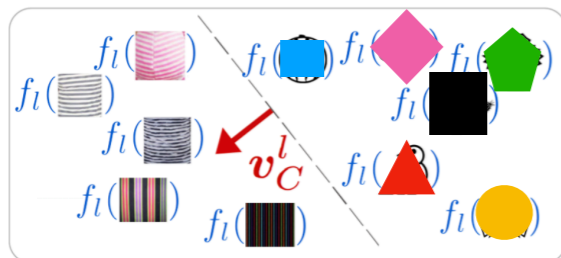
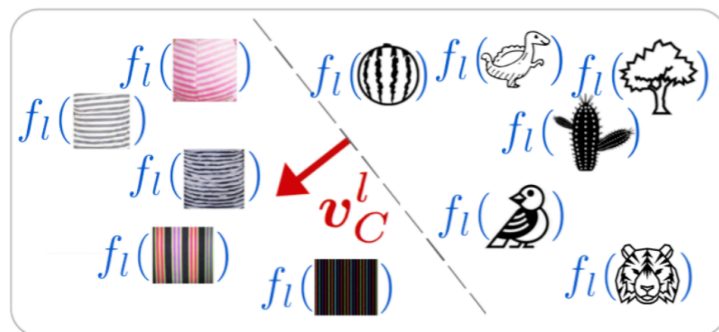
# Guarding against spurious CAV



Learn many stripes CAVs  
using different sets of  
random images

Quantitative validation:

# Guarding against spurious CAV



Zebra

→  $\text{TCAV}_{Q_C, k, l} :$

⋮

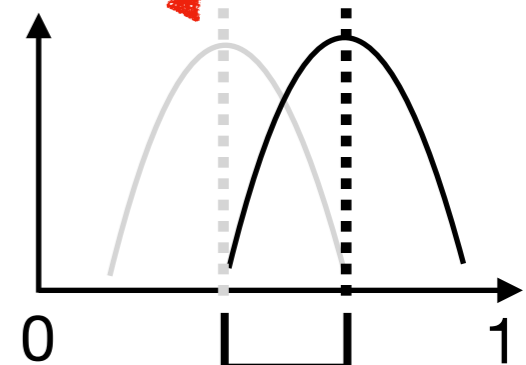
→  $\text{TCAV}_{Q_C, k, l} :$

→  $\text{TCAV}_{Q_C, k, l} :$

→  $\text{TCAV}_{Q_C, k, l} :$

⋮

TCAV score  
random



Check the distribution of  $\text{TCAV}_{Q_C, k, l}$  is statistically different from random using t-test

# Recap TCAV:

## Testing with Concept Activation Vectors

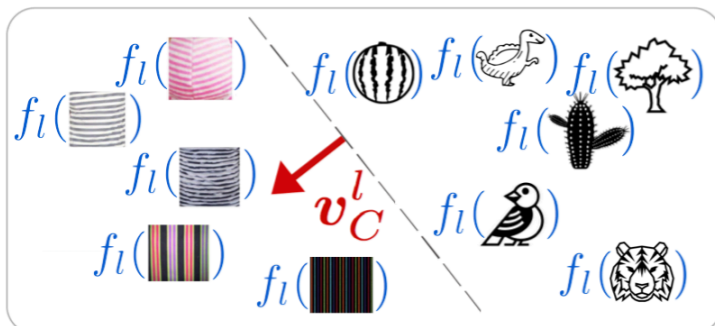


TCAV provides **quantitative importance** of a concept **if and only if** your network learned about it.

Even if your training data wasn't tagged with the **concept**

Even if your input feature did not include the **concept**

### 1. Learning CAVs



### 2. Getting TCAV score

$$\begin{matrix} S_{C,k,l}(\text{zebra}) \\ S_{C,k,l}(\text{zebra}) \\ S_{C,k,l}(\text{zebra}) \end{matrix} \left. \vphantom{\begin{matrix} S_{C,k,l}(\text{zebra}) \\ S_{C,k,l}(\text{zebra}) \\ S_{C,k,l}(\text{zebra}) \end{matrix}} \right\} \rightarrow \text{TCAV}_{QC,k,l}$$

### 3. CAV validation

Qualitative  
Quantitative

# Results

## 1. Sanity check experiment



## 2. Biases from Inception V3 and GoogleNet

## 3. Domain expert confirmation from Diabetic Retinopathy

DR level 4 Retina



TCAV for DR level 4



# Sanity check experiment

If we know the ground truth  
(important concepts),  
will TCAV match?



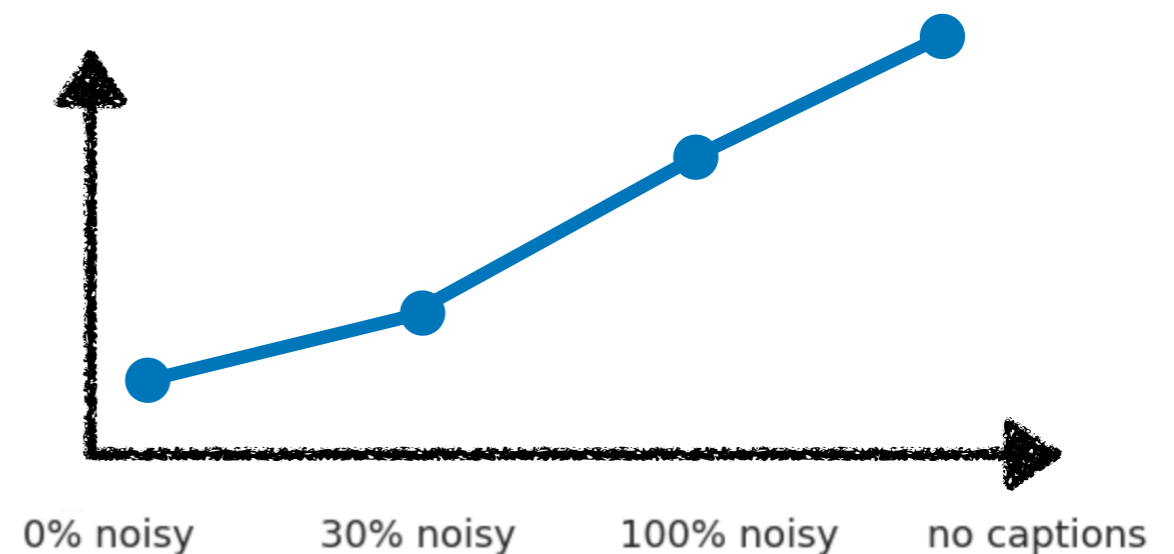
# Sanity check experiment setup



models can use either image or caption concept for classification.



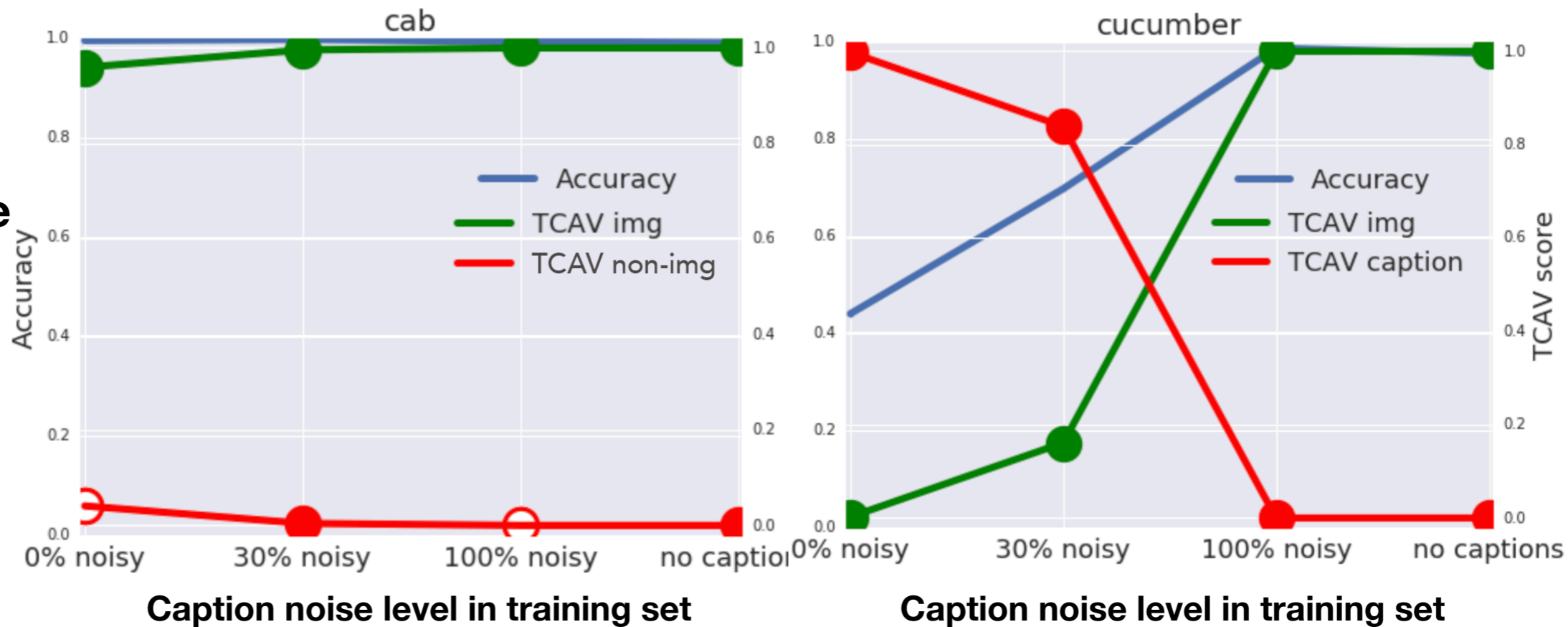
Test accuracy with no caption image = Importance of image concept



Caption noise level in training set of each model

# Sanity check experiment

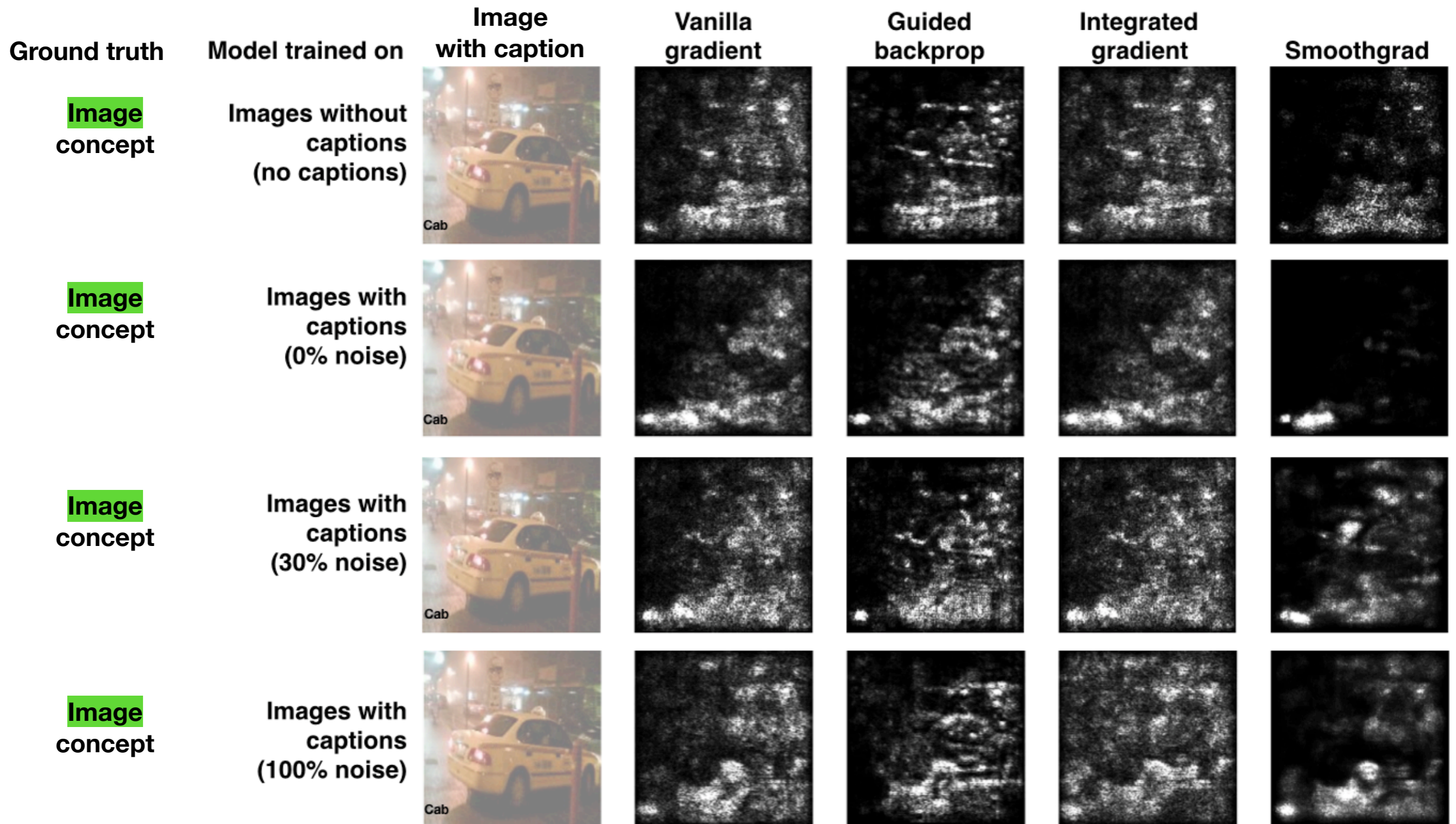
Test accuracy  
with  
no caption image





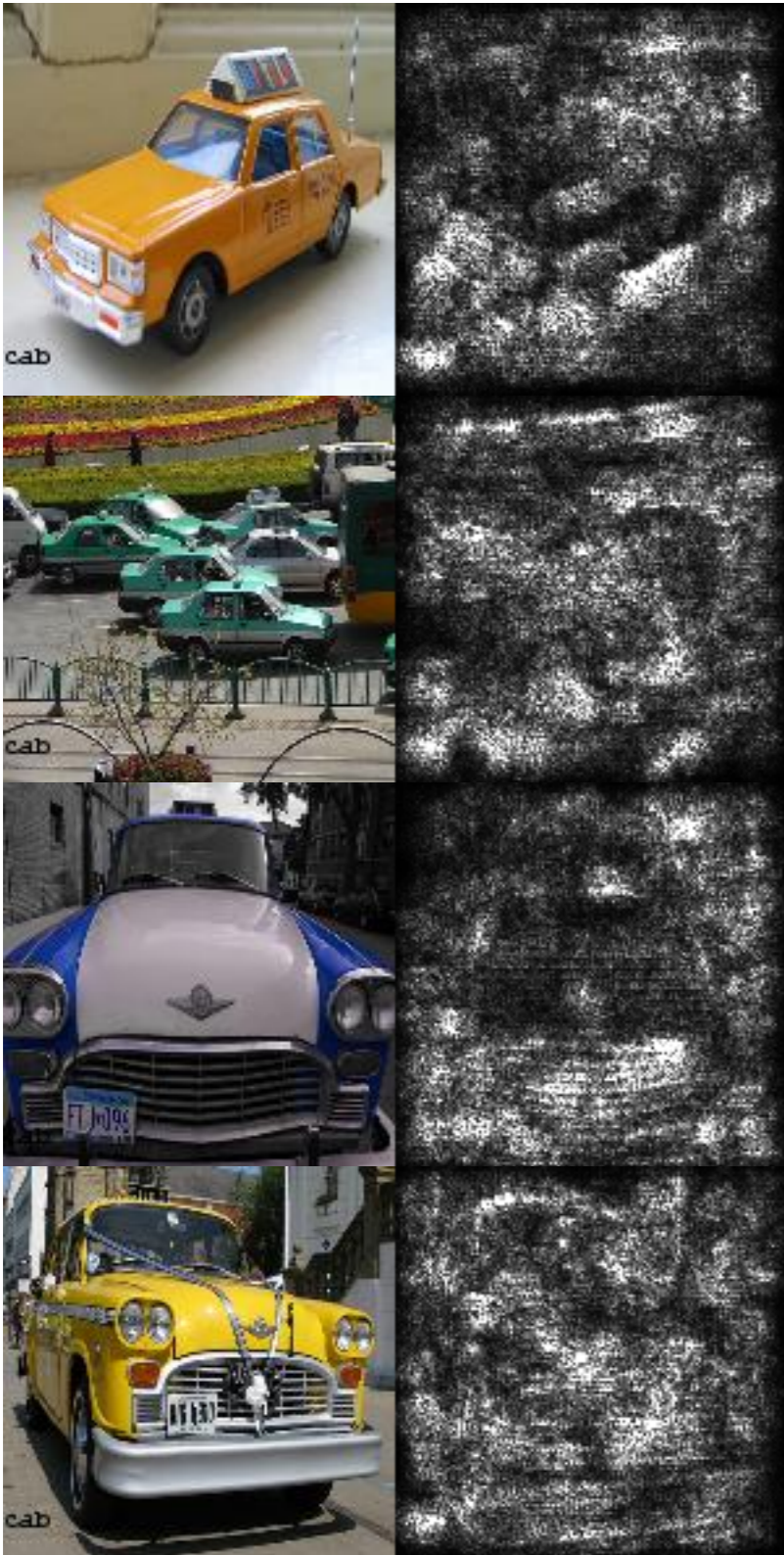
Cool, cool.  
Can saliency maps do this too?

# Can saliency maps communicate the same information?





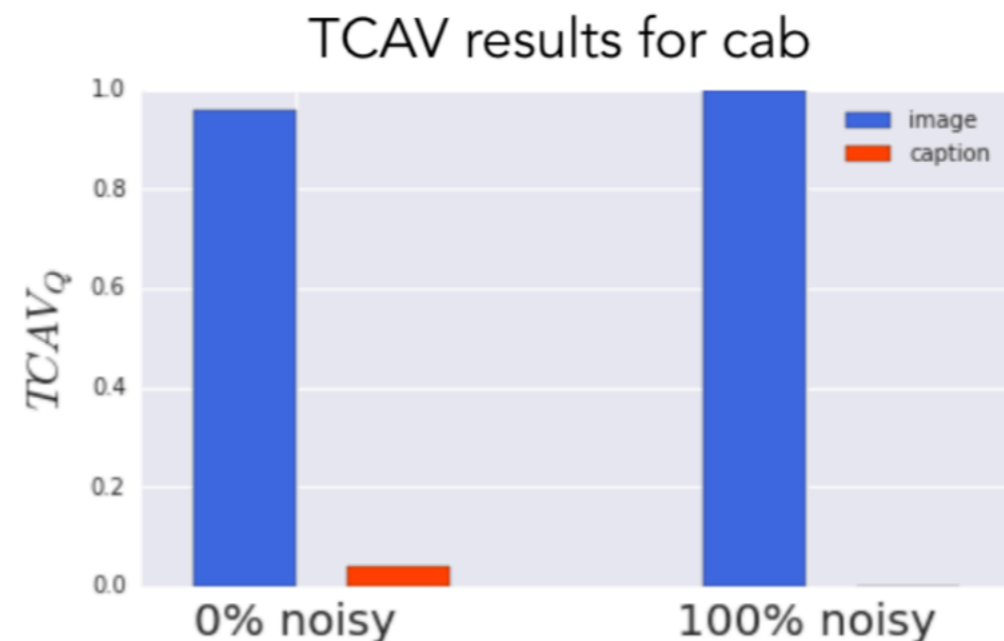
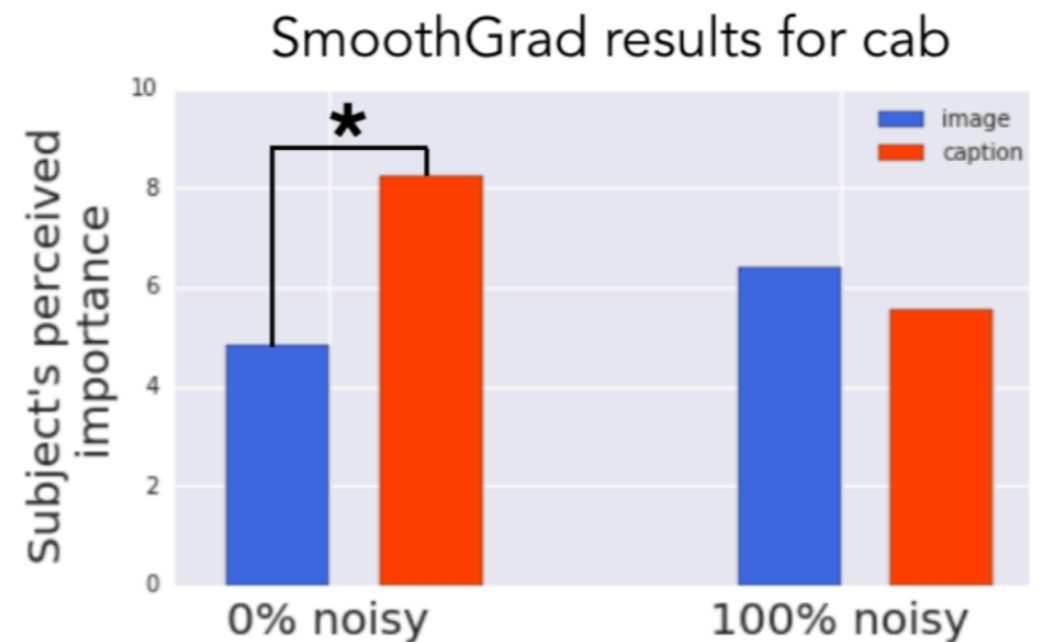
# Human subject experiment: Can saliency maps communicate the same information?



- 50 turkers are
- asked to judge importance of **image** vs. **caption** given saliency maps.
- asked to indicate their confidence
- shown 3 classes (cab, zebra, cucumber) x 2 saliency maps for one model

# Human subject experiment: Can saliency maps communicate the same information?

- Random chance: 50%
- Human performance with saliency map: 52%
- Humans can't agree: more than 50% no significant consensus





# Human subject experiment: Can saliency maps communicate the same information?

- Random chance: 50%
- Human performance with saliency map: 52%
- Humans can't agree: more than 50% no significant consensus
- Humans are **very** confident even when they are wrong.



# Results

1. Sanity check experiment



2. Biases from Inception V3 and GoogleNet

3. Domain expert confirmation from Diabetic Retinopathy

DR level 4 Retina

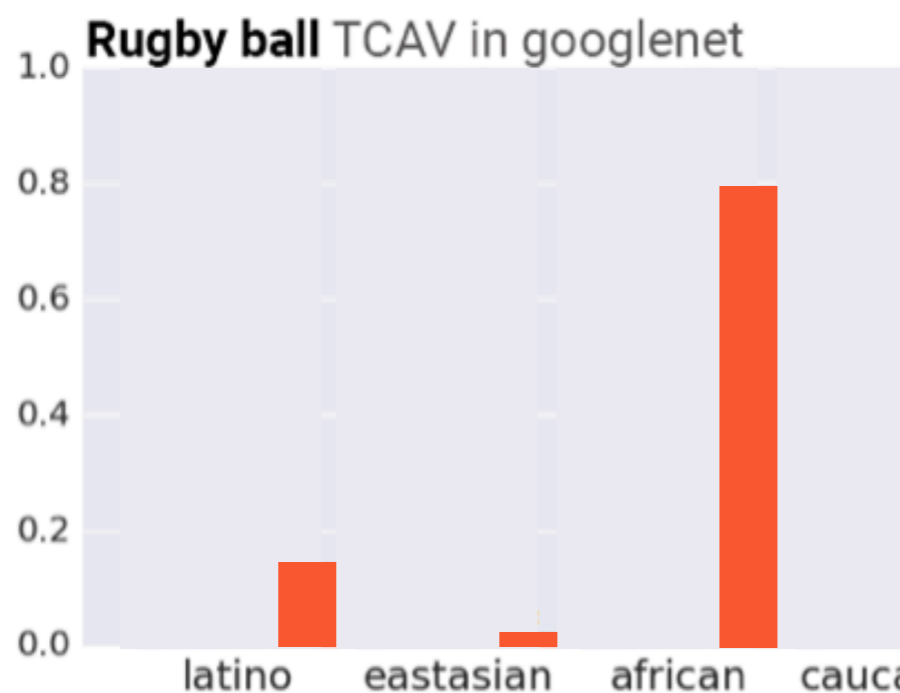
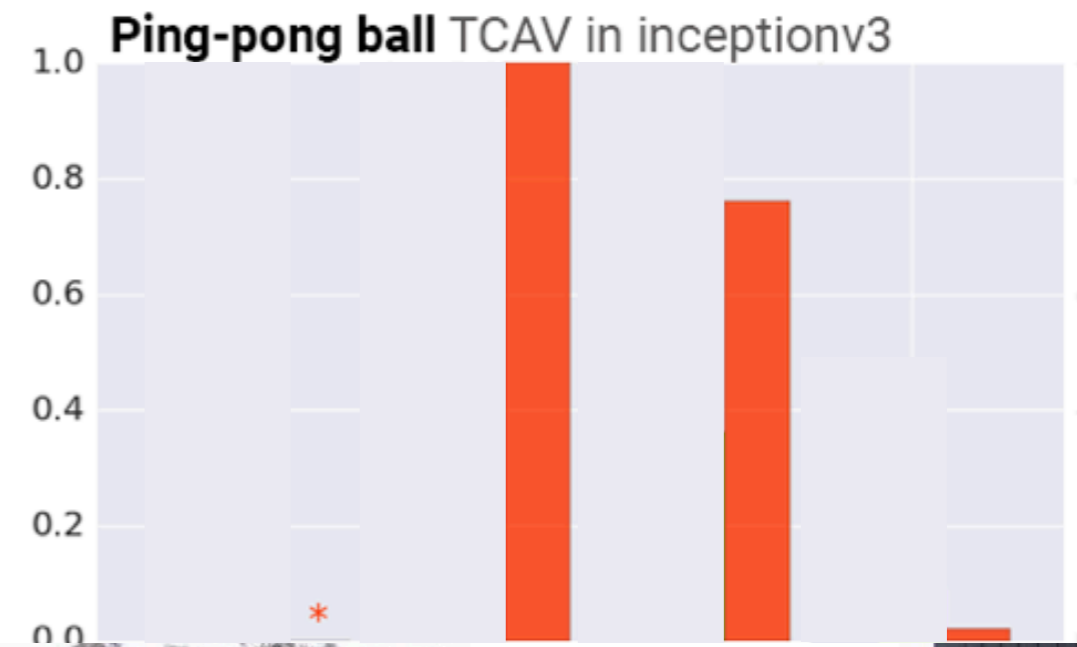
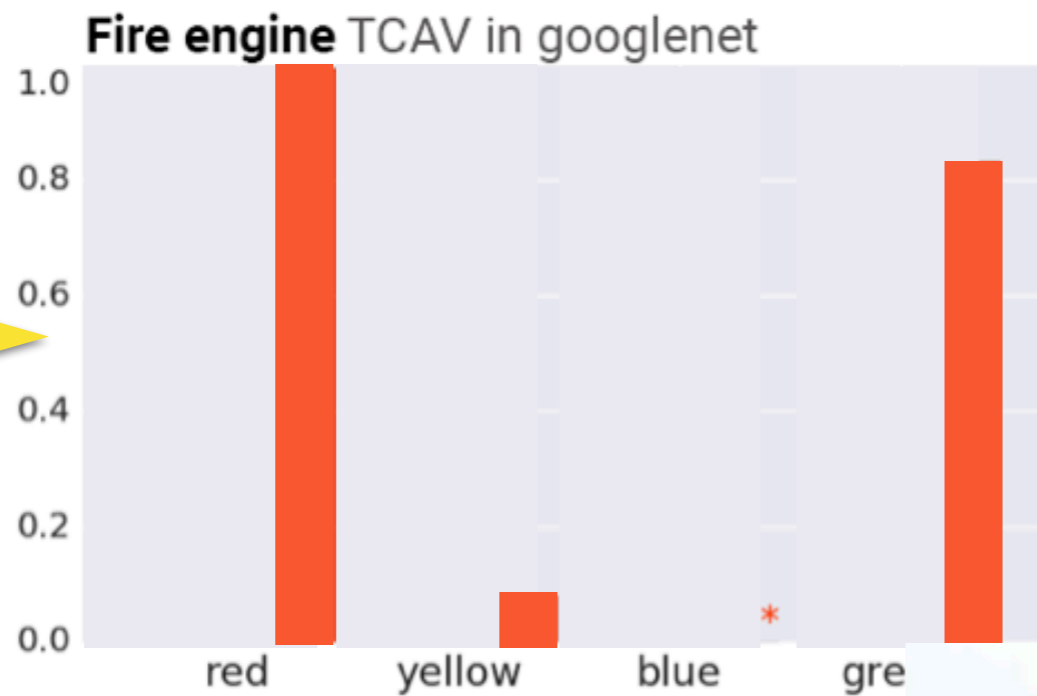


TCAV for DR level 4



# TCAV in Two widely used image prediction models

Geographical bias!

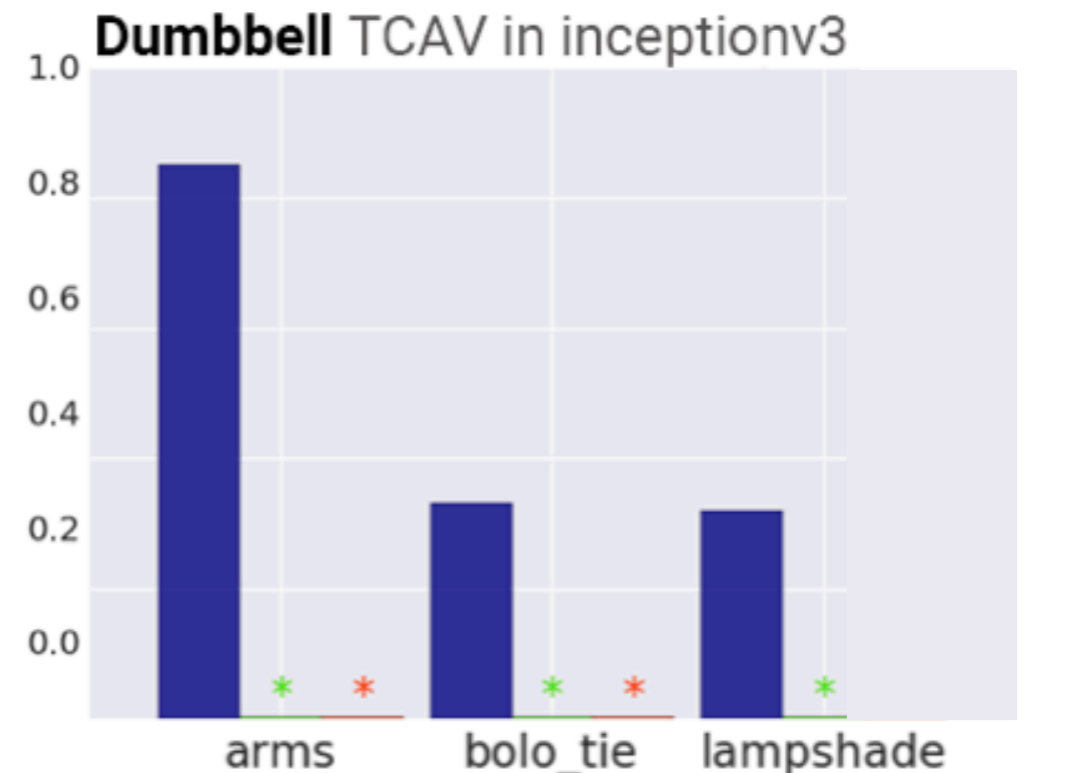
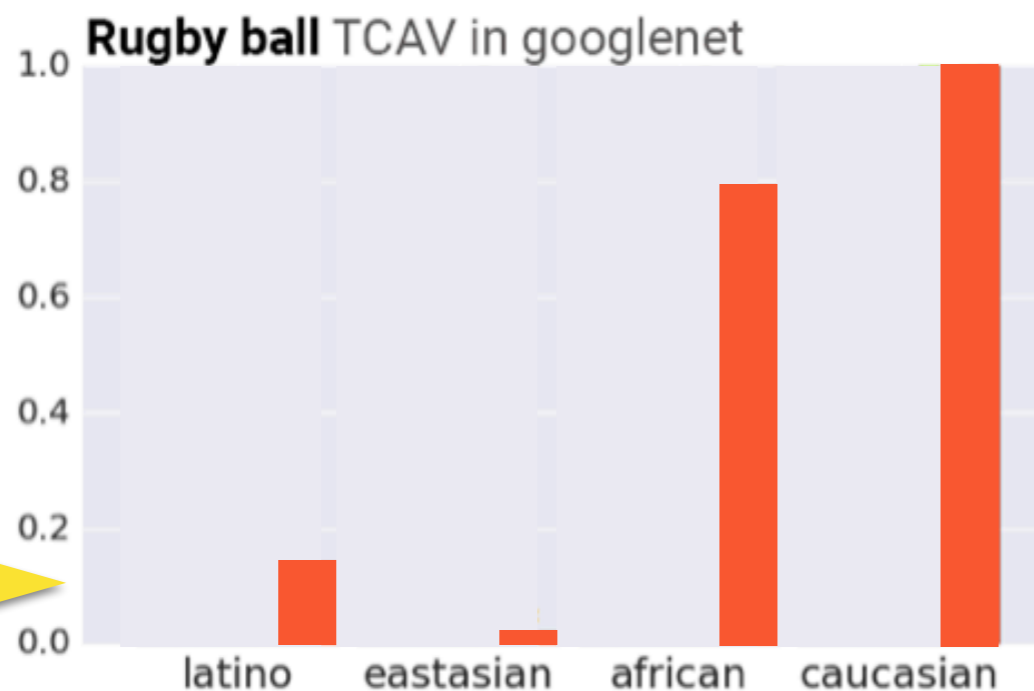
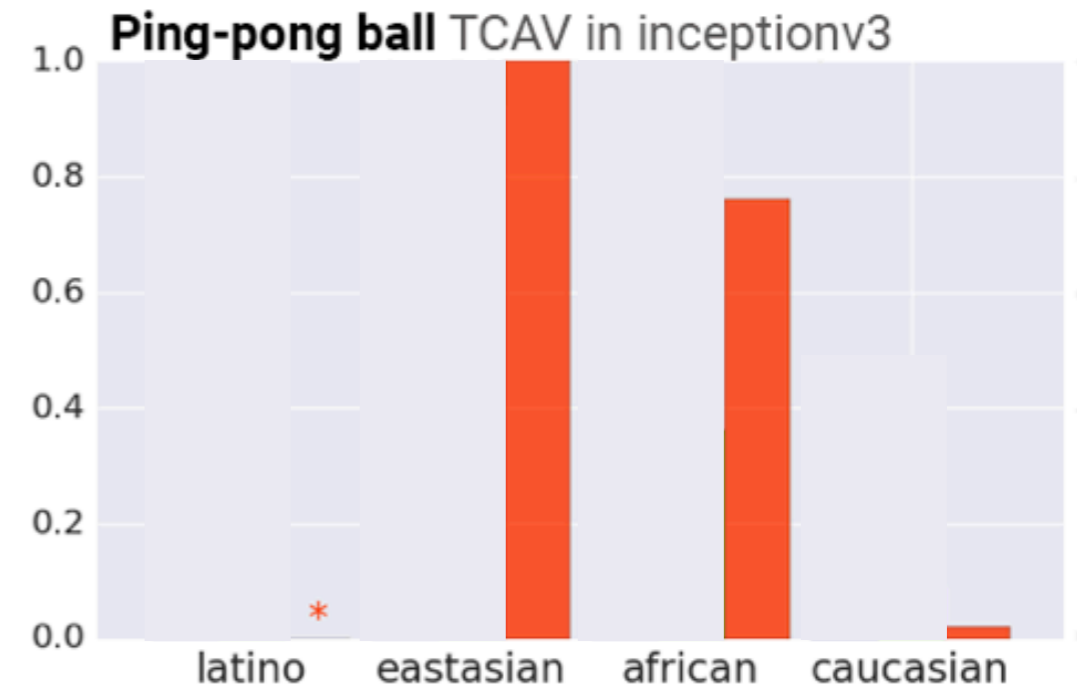
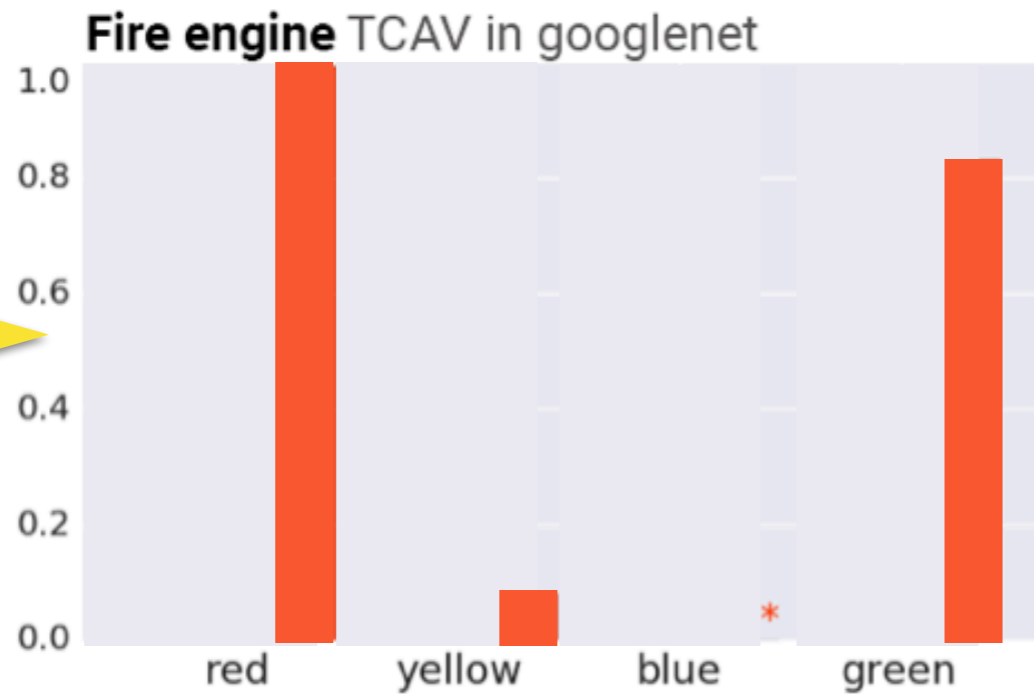


# TCAV in

# Two widely used image prediction models

Geographical bias?

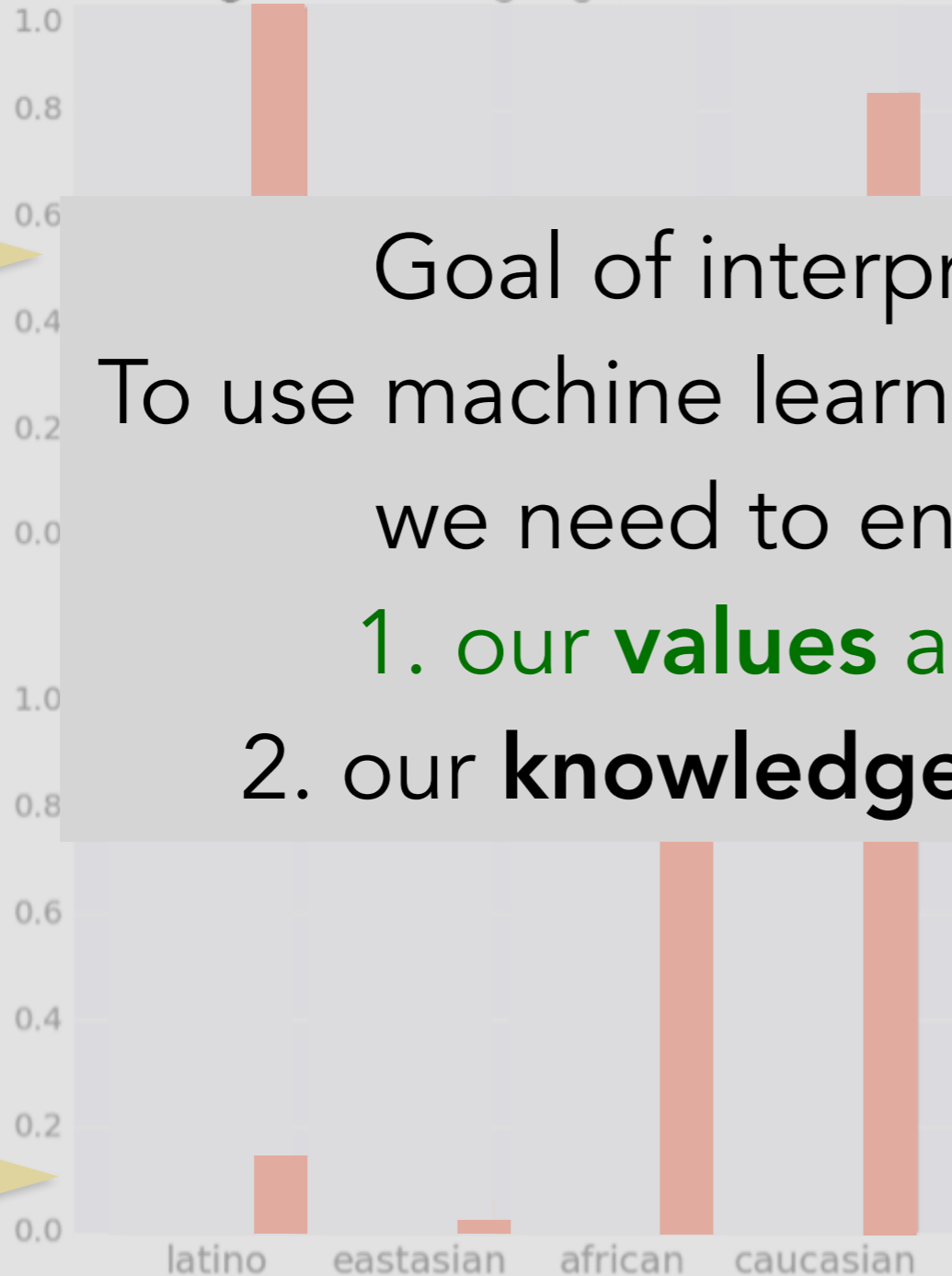
Quantitative confirmation to previously qualitative findings [Stock & Cisse, 2017]



TCAV in

# Two widely used image prediction models

Fire engine TCAV in googlenet



Ping-pong ball TCAV in inceptionv3



Geographical bias?

Quantitative confirmation to previously qualitative findings [Stock & Cisse, 2017]

Goal of interpretability:  
To use machine learning **responsibly**  
we need to ensure that

1. our **values** are aligned
2. our **knowledge** is reflected

# Results

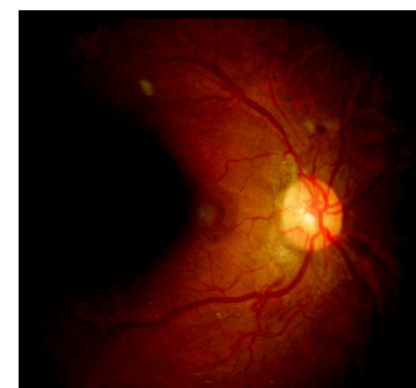
1. Sanity check experiment



2. Biases Inception V3 and GoogleNet

3. Domain expert confirmation from Diabetic Retinopathy

DR level 4 Retina



TCAV for DR level 4





# Diabetic Retinopathy

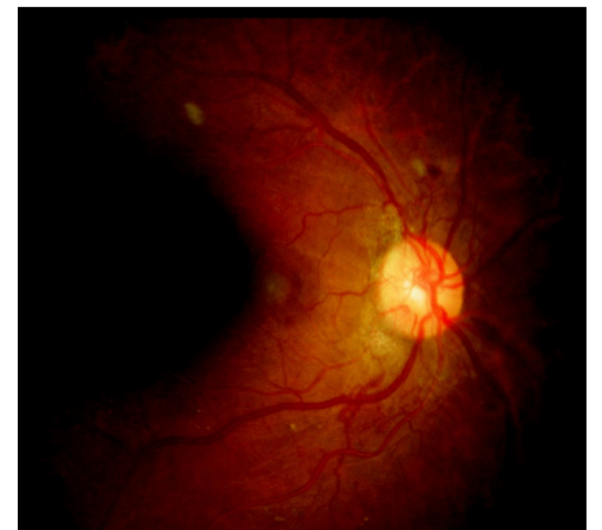
- Treatable but sight-threatening conditions
- Have model to with accurate prediction of DR (85%)  
[Krause et al., 2017]

Concepts the **ML model** uses

Vs

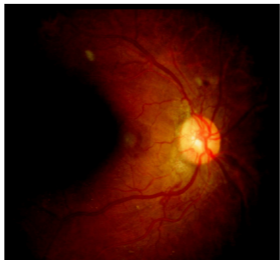

Diagnostic Concepts **human** doctors use

DR level 4 Retina



# Collect human doctor's knowledge



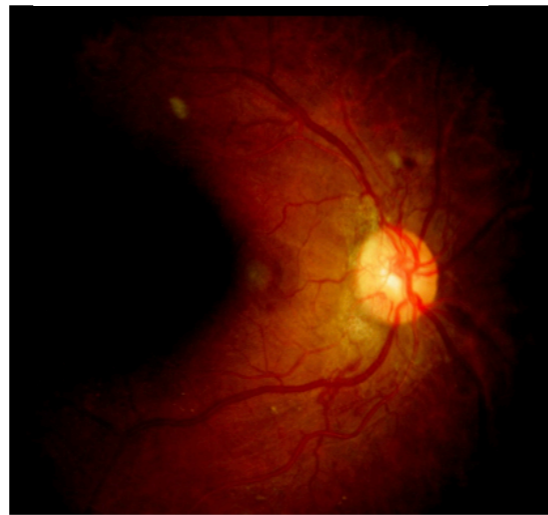
		Concepts belong to this level	Concepts do not belong to this level
DR level 4		<b>PRP</b> <b>PRH/VH</b> <b>NV/FP</b>	<b>VB</b>
DR level 1		<b>MA</b>	<b>HMA</b>

# TCAV for Diabetic Retinopathy

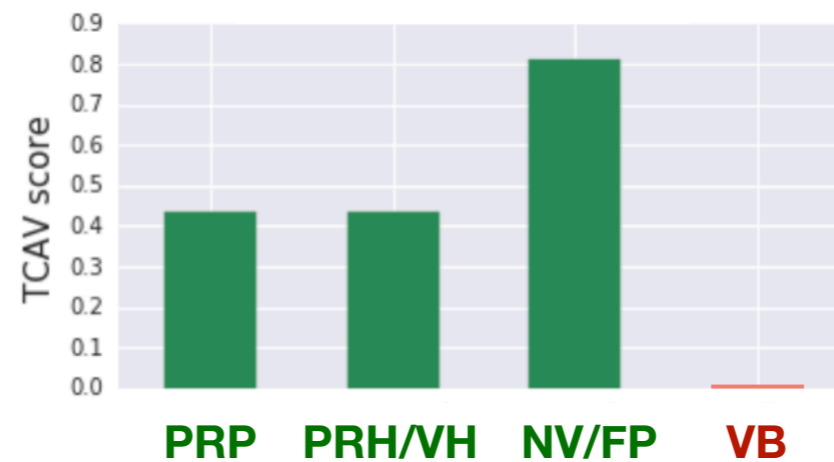
Prediction class    Prediction accuracy

DR level 4    High

Example



TCAV scores

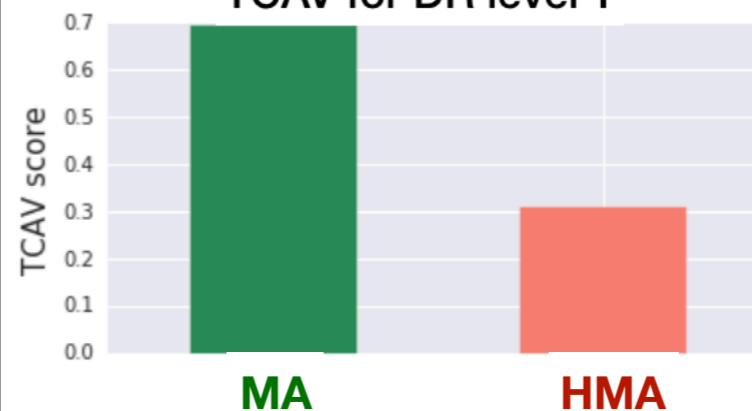


TCAV shows the model is **consistent** with doctor's knowledge when model is **accurate**

DR level 1    Med



TCAV for DR level 1



TCAV shows the model is **inconsistent** with doctor's knowledge for classes when model is less accurate

Green: domain expert's label on concepts belong to the level  
 Red: domain expert's label on concepts does not belong to the level

# TCAV for Diabetic Retinopathy

Prediction class  
Prediction accuracy

Example

Level 1 was often confused to level 2.

HMA distribution on predicted DR

DR level 4

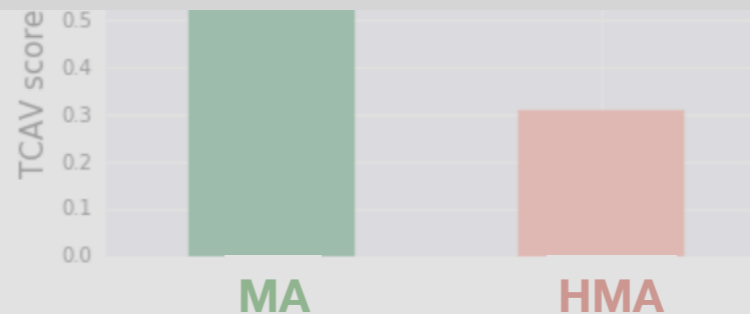
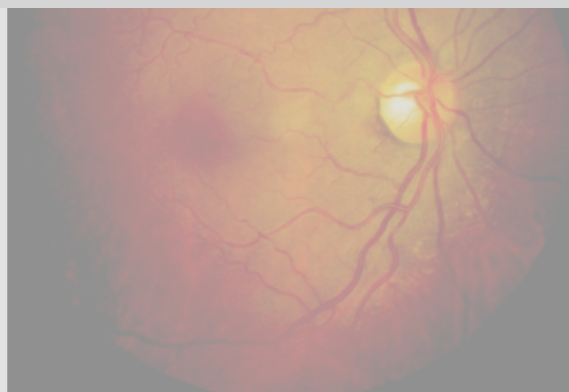
High

Goal of interpretability:  
To use machine learning **responsibly**  
we need to ensure that

1. our **values** are aligned
2. our **knowledge** is reflected

DR level 1

Low



TCAV shows the model is **inconsistent** with doctor's knowledge for classes when model is less accurate

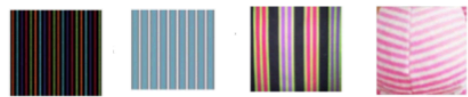
Green: domain expert's label on concepts belong to the level


Red: domain expert's label on concepts does not belong to the level

# Summary:

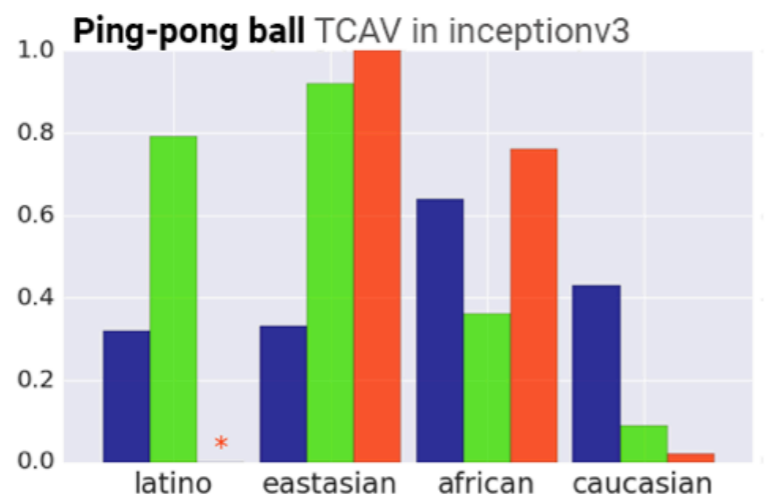
## Testing with Concept Activation Vectors

Joint work with Wattenberg, Gilmer, Cai, Wexler, Viegas, Sayres

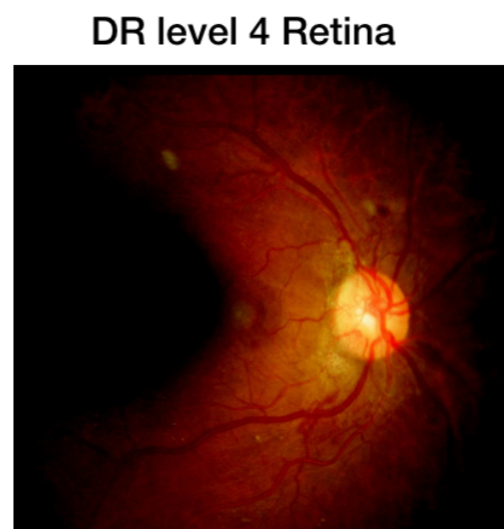


**stripes** concept (score: 0.9)  
was important to **zebra** class  
for this trained network. 

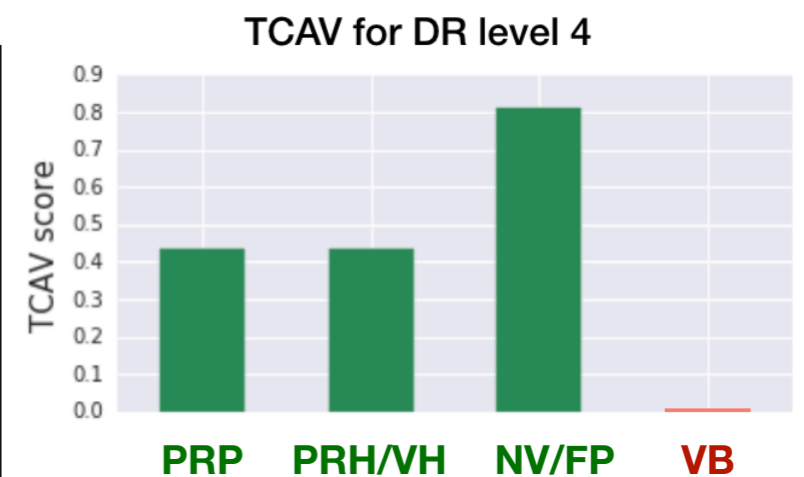
TCAV provides  
**quantitative importance** of  
a concept **if and only if** your  
network learned about it.



Our values



Our knowledge



# Agenda

1. Revisit some existing methods:  
Sanity check questions

2. Make explanations  
that work for lay people.

$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

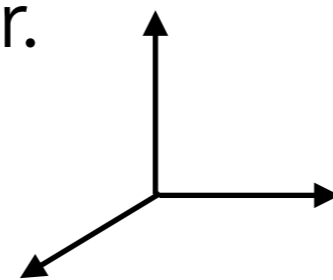
3. Understand how humans  
understand explanations

4. Make explanations to detect  
trustworthy predictions.



# What makes explanations hard or easy for humans?

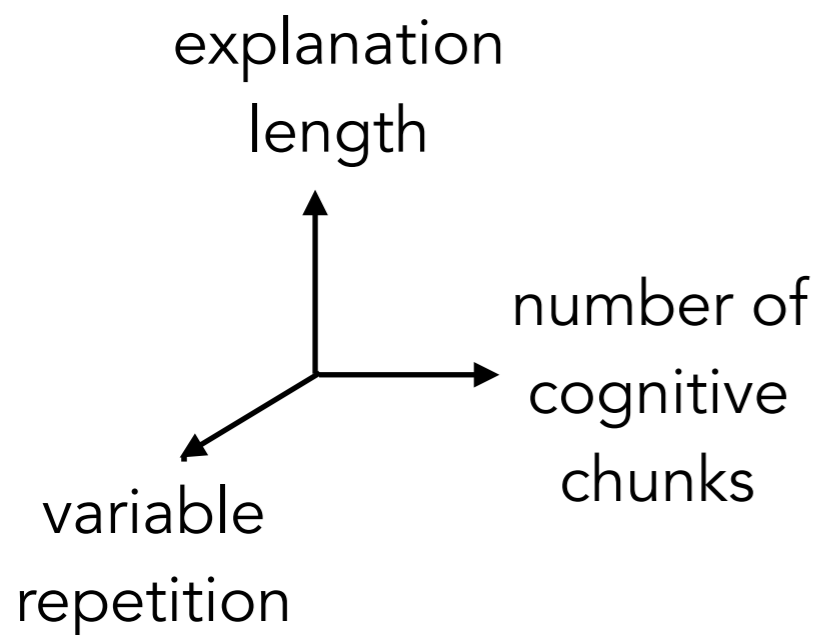
- How do humans' understanding changes as we vary factors in explanations?  $\operatorname{argmax}_E Q(E|M, H, D, T)$
- Among many explanations, we choose a rule-set.
- Among many factors, we choose a subset based on what prior literatures assumed to matter.



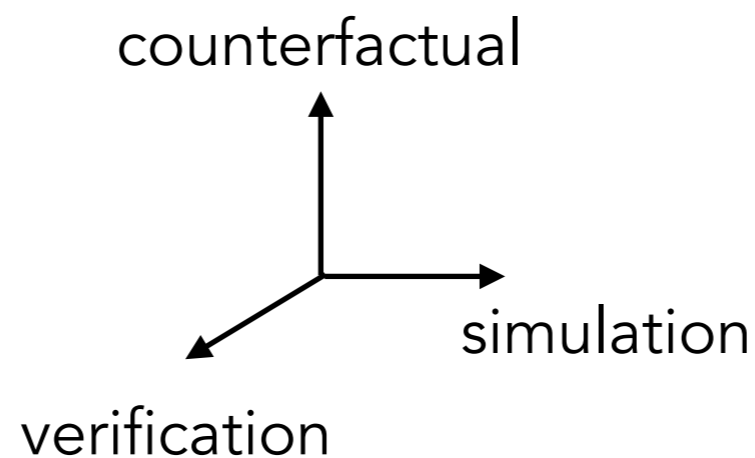
factors

# What makes explanations hard or easy for humans?

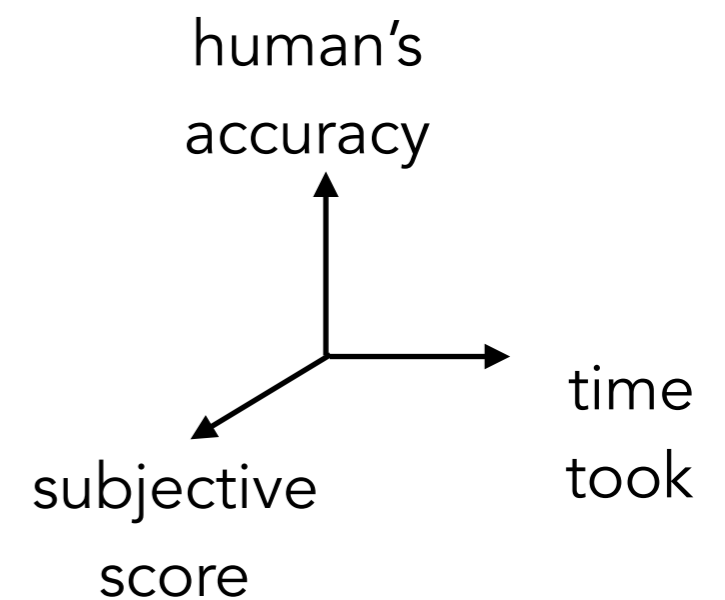
- How do humans' understanding changes as we vary factors in explanations?  $\operatorname{argmax}_E Q(E|M, H, D, T)$



we vary these factors



humans do these tasks



measures of interpretability

# Controlling for prior knowledge.

Using a made-up 'alien world' to control for prior knowledge.

## The alien's preferences:

checking the news and coughing → windy  
snowing or humid and weekend → **spices or vegetables** and **grains**  
embarrassed and grouchy or raining → **dairy** or **vegetables**  
snowing or windy and energetic → **candy or dairy** and **fruit**  
grouchy or weekend and windy → **spices or grains** and **fruit**



Is the alien happy with his meal?

Yes  No

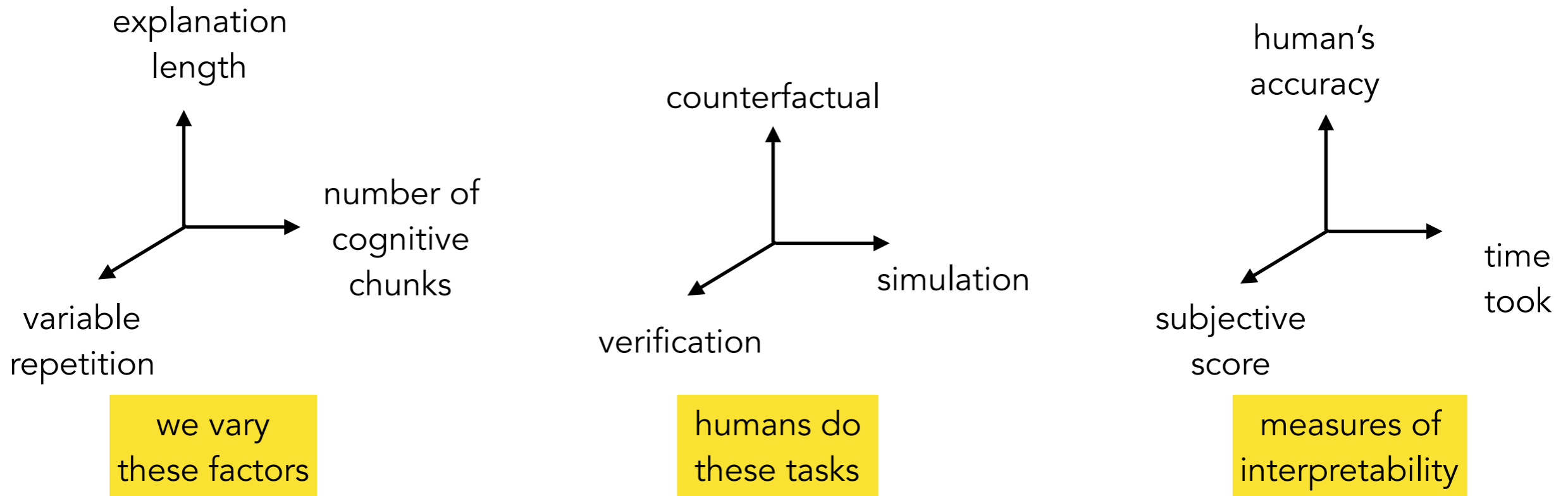
Observations: **Saturday, coughing, checking the news**

Recommendation: **bagel, rice, strawberry**

## Ingredients:

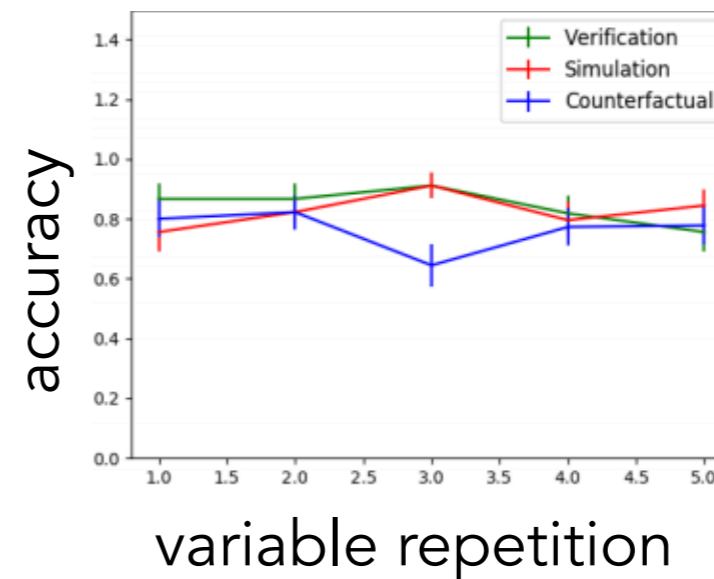
- **Vegetables:** okra, carrot, spinach
- **Spices:** turmeric, thyme, cinnamon
- **Dairy:** milk, butter, yogurt
- **Fruit:** mango, strawberry, guava
- **Candy:** chocolate, taffy, caramel
- **Grains:** bagel, rice, pasta

# (a small subset of) results



Variable repetition mattered less for accuracy than other factors.

\*all repeated variables are needed for task completion



# Agenda

1. Revisit some existing methods:  
Sanity check questions

2. Make explanations  
that work for lay people.

$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

3. Understand how humans  
understand explanations

4. Make explanations to detect  
trustworthy predictions.

# Ultimate goal is to use ML more responsibly.

Goal of interpretability:

To use machine learning **responsibly**

we need to ensure that

1. our **values** are aligned
2. our **knowledge** is reflected

Not using the classifier when it's suspicious.



Improve confidence measure coming from a classifier

$$\operatorname{argmax}_E Q(E|M, H, D, T)$$

Simply confidence scores

**Problem:**  
precision of  
"definitely trustworthy  
(correct)" predictions

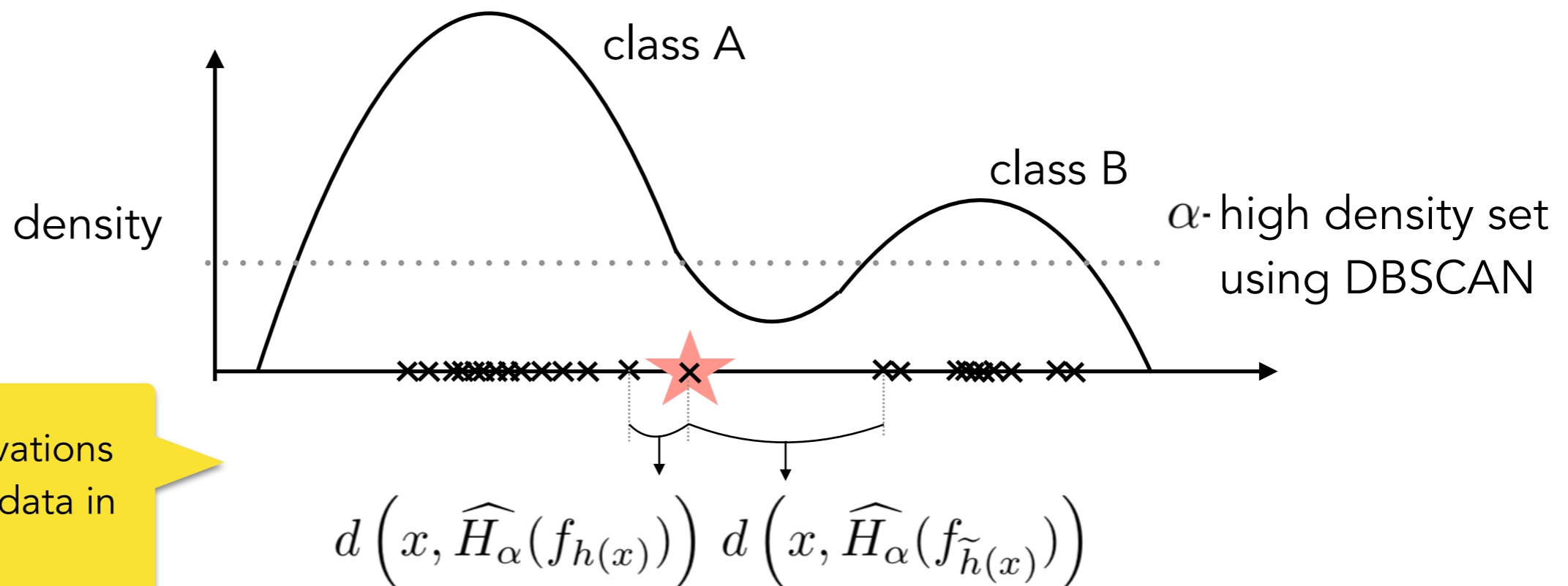
"definitely suspicious  
(incorrect)"  
predictions



# Trust score: a super simple method

★ was predicted as class A. Can we trust this?

$$\text{Trust score} := \frac{d\left(x, \widehat{H}_\alpha(f_{\tilde{h}(x)})\right)}{d\left(x, \widehat{H}_\alpha(f_{h(x)})\right)}$$



We can use activations instead of input data in NN!

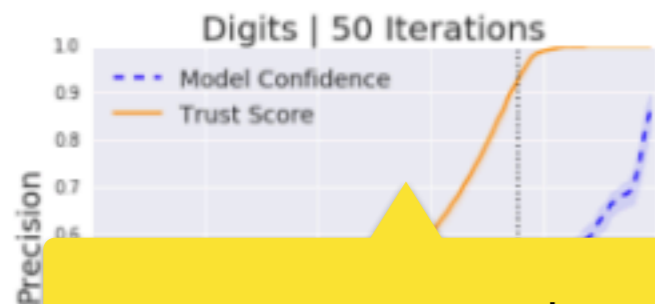
results:

We can detect trustworthy and suspicious predictions with high precision.

Detect trustworthy (correct)



Detect suspicious (incorrect)



Theoretical results: why does this work?

The trust score reveals the signal from a Bayes optimal classifier (with high probability).

# Summary, future work

1. Revisit some existing methods:

Sanity check questions

Sanity Checks for Saliency Maps

Joint work with Adebayo, Gilmer, Goodfellow, Hardt  
NIPS 2018

2. Make explanations that work for lay people.

TCAV: Testing with concept activation vectors

Joint work with Wattenberg, Gilmer, Cai, Wexler, Viegas, Sayres  
ICML 2018

$$\operatorname{argmax}_E Q(\mathbf{Explanation} | \mathbf{Model}, \mathbf{Human}, \mathbf{Data}, \mathbf{Task})$$

3. Understand how humans understand explanations

How do Humans Understand Explanations from Machine Learning Systems?

An Evaluation of the Human-Interpretability of Explanation

joint work with Narayanan, Chen, He, Gershman, and Doshi-Velez 2017

4. Make explanations to detect trustworthy predictions.

To trust or not to trust a classifier

joint work with Jiang and Gupta

NIPS 2018

Understanding superhuman performance networks

Understanding models under production @ Google

Detect 'different types of mistakes' that a model makes.

...lots of others.