



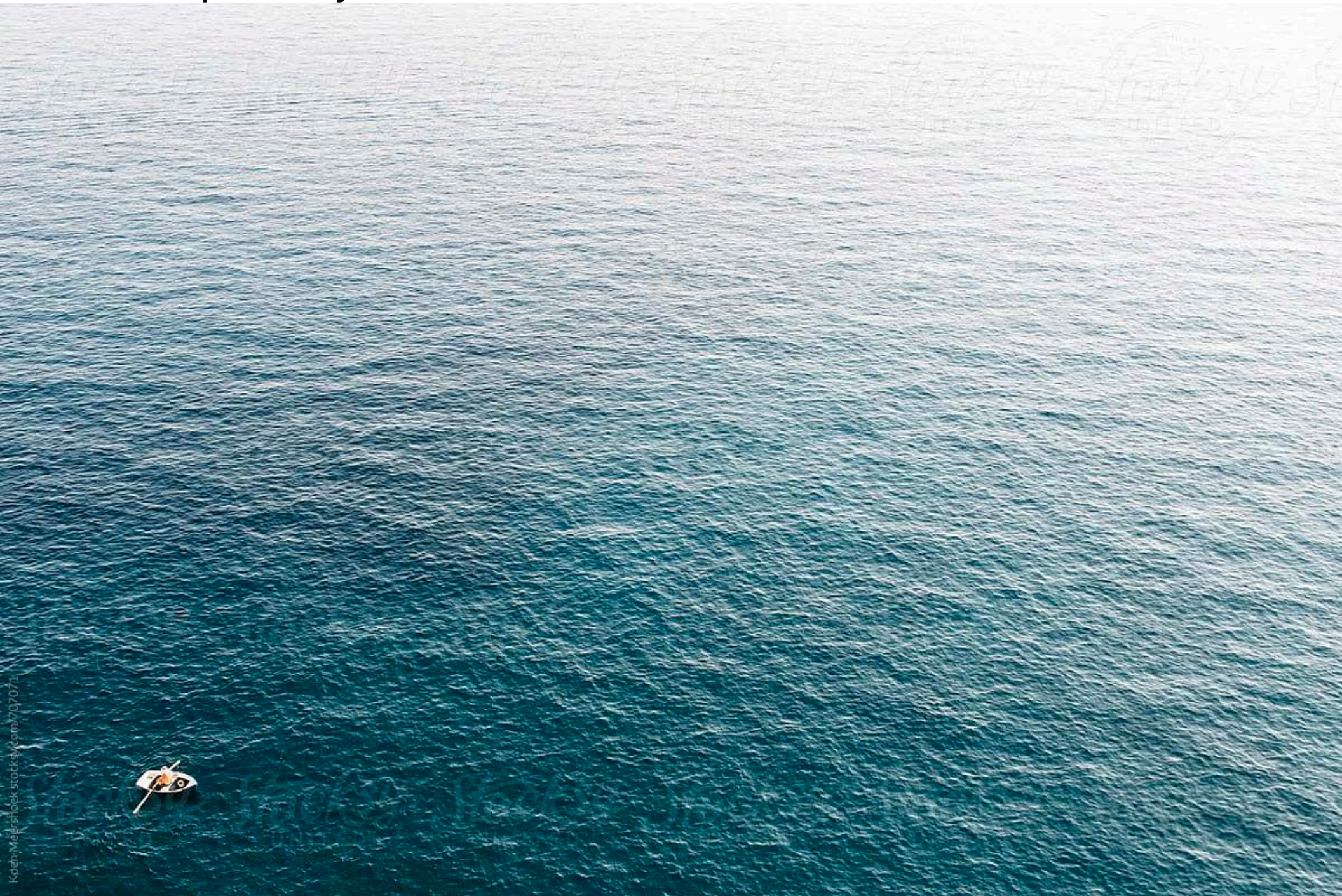
Interpretability: what now?

Been Kim

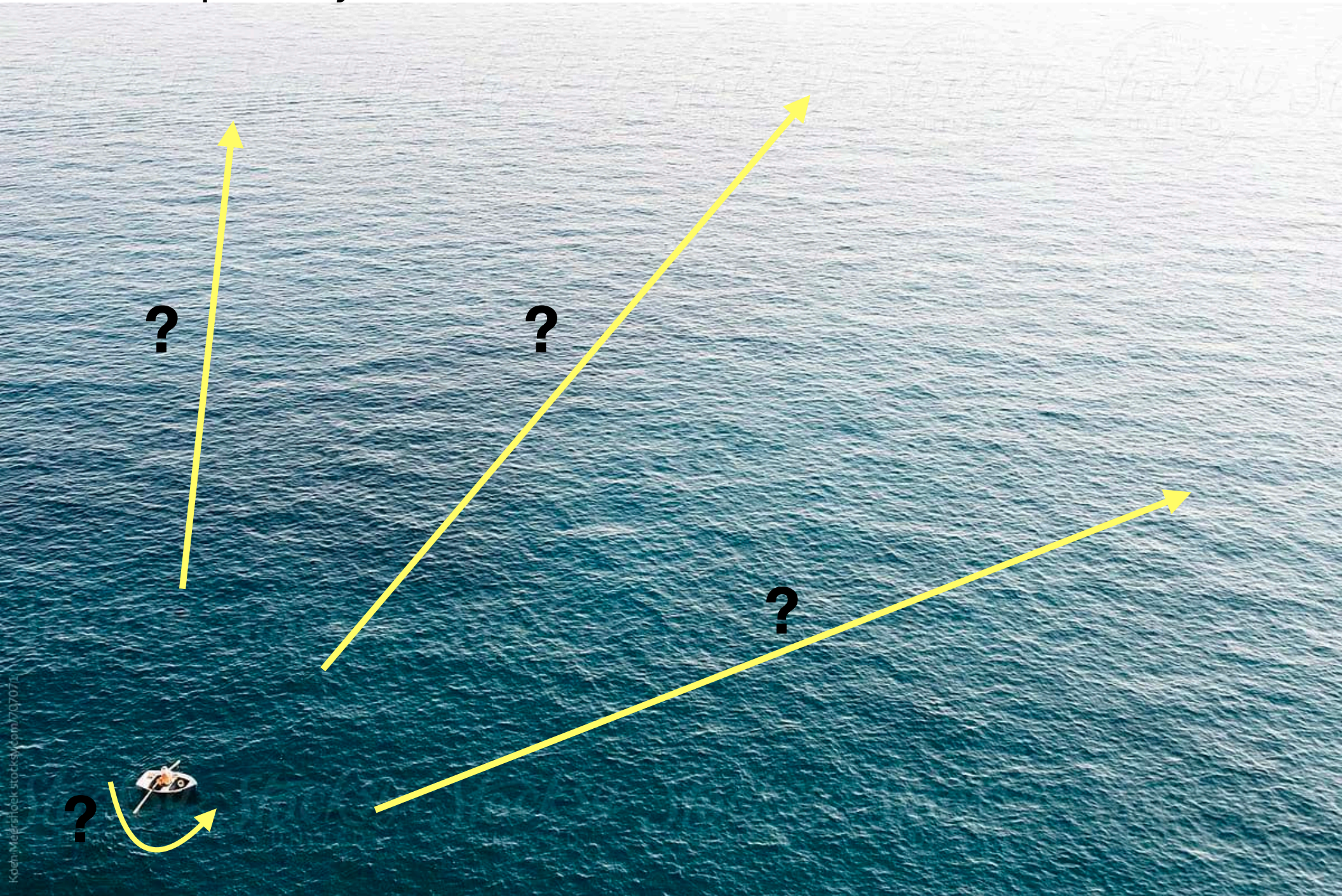
Presenting work with a lot of awesome people inside and outside of Google

Julius Adebayo, Sherry Yang, Justin Gilmer, Martin Wattenberg, Carrie Cai, James Wexler,
Fernanda Viegas, Rory Sayres, Ian Goodfellow, Mortiz Hardt, Michael Muelly

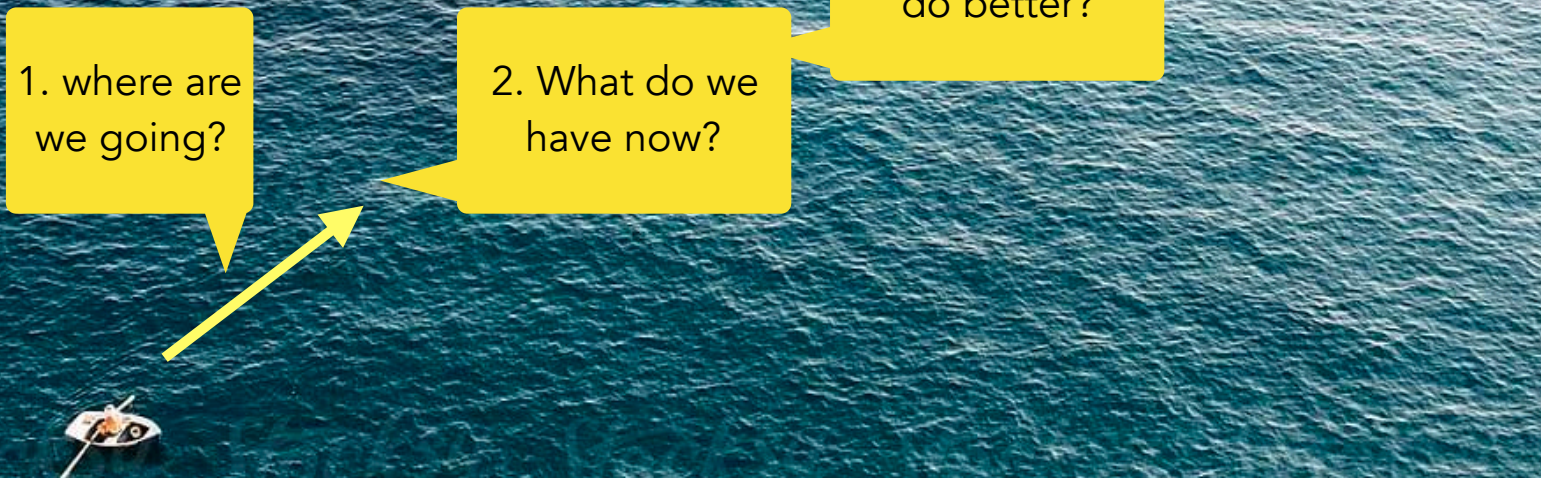
Sea of interpretability



Sea of interpretability



Sea of interpretability



Sea of interpretability



1. where are
we going?



My goal

interpretability

To use machine learning **responsibly**

we need to ensure that

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

My goal

interpretability

To use machine learning **responsibly**

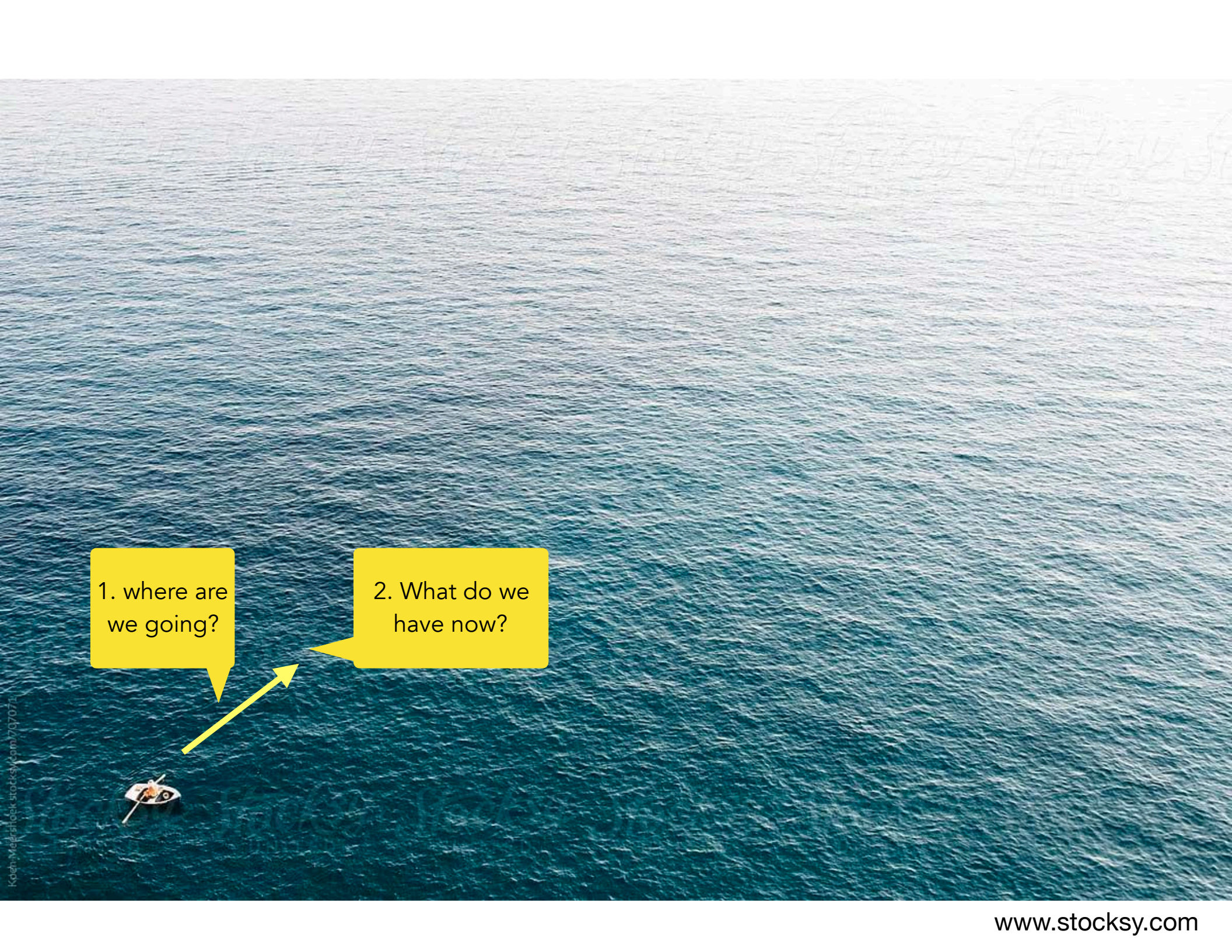
we need to ensure that

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

NON-goals

Interpretability is NOT...

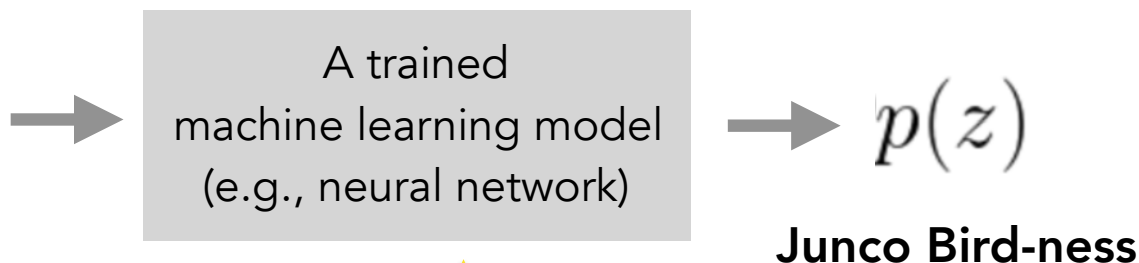
- about making ALL models interpretable.
- about understanding EVERY SINGLE BIT about the model
- against developing highly complex models.
- only about gaining user trust or fairness



1. where are
we going?

2. What do we
have now?

Investigating post-training interpretability methods.

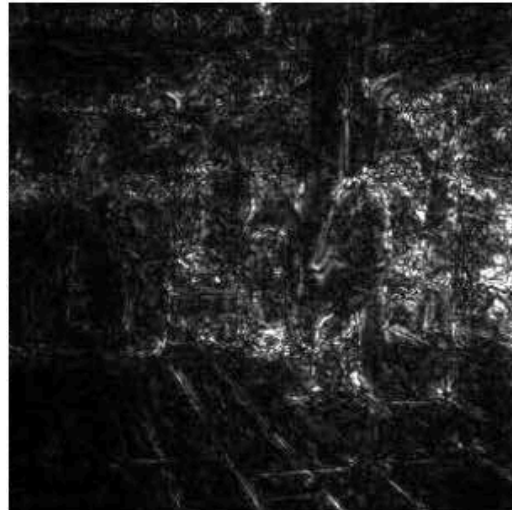


Given a fixed model, find the **evidence** of **prediction**.

Why was this a Junco bird?

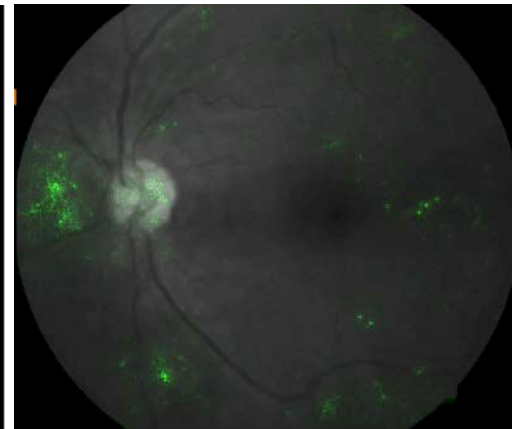
One of the most popular interpretability methods for images:

Saliency maps



Caaaaan do! We've got saliency maps to measure importance of each pixel!

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



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:

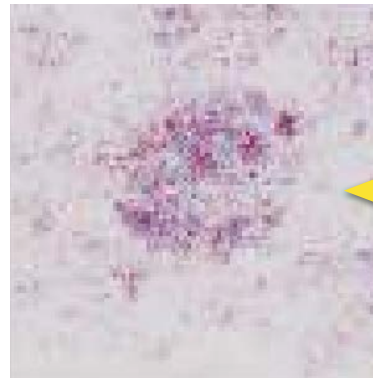
Saliency maps



A trained
machine learning model
(e.g., neural network)

→ $p(z)$

Junco Bird-ness



The promise:
these pixels are the
evidence of
prediction.

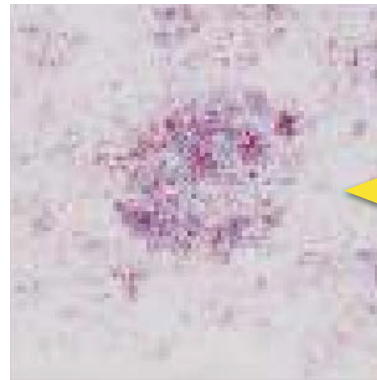
Sanity check question.



A trained
machine learning model
(e.g., neural network)

→ $p(z)$

Junco Bird-ness



The promise:
these pixels are the
**evidence of
prediction.**

Sanity check question.



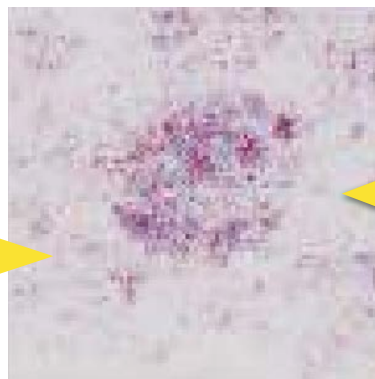
A trained
machine learning model
(e.g., neural network)

→ $p(z)$

Junco Bird-ness

If so, when **prediction** changes,
the explanation should change.

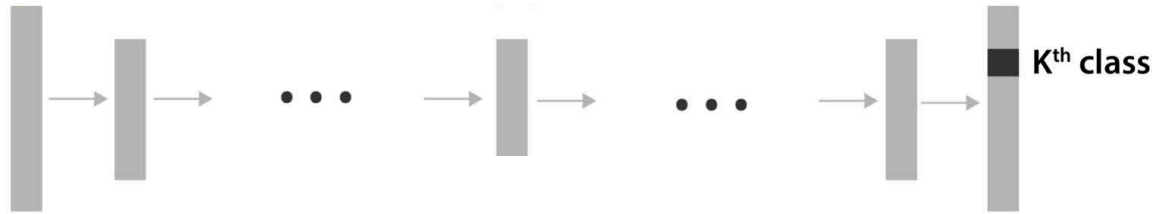
Extreme case:
If **prediction** is random,
the **explanation** should
REALLY change.



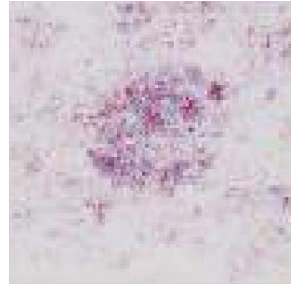
The promise:
these pixels are the
evidence of
prediction.

Some confusing behaviors of saliency maps.

Original Image

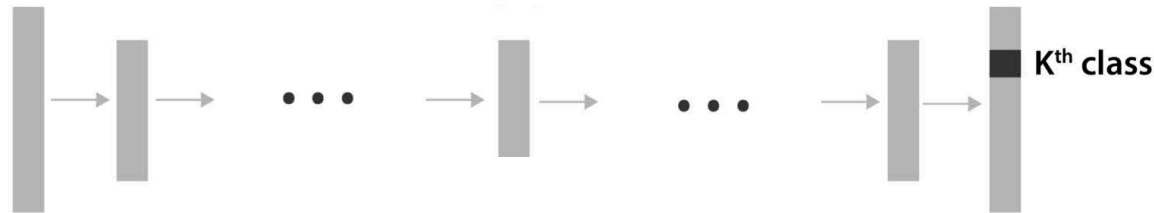


Saliency map

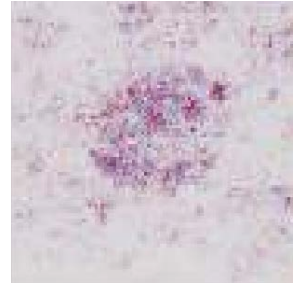


Some confusing behaviors of saliency maps.

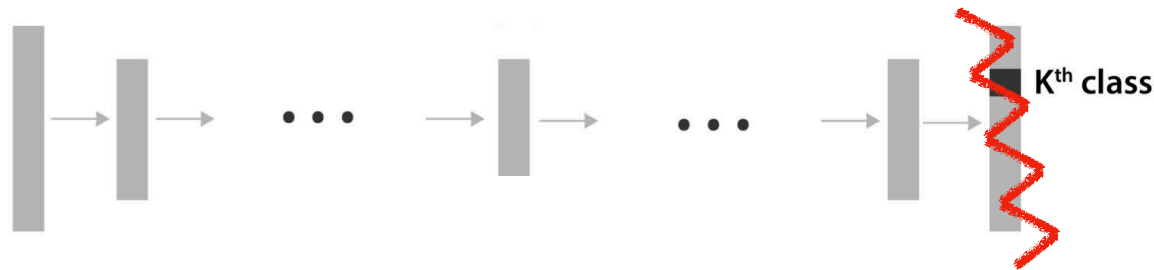
Original Image



Saliency map

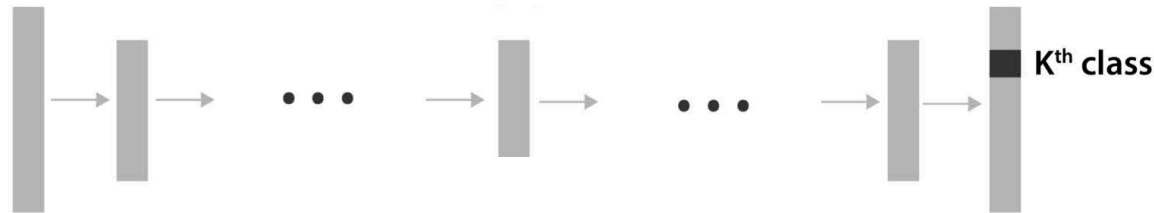


Randomized weights!
Network now makes garbage prediction.

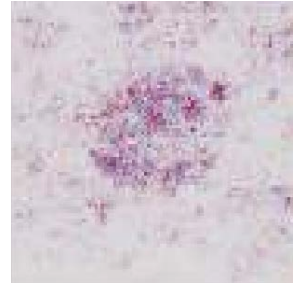


Some confusing behaviors of saliency maps.

Original Image



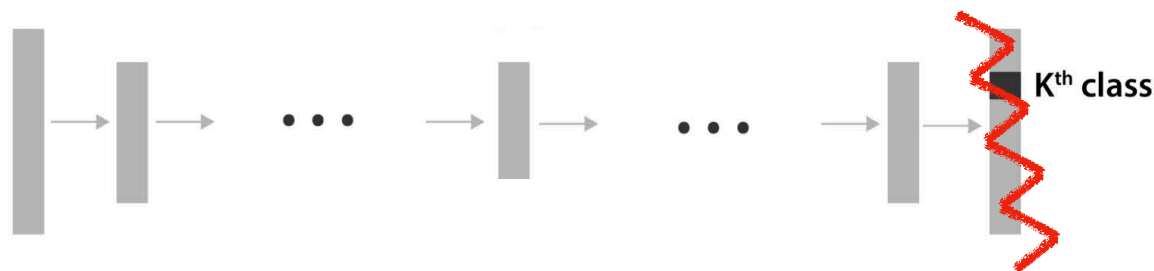
Saliency map



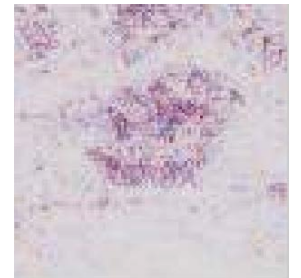
Original Image



Randomized weights!
Network now makes garbage prediction.



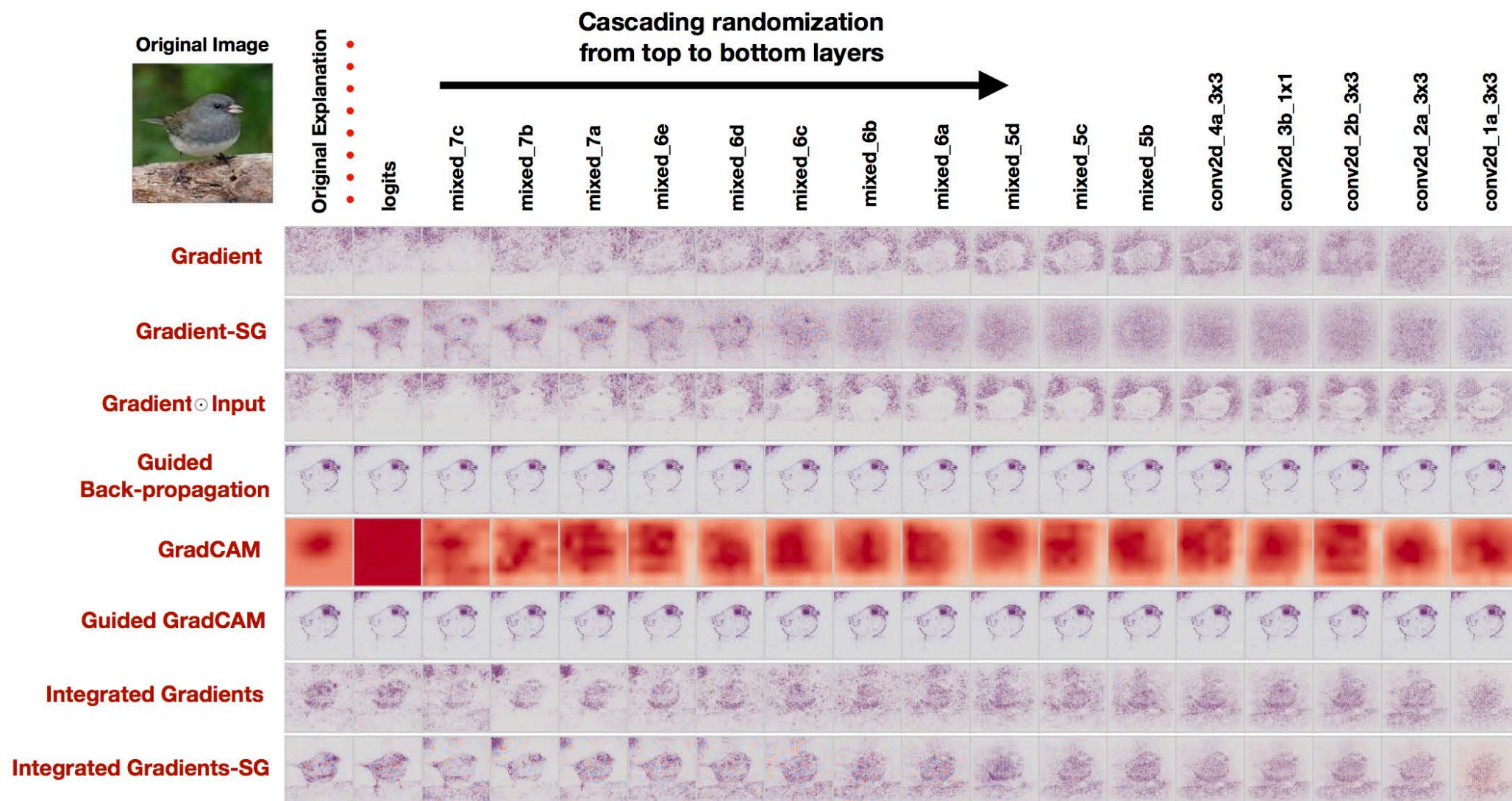
!!!!!!????!?



Sanity check1:

When prediction changes, do explanations change?

No!

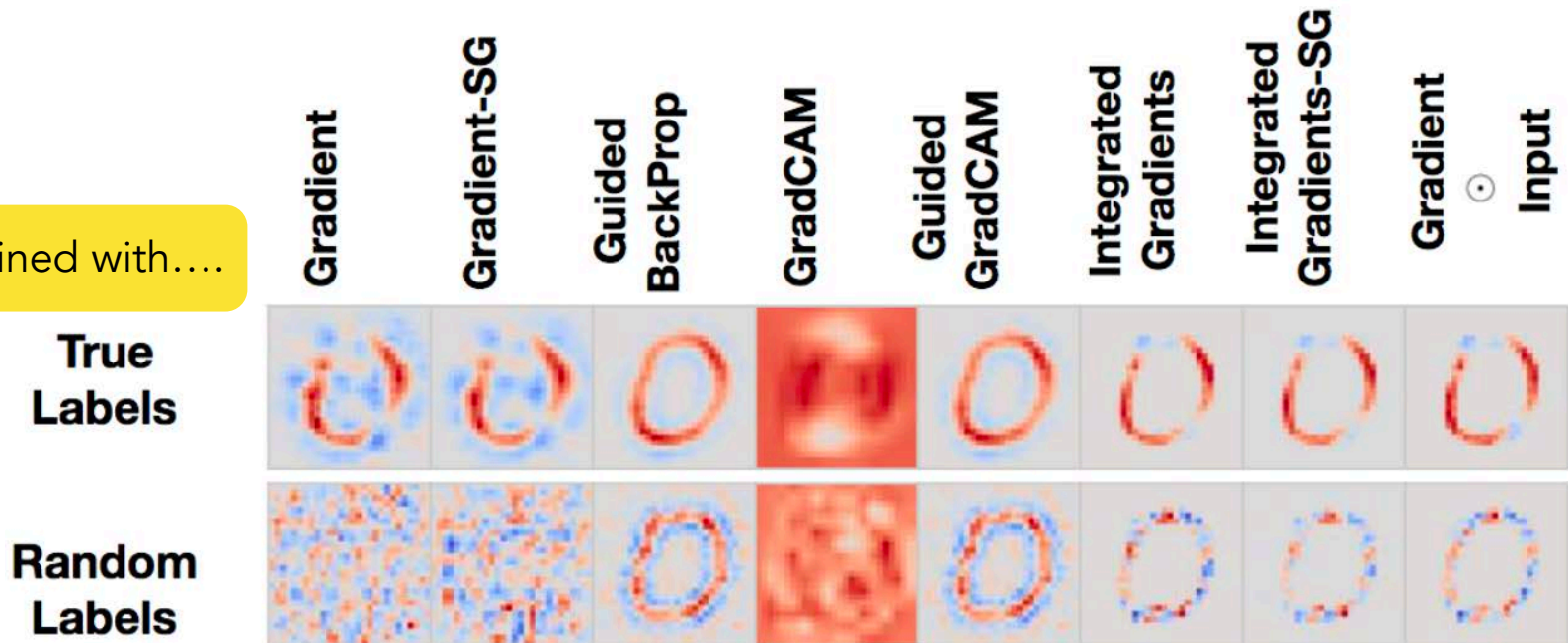


Sanity check2:

Networks trained with true and random labels,
Do explanations deliver different messages?

No!

Networks trained with....



What can we learn from this?

- **Confirmation bias:** Just because it “makes sense” to humans, doesn’t mean it reflects the evidence for prediction.
- Others who independently reached the same conclusions:
[Nie, Zhang, Patel '18] [Ulyanov, Vedaldi, Lempitsky '18]
- Some of these methods have been shown to be useful for humans. Why? More studies needed.



This was a low bar test.

Can we put interpretability methods
on a harder test?

Benchmarking interpretability methods (BIM)

Forest



A thing



Benchmarking interpretability methods (BIM)

Forest



Forest



A thing



Bedroom



Kitchen



Benchmarking interpretability methods (BIM)

Forest



Forest



A thing





Bedroom



Kitchen



 is NOT important for predicting scene classes.
↓
 should NOT Be part of explanation

Benchmarking interpretability methods (BIM)

Forest



Forest



A thing



Bedroom



Kitchen



is NOT important for predicting scene classes.



should NOT
Be part of explanation

We can also make



more important
to some classes by
controlling when it appears.



should be more
important explanation in
some classes than others.

Benchmarking interpretability methods (BIM)

Forest



Forest



A thing





Bedroom



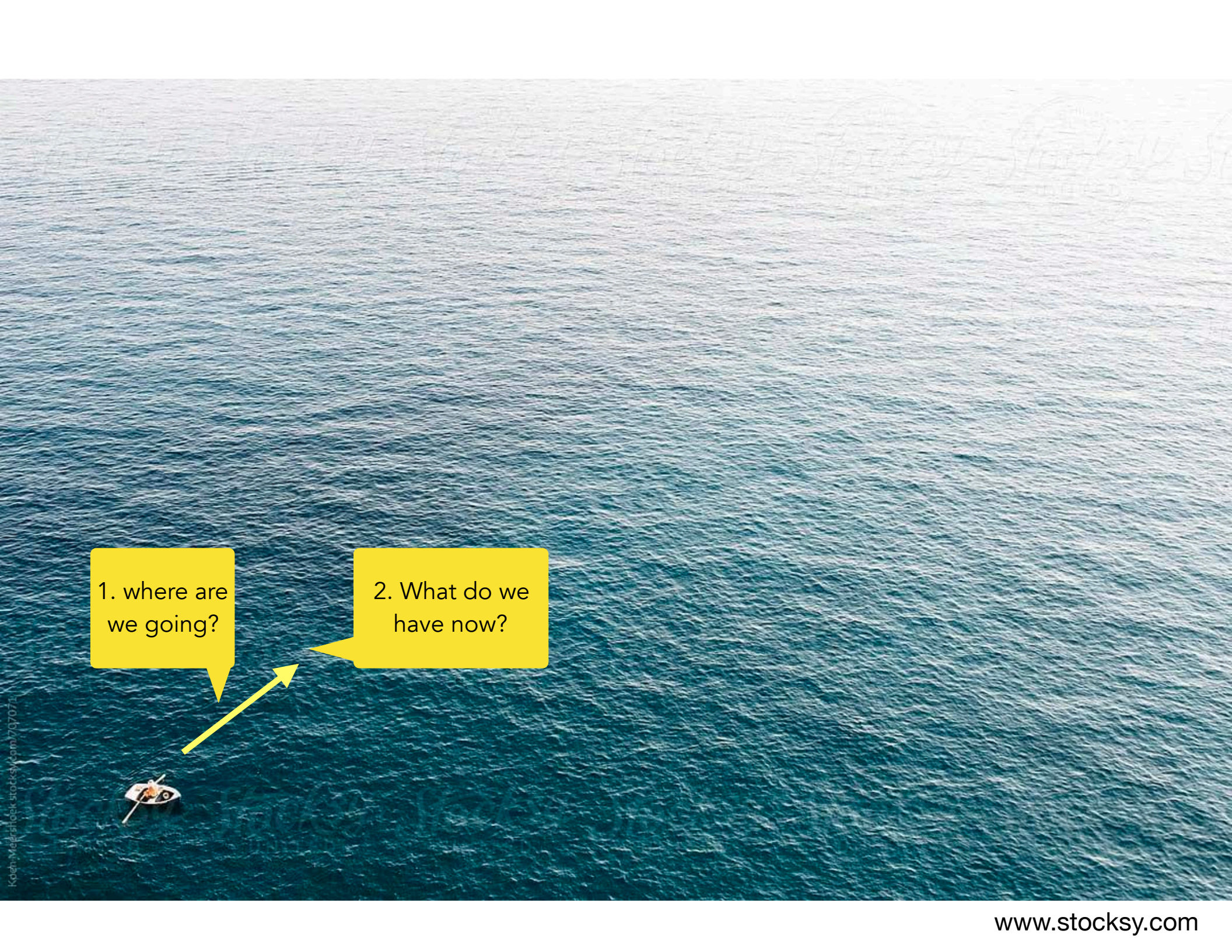
coming soon!
data, model and metrics
for existing
interpretability methods

Kitchen



 is NOT important for predicting scene classes.
↓
 should NOT Be part of explanation

github.com/google-research-datasets/bim
work with Sherry Yang



1. where are we going?

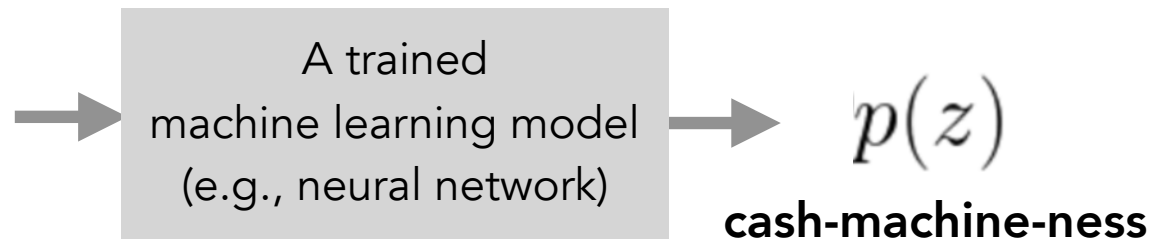
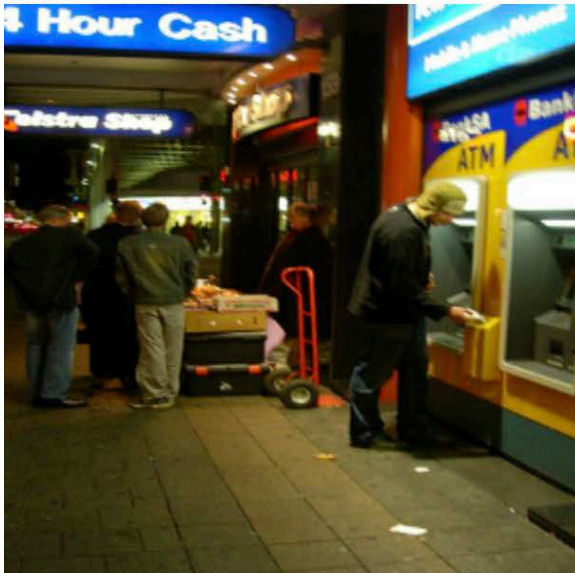
2. What do we have now?

3. What can we do better?

Problem:

Post-training explanation

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



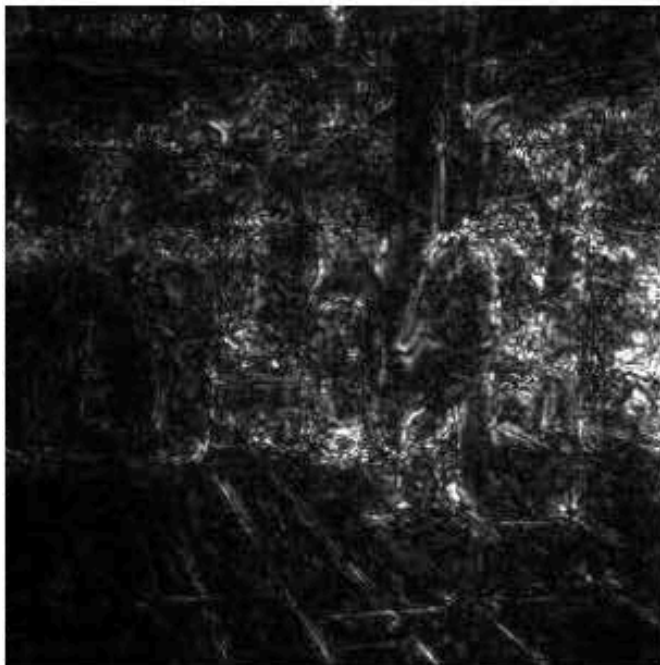
Why was this a cash machine?

Common solution: Saliency map

prediction:
Cash machine



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



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

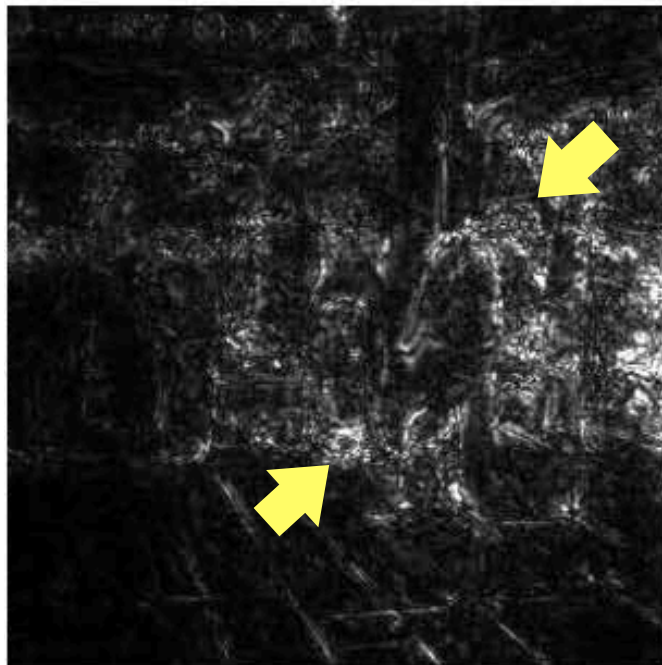
What we really want to ask...

prediction:
Cash machine



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

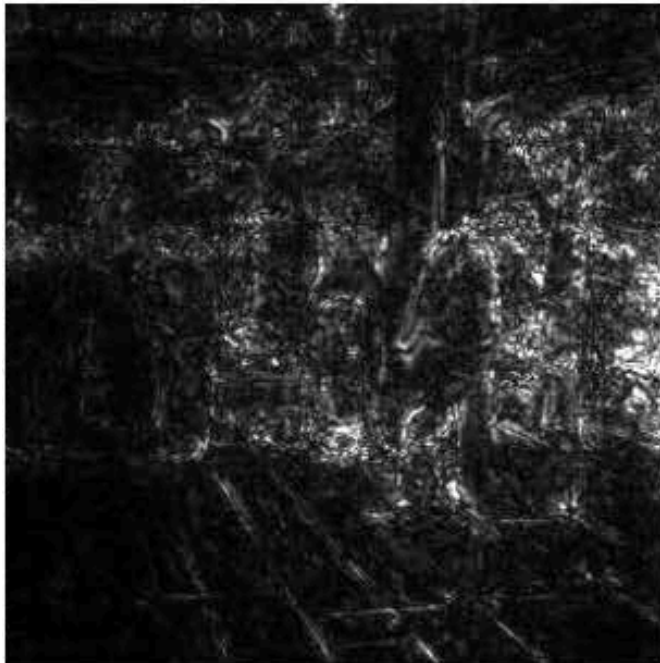
Did the 'human' concept matter?
Did the 'wheels' concept matter?



<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?

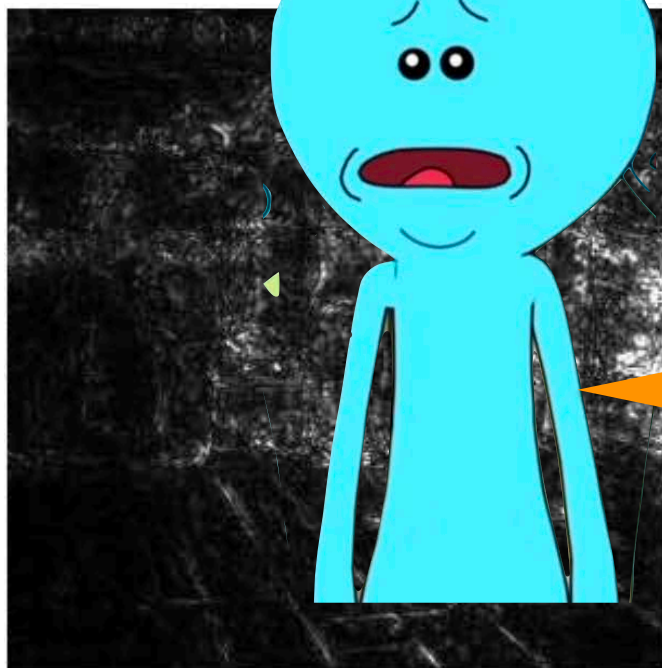
Which concept mattered more?

Is this true for all other cash machine predictions?

<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?

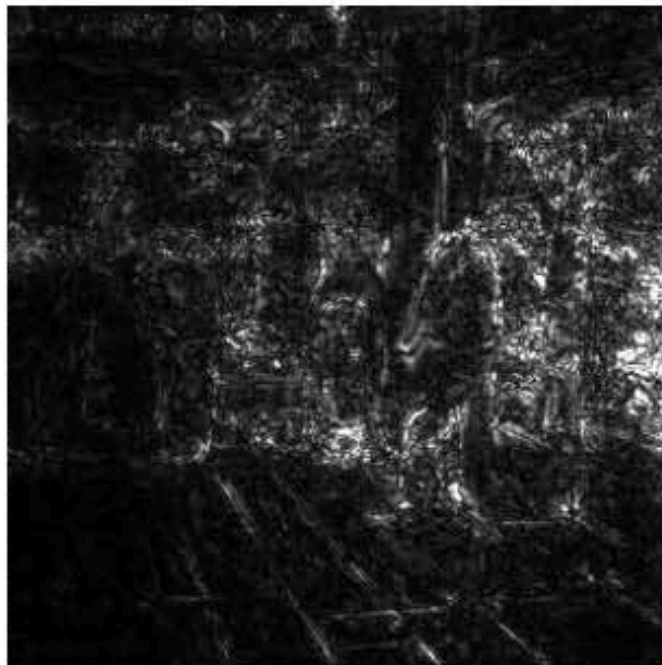
Is this true for all other cash machine predictions?

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?

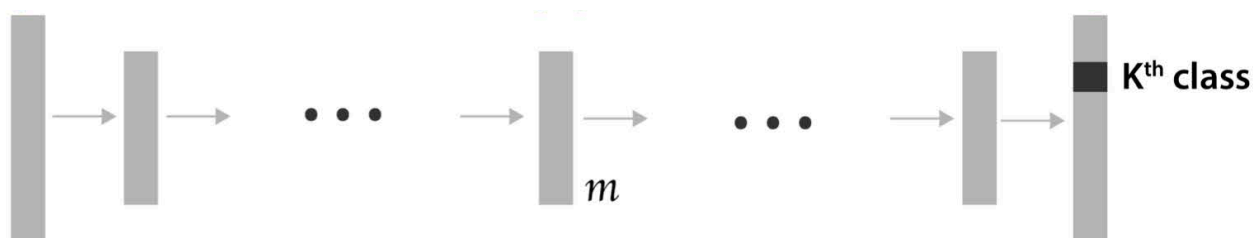
Which concept mattered more?

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/>

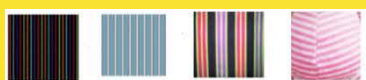
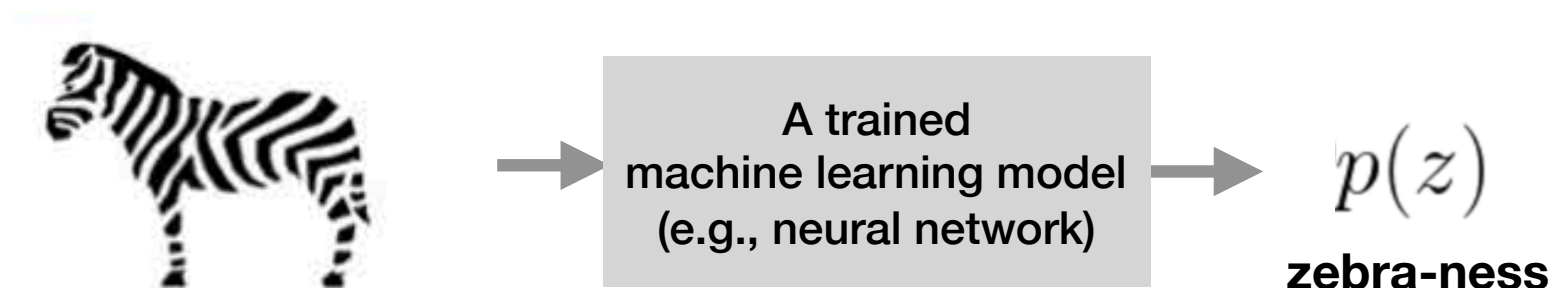
Goal of TCAV: Testing with Concept Activation Vectors



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

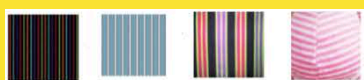
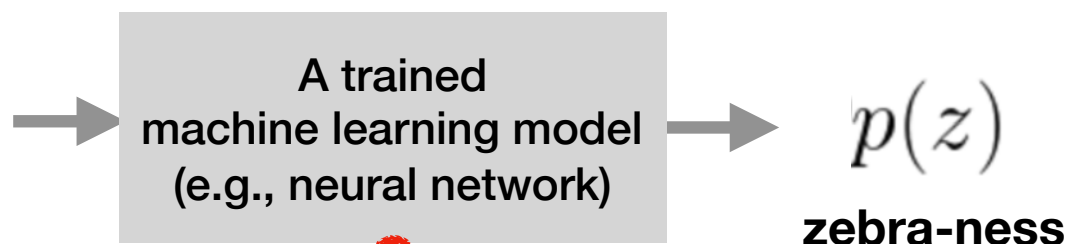
...even if the **concept** was not part of the training.

Goal of TCAV: Testing with Concept Activation Vectors



Was striped concept important
to this zebra image classifier?

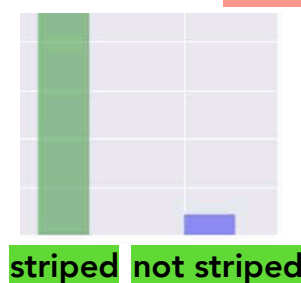
Goal of TCAV: Testing with Concept Activation Vectors



Was **striped** concept important to this **zebra** image classifier?



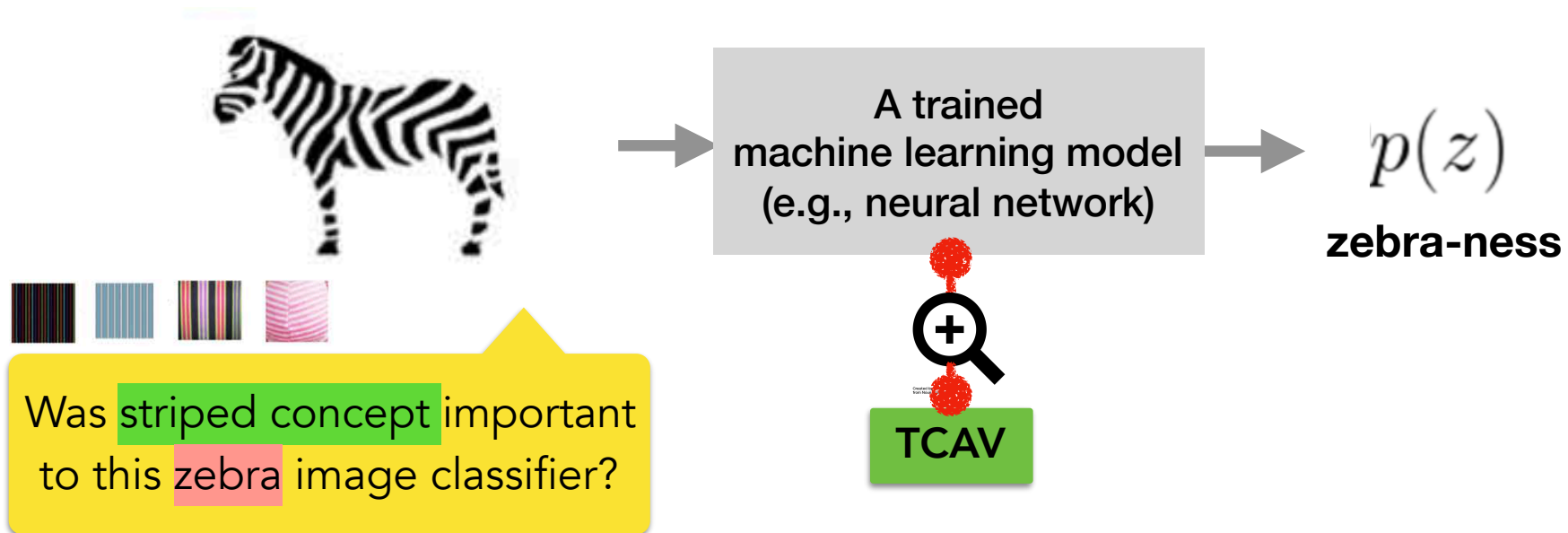
TCAV score for **Zebra**



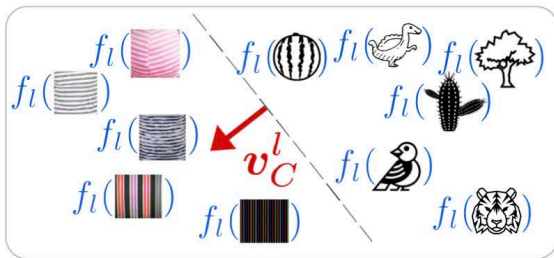
TCAV provides **quantitative importance** of a concept **if and only if** your network learned about it.

TCAV:

Testing with Concept Activation Vectors



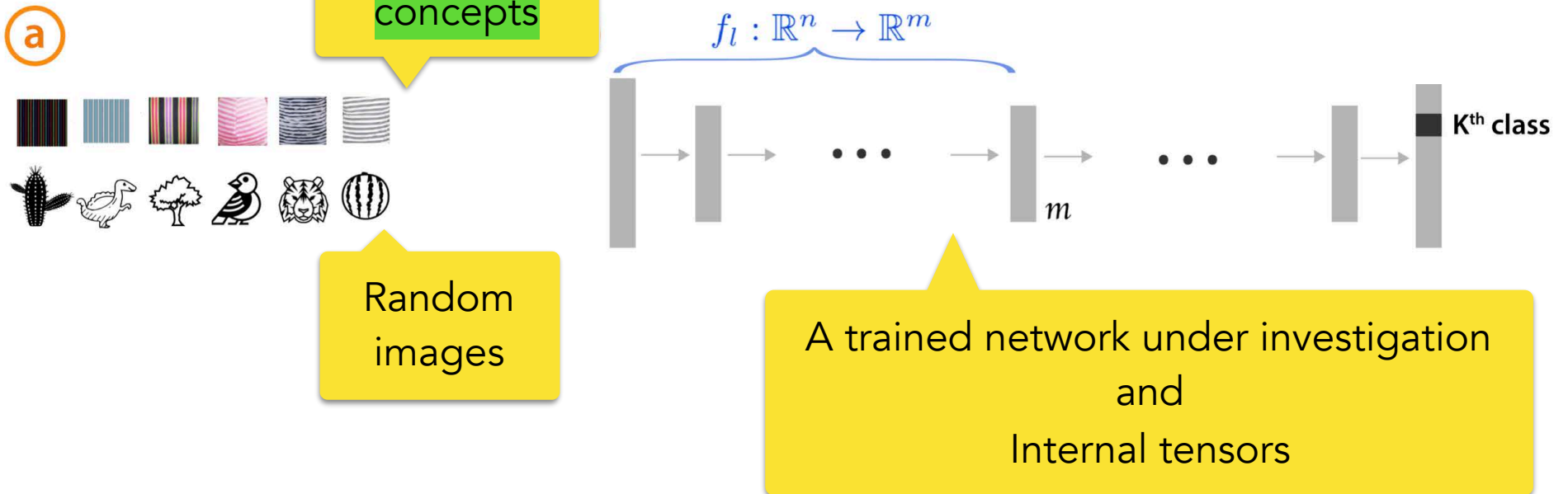
1. Learning CAVs



1. How to define concepts?

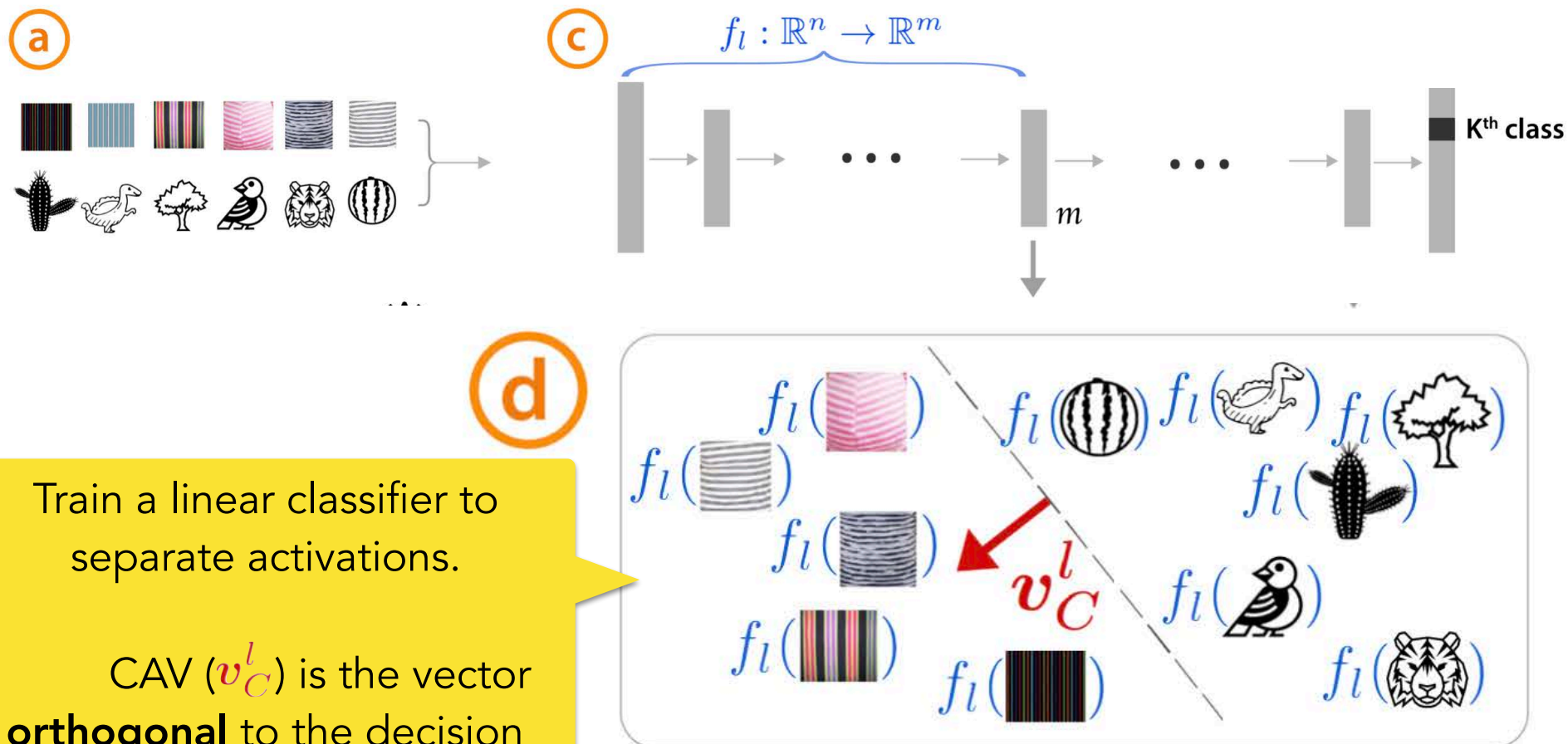
Defining concept activation vector (CAV)

Inputs:



Defining concept activation vector (CAV)

Inputs:



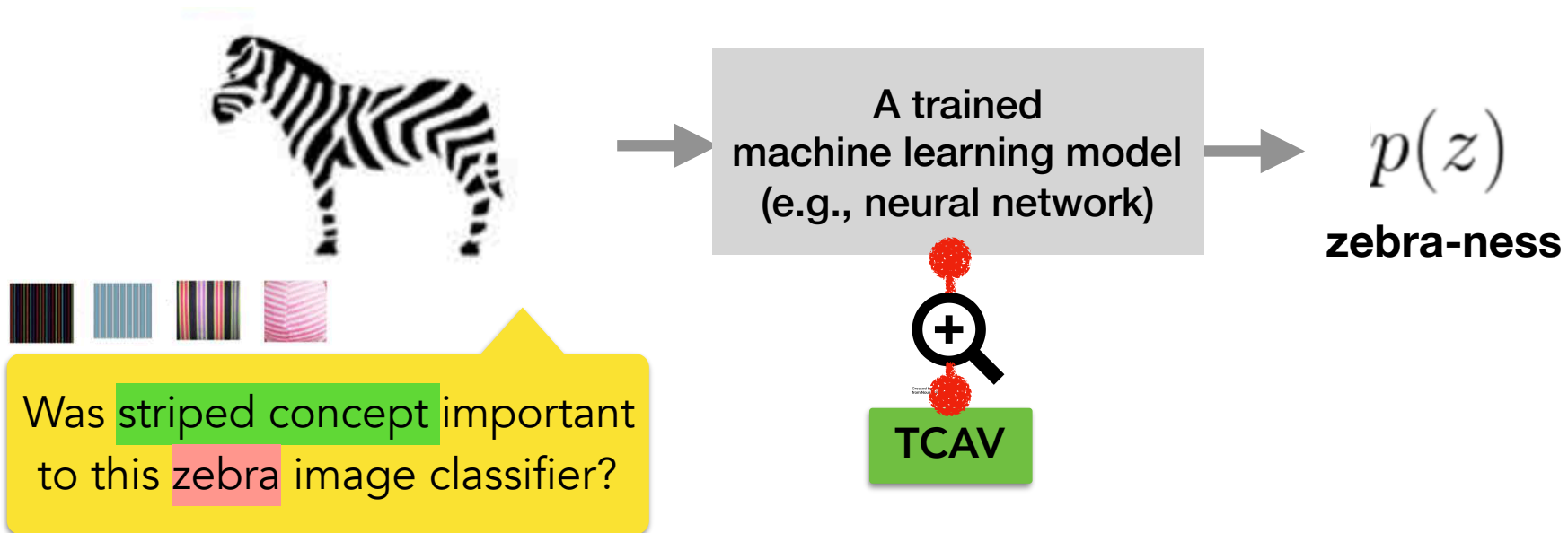
Train a linear classifier to separate activations.

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

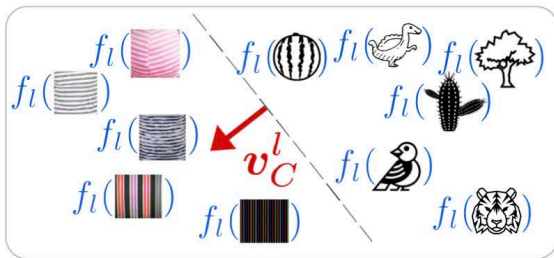
[Smilkov '17, Bolukbasi '16, Schmidt '15]

TCAV:

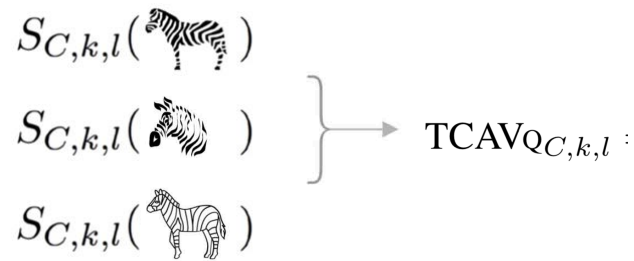
Testing with Concept Activation Vectors



1. Learning CAVs



2. Getting TCAV score

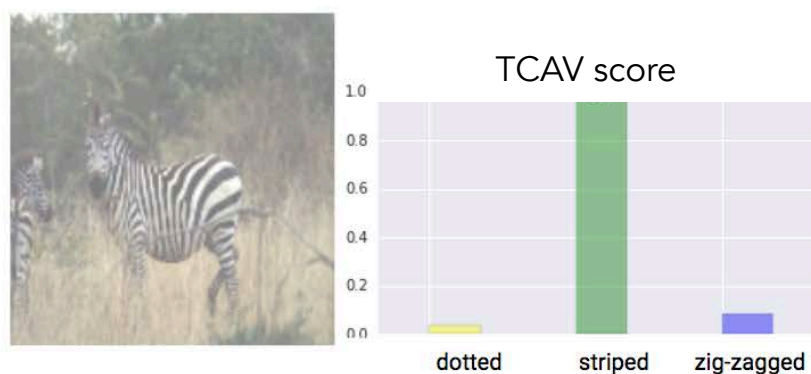


2. How are the CAVs useful to get explanations?

TCAV core idea:

Derivative with CAV to get prediction sensitivity

TCAV



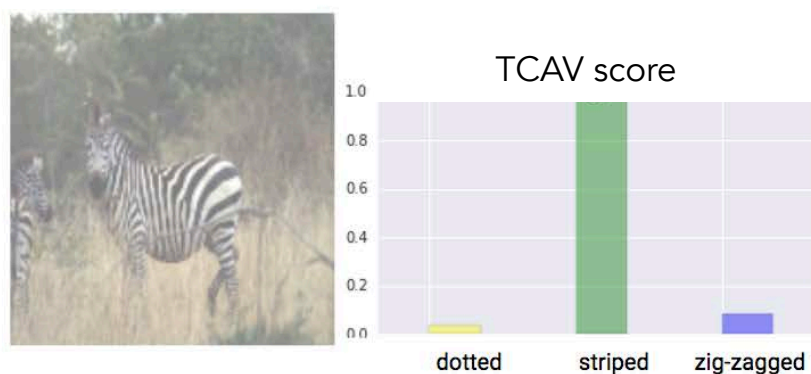
$$\begin{aligned} \text{zebra-ness} &\rightarrow \frac{\partial p(z)}{\partial \mathbf{v}_C^l} = S_{C,k,l}(\mathbf{x}) \\ \text{striped CAV} &\rightarrow \end{aligned}$$

Directional derivative with CAV

TCAV core idea:

Derivative with CAV to get prediction sensitivity

TCAV



$$\begin{aligned} \text{zebra-ness} &\rightarrow \frac{\partial p(z)}{\partial \mathbf{v}_C^l} = S_{C,k,l}(\mathbf{x}) \\ \text{striped CAV} &\rightarrow \end{aligned}$$

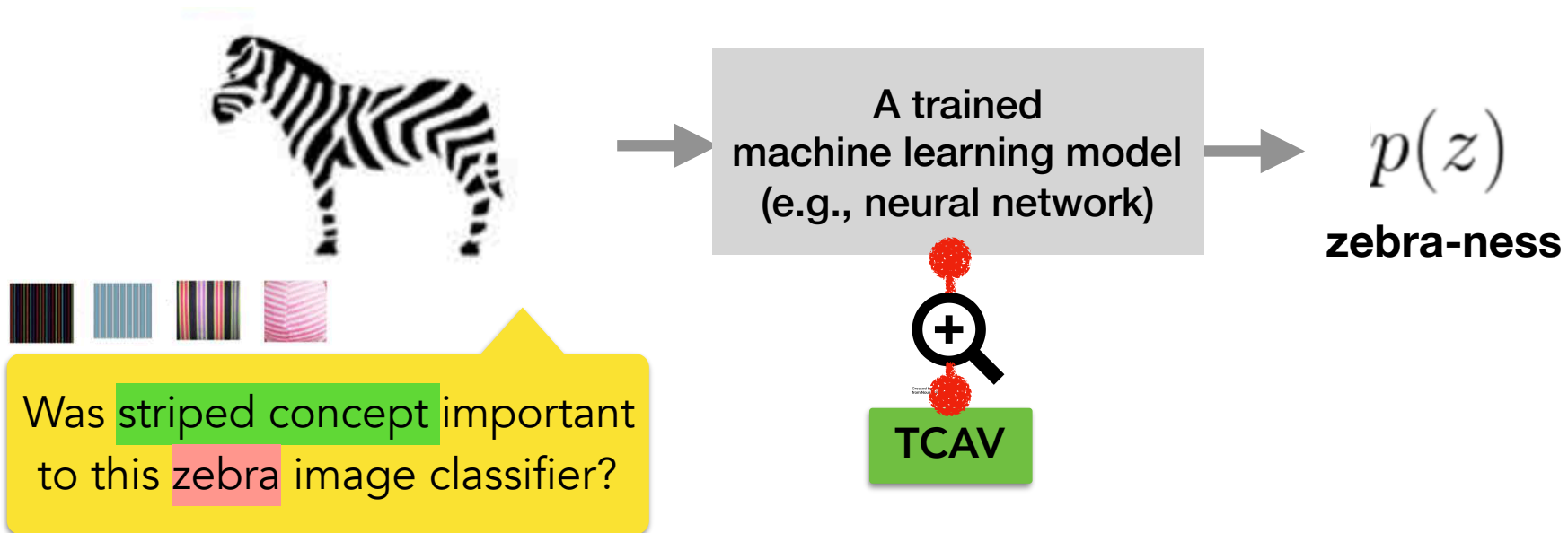
$$\left. \begin{aligned} S_{C,k,l}(\text{zebra}) \\ S_{C,k,l}(\text{zebra head}) \\ S_{C,k,l}(\text{zebra body}) \\ S_{C,k,l}(\text{zebra tail}) \end{aligned} \right\}$$

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

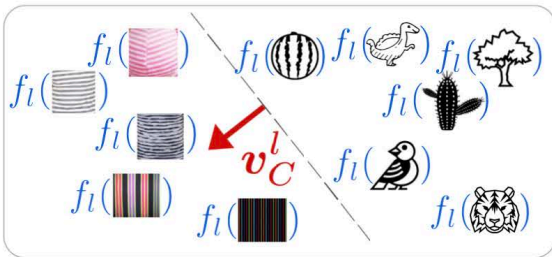
Directional derivative with CAV

TCAV:

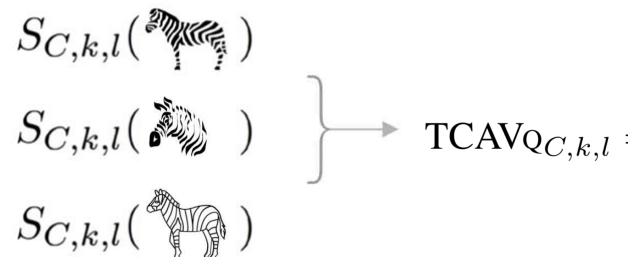
Testing with Concept Activation Vectors



1. Learning CAVs

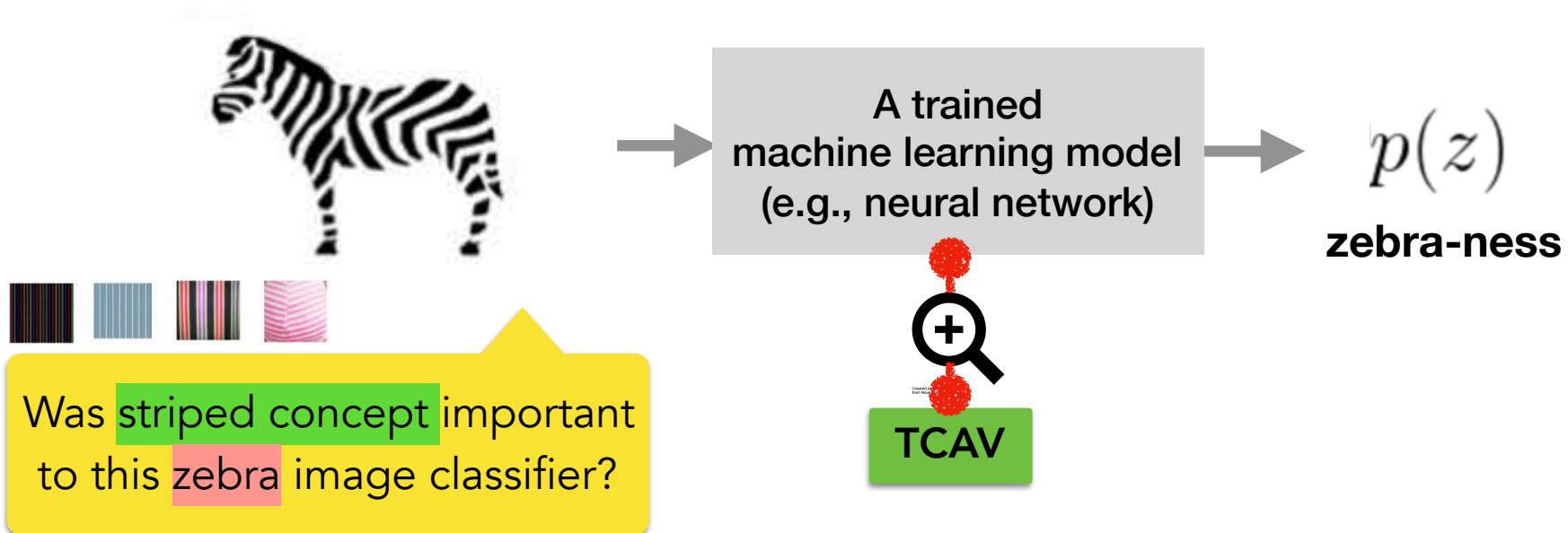


2. Getting TCAV score

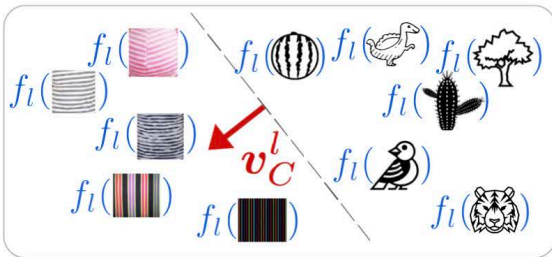


TCAV:

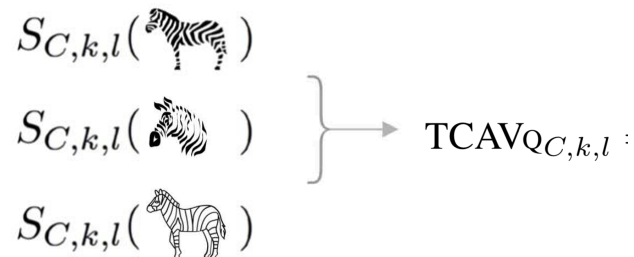
Testing with Concept Activation Vectors



1. Learning CAVs



2. Getting TCAV score



3. CAV validation

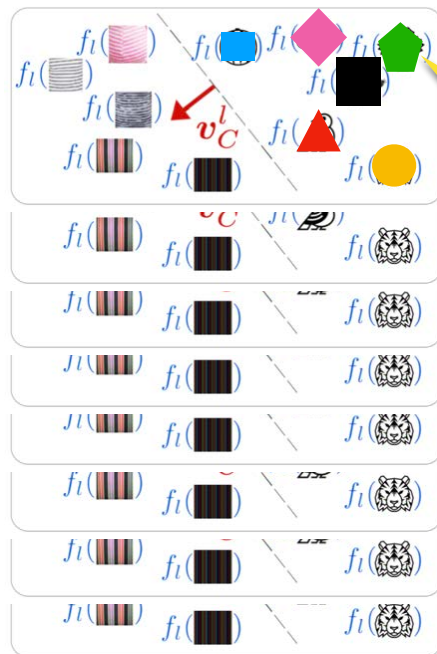
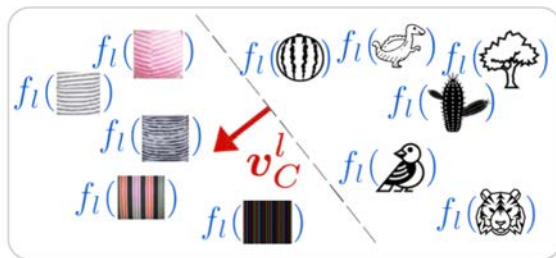
Qualitative
Quantitative

Quantitative validation:

Guarding against spurious CAV

Did my CAVs returned high sensitivity by chance?

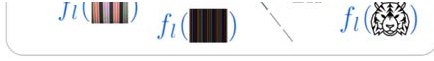
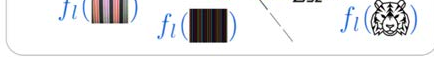
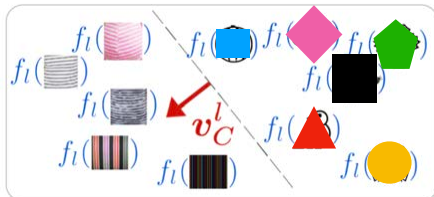
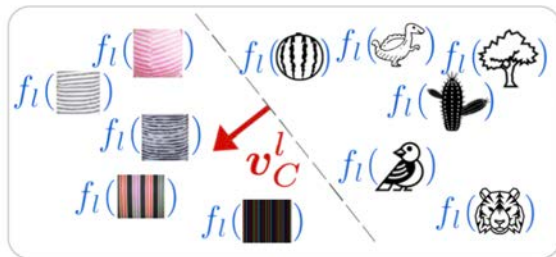
Quantitative validation: Guarding against spurious CAV



Learn many stripes CAVs
using different sets of
random images

Quantitative validation:

Guarding against spurious CAV



Zebra

→ $\text{TCAV}_{Q_{C,k,l}}$:

⋮

→ $\text{TCAV}_{Q_{C,k,l}}$:

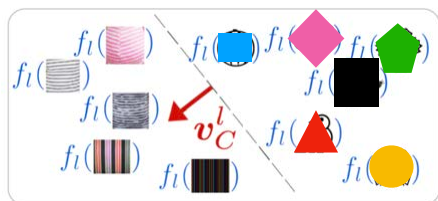
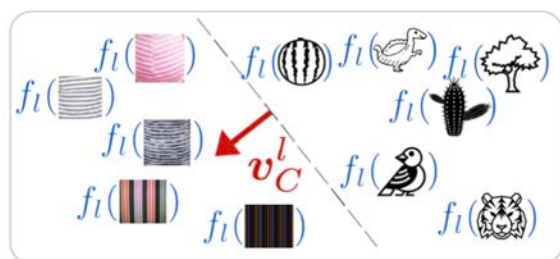
→ $\text{TCAV}_{Q_{C,k,l}}$:

→ $\text{TCAV}_{Q_{C,k,l}}$:

⋮

Quantitative validation:

Guarding against spurious CAV



Zebra

→ $\text{TCAV}_{Q_C, k, l} :$

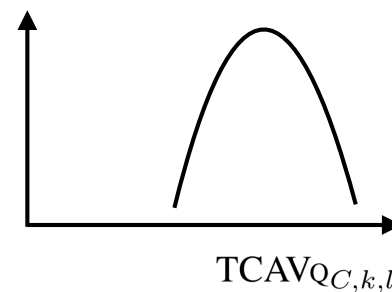
⋮

→ $\text{TCAV}_{Q_C, k, l} :$

→ $\text{TCAV}_{Q_C, k, l} :$

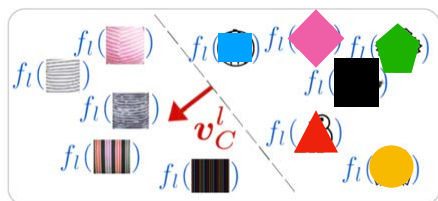
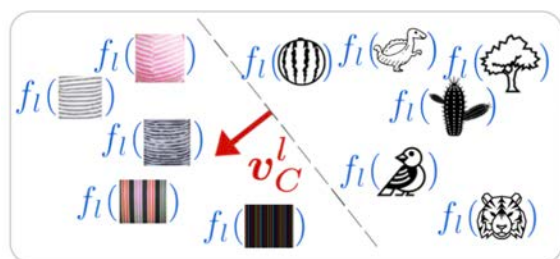
→ $\text{TCAV}_{Q_C, k, l} :$

⋮



Quantitative validation:

Guarding against spurious CAV



Zebra

→ $\text{TCAV}_{Q_{C,k,l}} :$

⋮

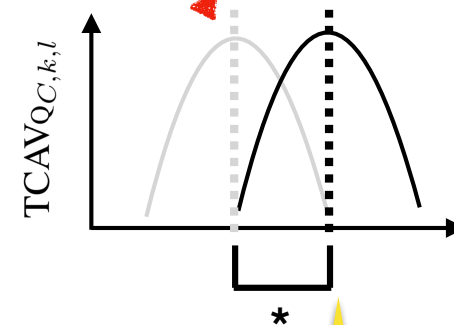
→ $\text{TCAV}_{Q_{C,k,l}} :$

→ $\text{TCAV}_{Q_{C,k,l}} :$

→ $\text{TCAV}_{Q_{C,k,l}} :$

⋮

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

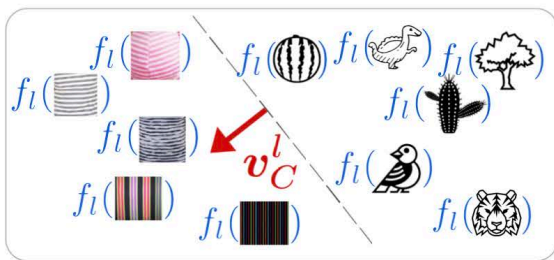


TCAV provides **quantitative importance** of a concept **if and only if** your network learned about it.

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

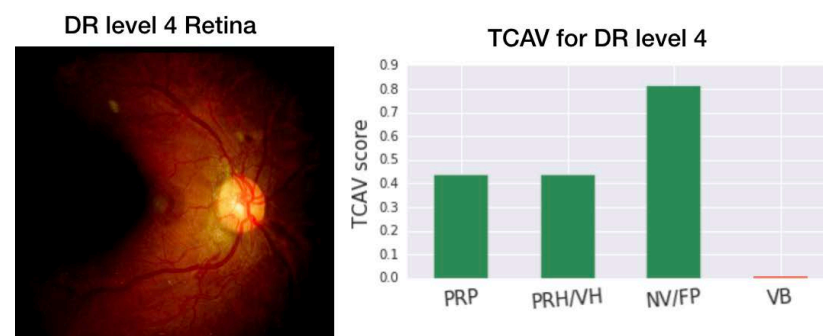
$$\left. \begin{array}{l} S_{C,k,l}(\text{zebra}) \\ S_{C,k,l}(\text{zebra}) \\ S_{C,k,l}(\text{zebra}) \end{array} \right\} \rightarrow \text{TCAV}_{QC,k,l}$$

3. CAV validation

Qualitative
Quantitative

Results

1. Sanity check experiment
2. Biases in Inception V3 and GoogleNet
3. Domain expert confirmation from Diabetic Retinopathy



Results

BIM results

1. Sanity check experiment



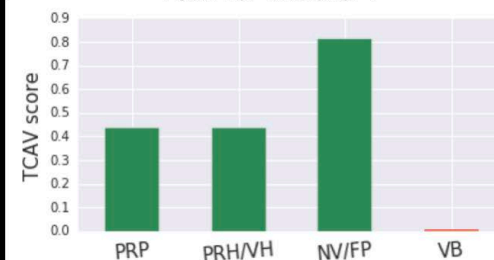
2. Biases from Inception V3 and GoogleNet

3. Domain expert confirmation from Diabetic Retinopathy

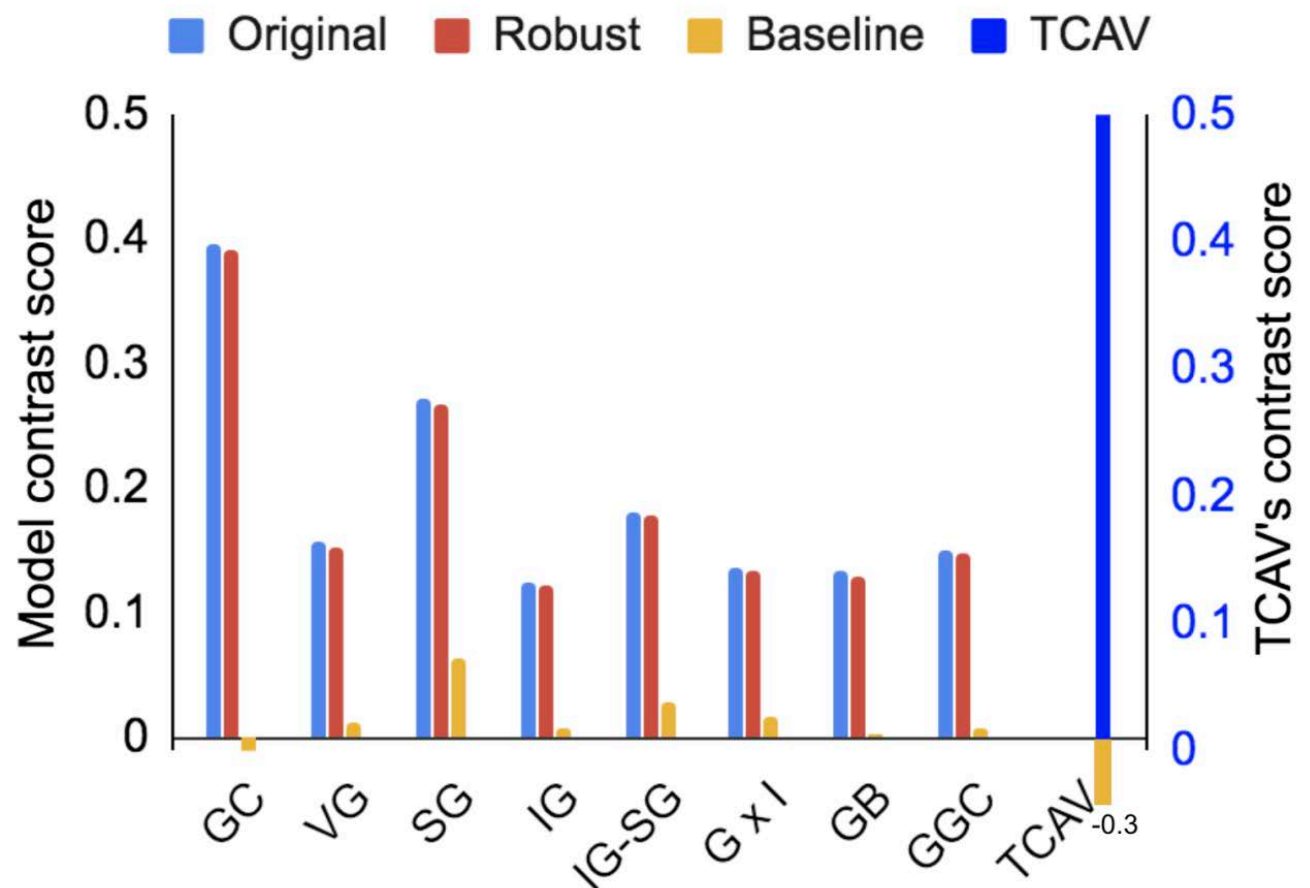
DR level 4 Retina



TCAV for DR level 4



Evaluating with BIM dataset



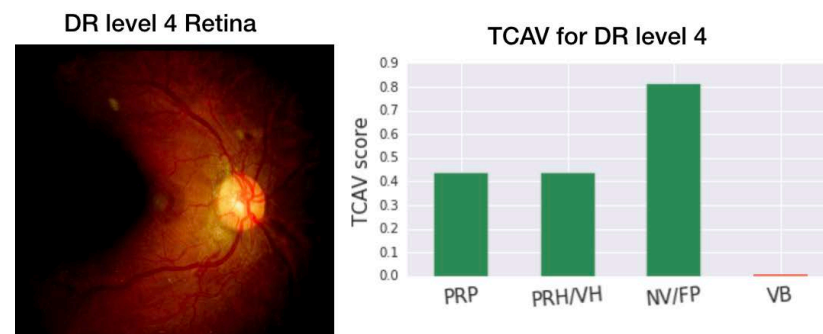
Results

1. Sanity check experiment

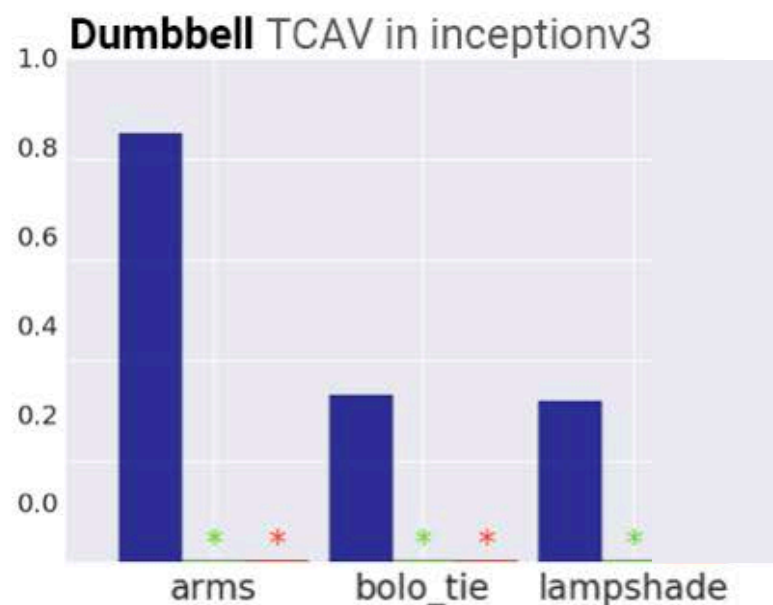
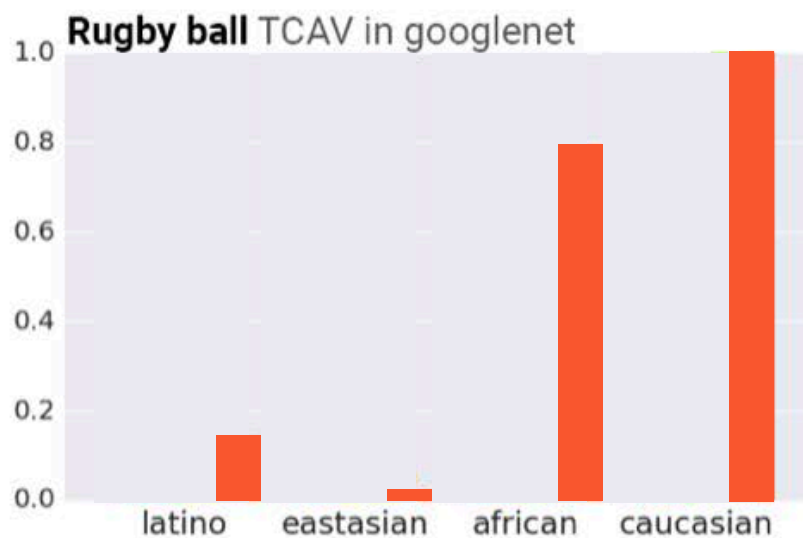
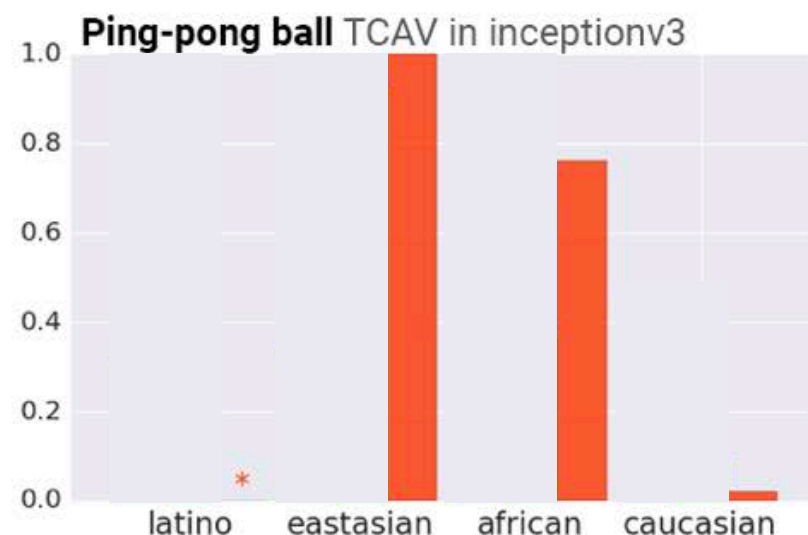
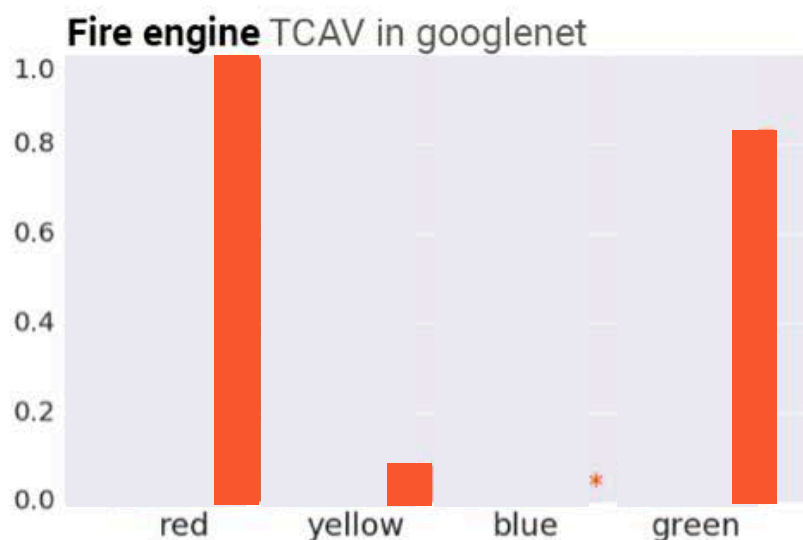


2. Biases from Inception V3 and GoogleNet

3. Domain expert confirmation from Diabetic Retinopathy

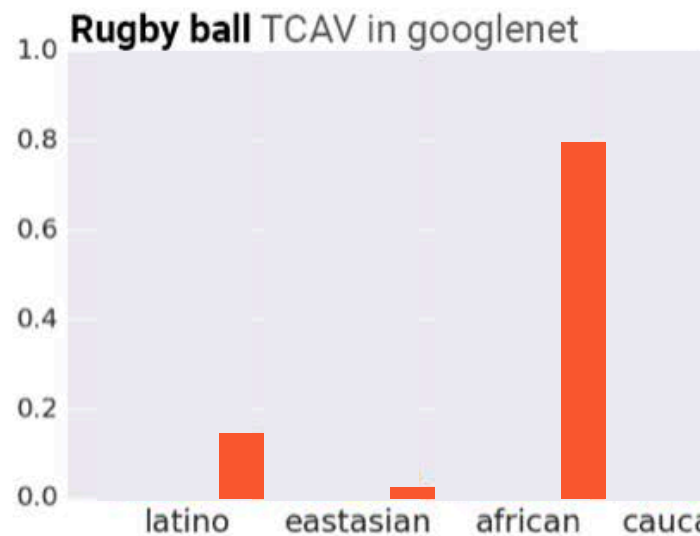
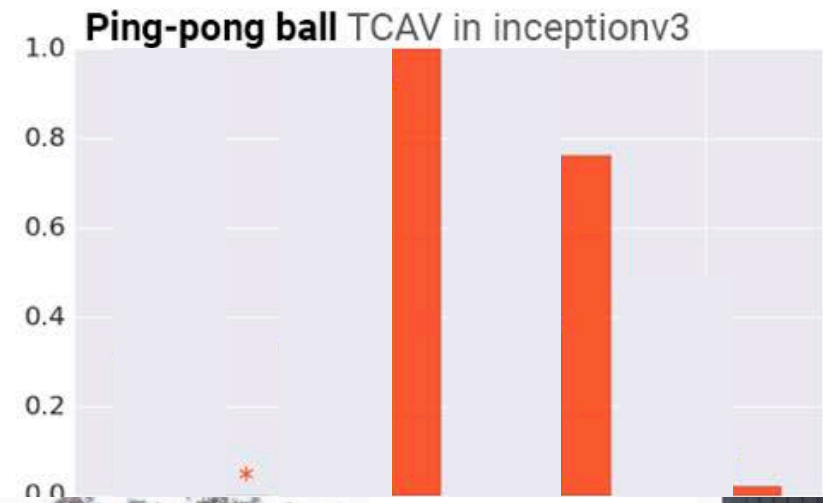
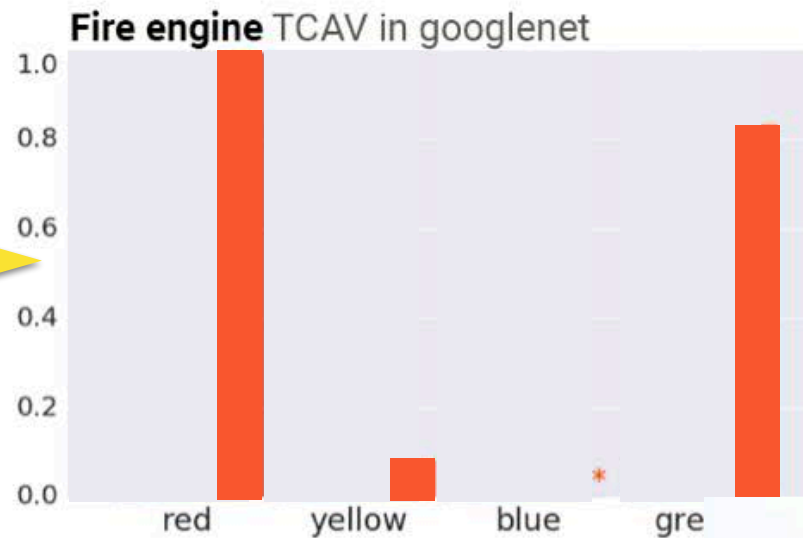


TCAV in Two widely used image prediction models



TCAV in Two widely used image prediction models

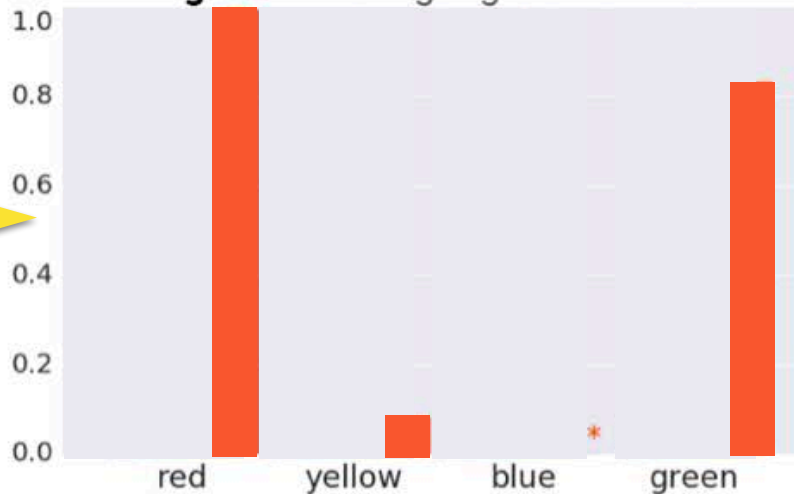
Geographical
bias!



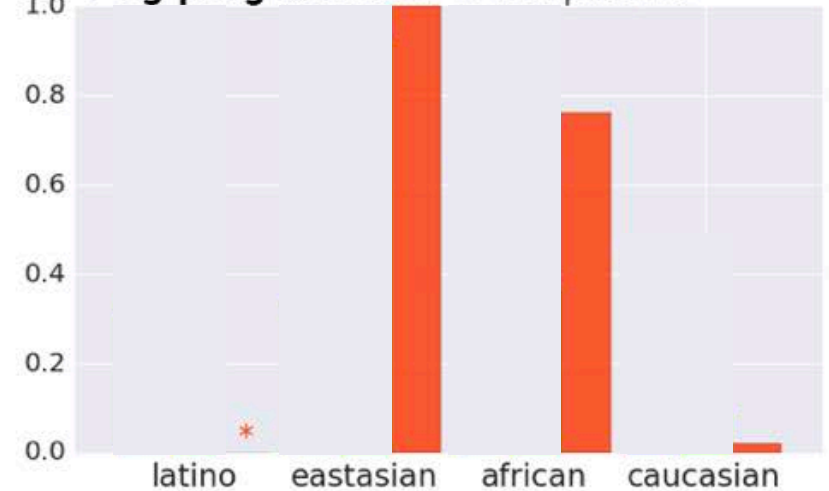
TCAV in Two widely used image prediction models

Geographical
bias?

Fire engine TCAV in googlenet

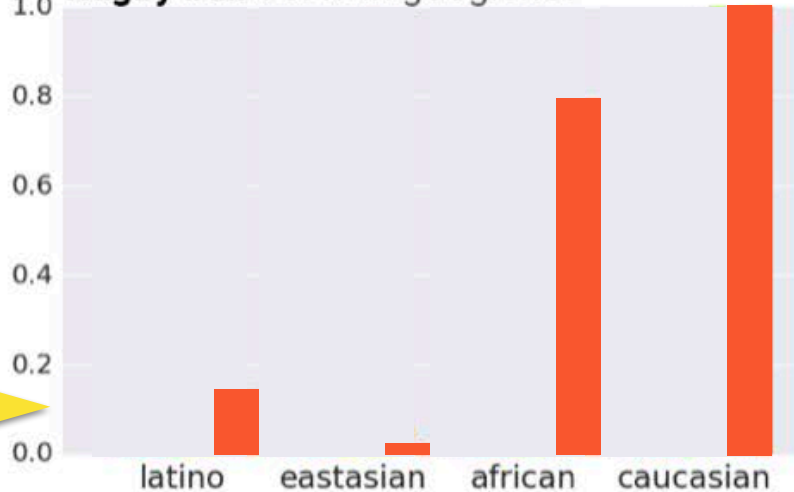


Ping-pong ball TCAV in inceptionv3

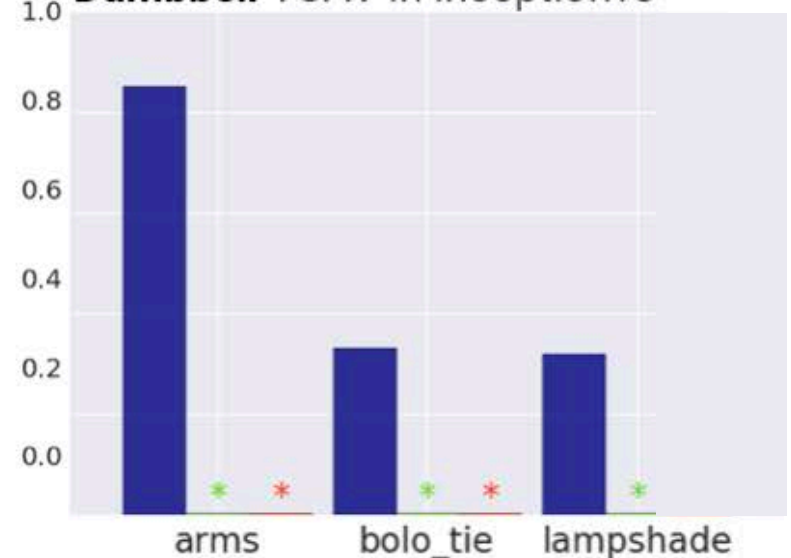


Quantitative
confirmation to
previously
qualitative
findings
[Stock & Cisse,
2017]

Rugby ball TCAV in googlenet



Dumbbell TCAV in inceptionv3

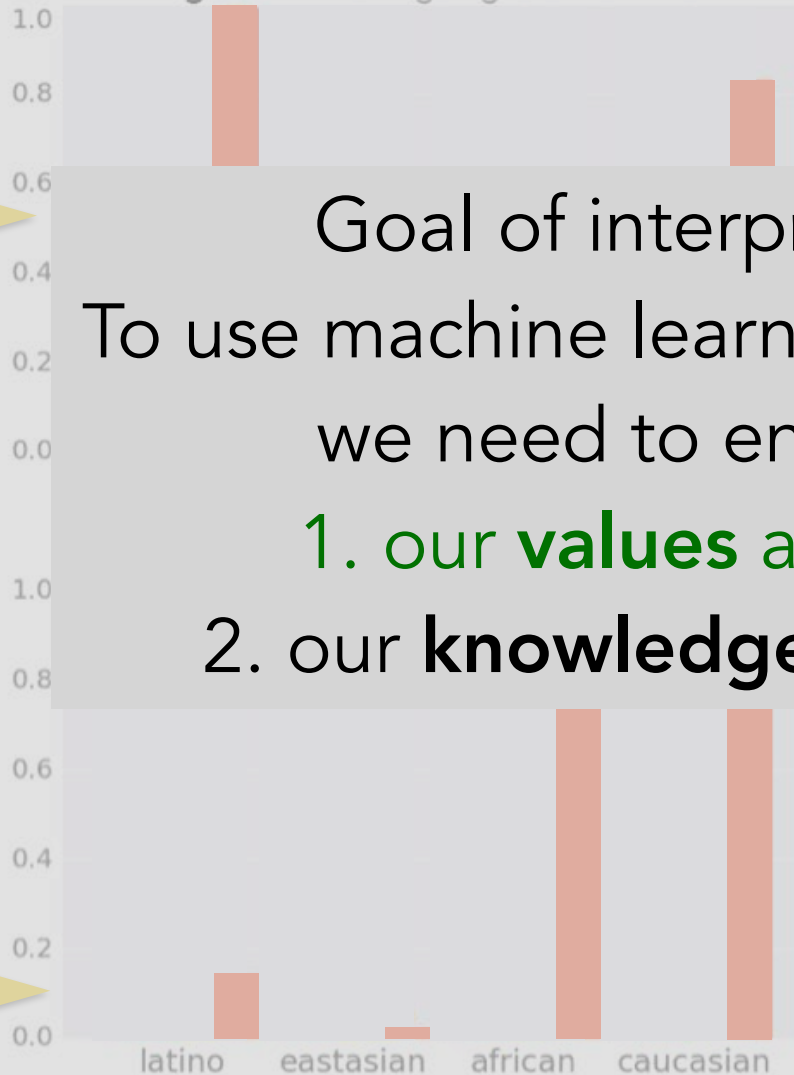


TCAV in

Two widely used image prediction models

Geographical
bias?

Fire engine TCAV in googlenet



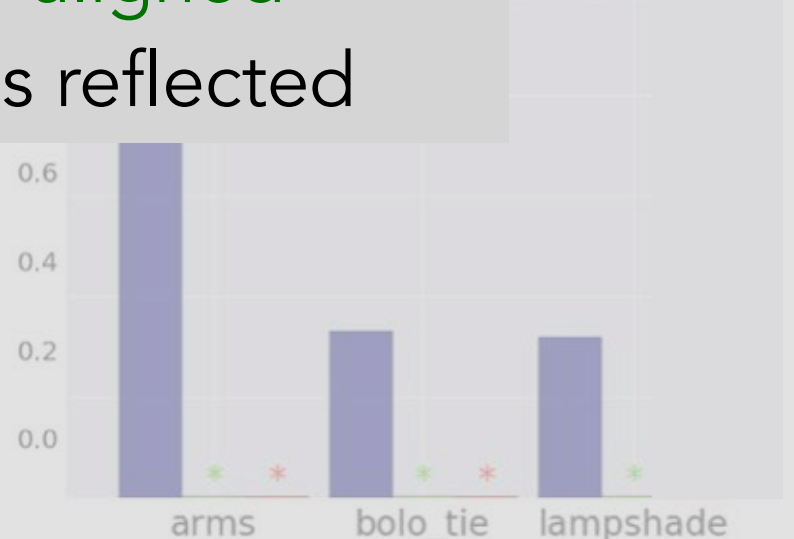
Ping-pong ball TCAV in inceptionv3



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

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

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



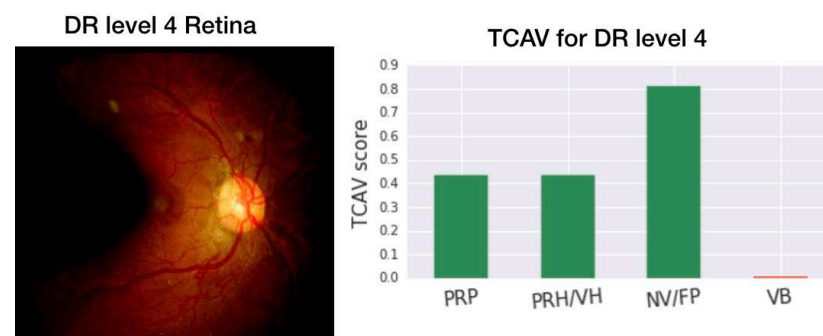
Results

1. Sanity check experiment



2. Biases Inception V3 and GoogleNet

3. Domain expert confirmation from Diabetic Retinopathy



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

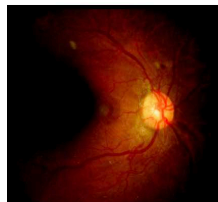
DR level 4 Retina



Collect human doctor's knowledge



DR level 4



Concepts
belong to
this level

PRP
PRH/VH
NV/FP

Concepts do not
belong to
this level

VB

DR level 1



MA

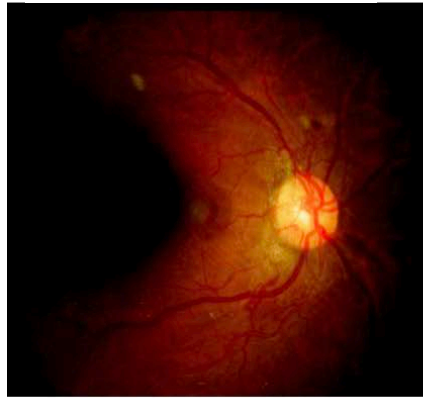
HMA

TCAV for Diabetic Retinopathy

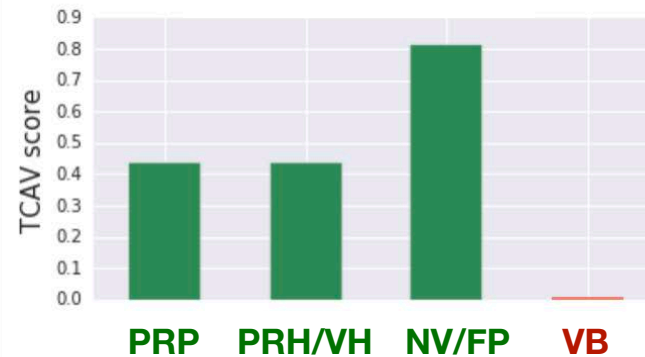
Prediction
class Prediction
accuracy

DR level 4 High

Example



TCAV scores



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

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

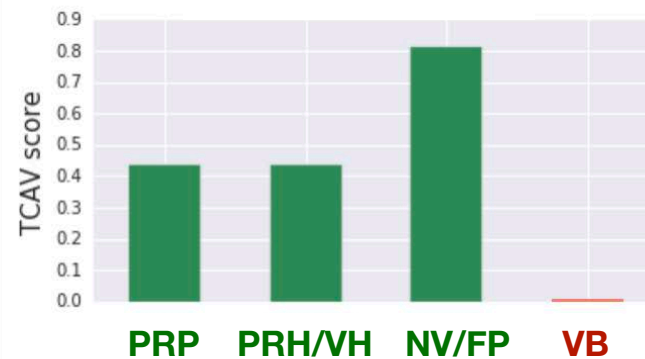
Prediction class Prediction accuracy

DR level 4 High

Example



TCAV scores

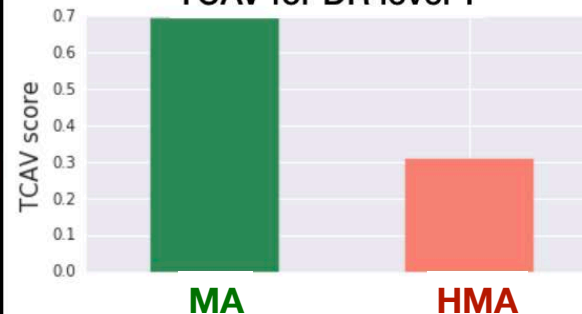


TCAV shows the model is **consistent** with doctor's knowledge when model is **accurate**

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

DR level 4

Hi

Level 1 was often confused to level 2.

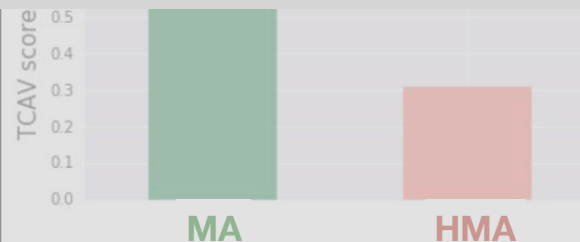
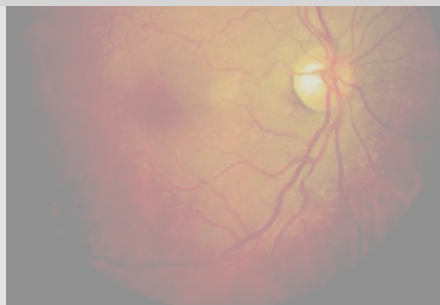
HMA distribution on predicted DR

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

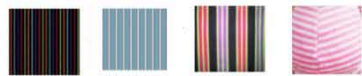
Green: domain expert's label on concepts belong to the level


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

Summary:

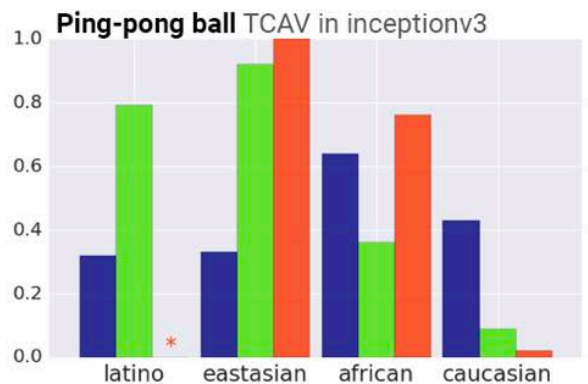
Testing with Concept Activation Vectors

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

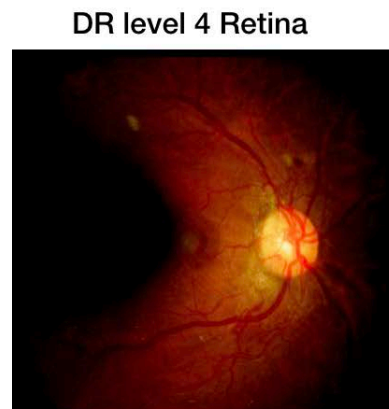


stripes concept (score: 0.9)
was important to **zebra** class
for this trained network. 

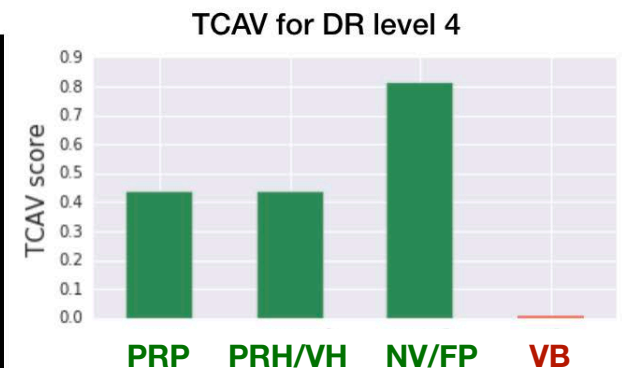
TCAV provides
quantitative importance of
a concept **if and only if** your
network learned about it.



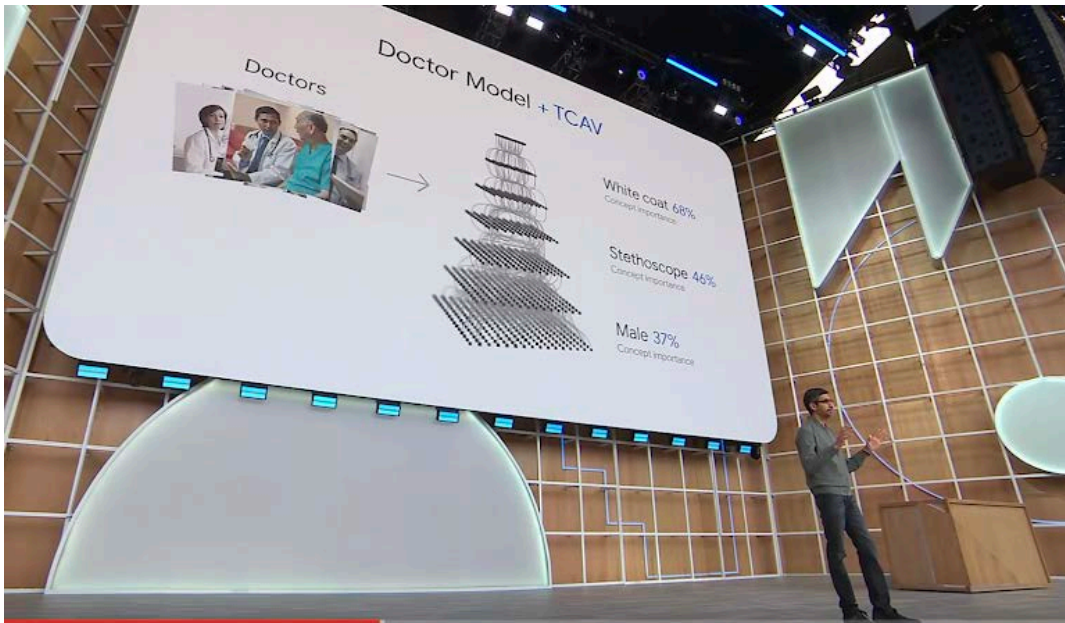
Our values



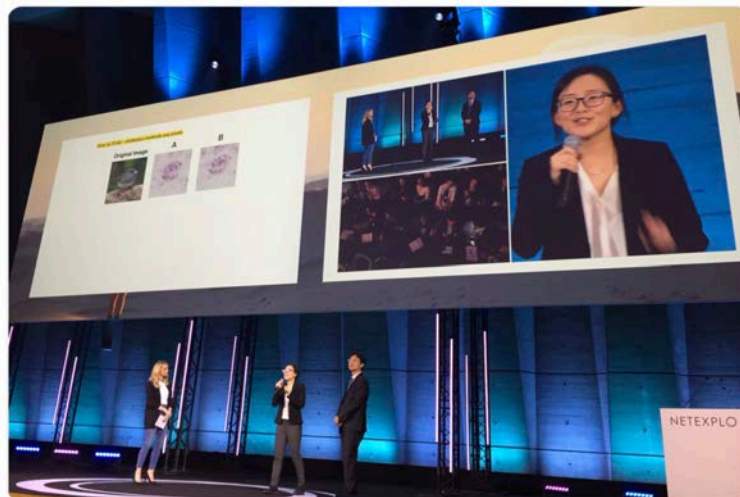
Our knowledge



Responses from outside of academia



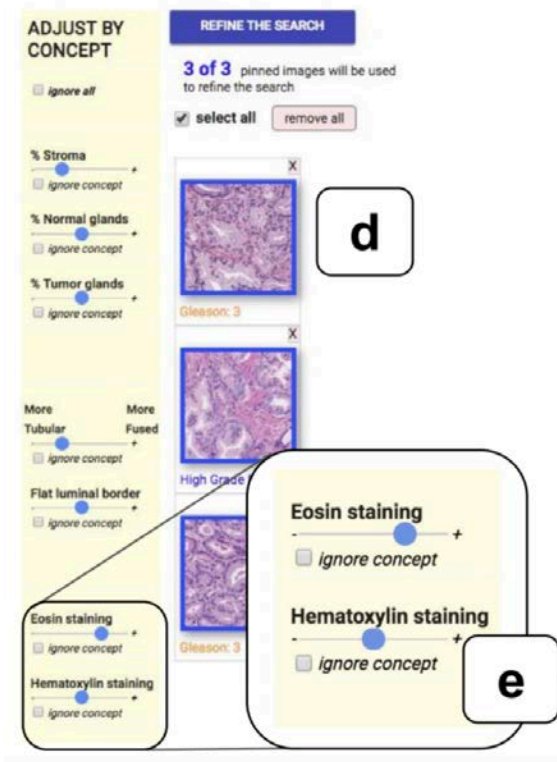
Sundar (CEO of Google) explaining how TCAV works in his keynote at Google I/O 2019



UNESCO NetExplo award 2019

Selected as one of ten "cutting-edge digital innovations with the potential of profound and lasting impact."

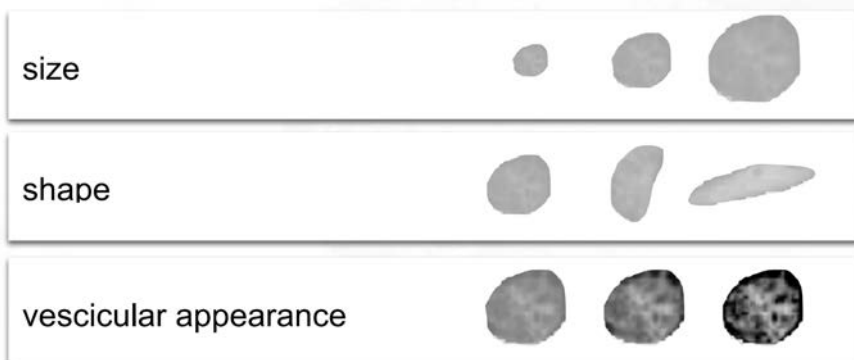
Responses from inside of academia



Using CAVs to help doctors find more diagnostically relevant images

Human-Centered Tools for Coping with Imperfect Algorithms during Medical Decision-Making

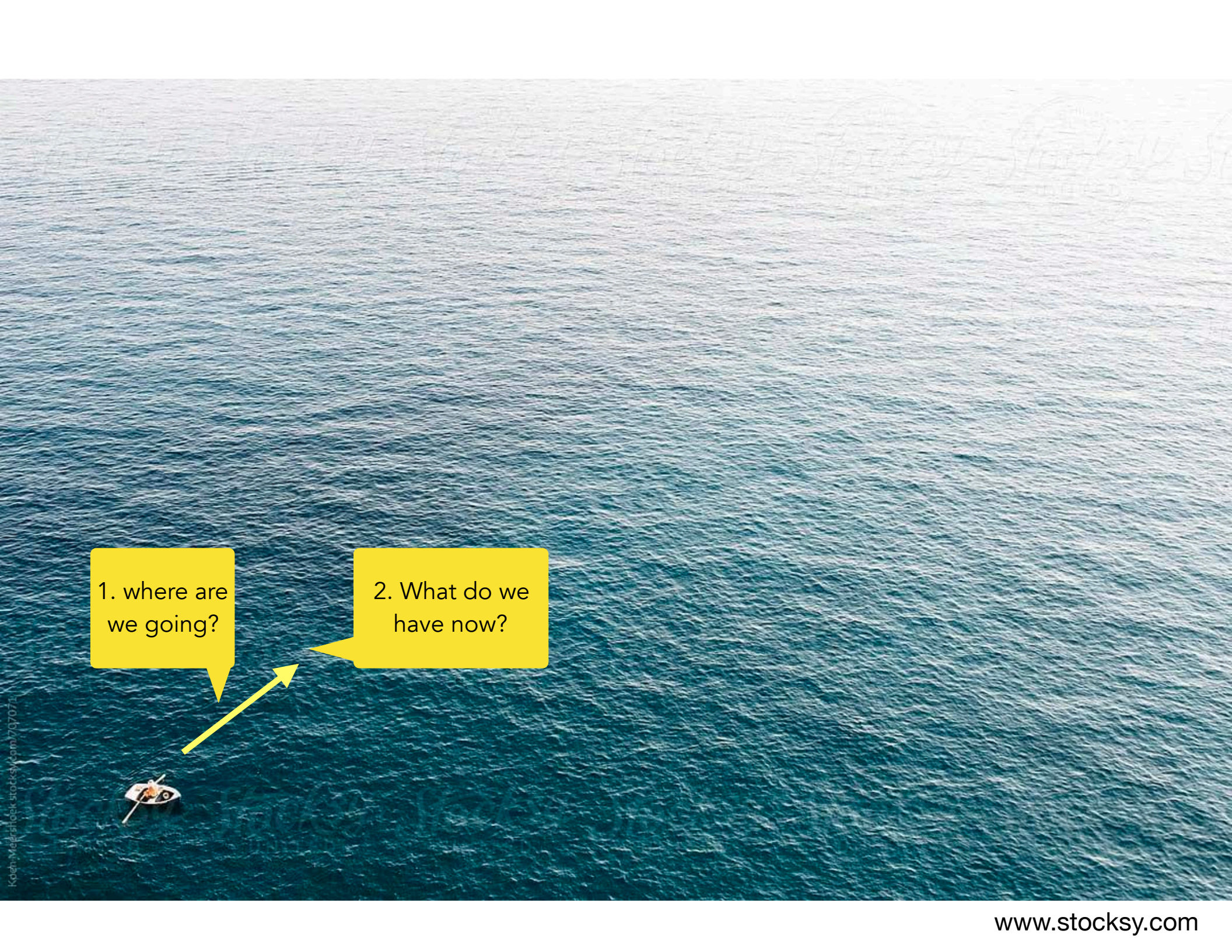
Work by Carrie J. Cai, Emily Reif, Narayan Hegde, Jason Hipp, K., Daniel Smilkov, Martin Wattenberg, Fernanda Viegas, Greg S. Corrado, Martin C. Stumpe, Michael Terry **CHI conference, best paper honorable mention**



Extending TCAV to regression models

"Regression Concept Vectors for Bidirectional Explanations in Histopathology"

Work by Mara Graziani, Vincent Andrearczyk, Henning Muller



1. where are we going?

2. What do we have now?

3. What can we do better?

Limitations of TCAV

- Concept has to 'expressible' using examples (e.g., "love" concept might be hard).
- User needs to know which concepts they want to test, and have examples for it. Follow-up work to automatically discover concepts for images (submitted), but many more directions are possible.
- Explanations provided by TCAV are not-causal
 - Follow-up work on causal TCAV (submitted)



1. where are we going?

2. What do we have now?

3. What can we do better?

4. What should we be careful?





Things to keep in mind during our journey.

- Proper evaluations
 - Sanity check and ground-truth-based evaluations
 - Test with humans!
- Remember that humans are biased and irrational.
- Importance of designing the right interaction - HCI.
- Try to criticize - think about what wasn't talked about in this talk but should have!
- Keep checking if we are going to the right direction!

1. where are we going?

Tool that can help more responsible AI

2. What do we have now?

Some existing methods fail a simple sanity check.

3. What can we do better?

TCAV

(btw it passes sanity check)

4. What should we be careful?

Evaluation

HCI

Human biases.

