# Interpretability for philosophical and skeptical minds

Been Kim

# What is the best explanation?
# From philosophers

- Many attempts to come up with **a single model of explanation**

  - Deductive Nomological (1942, Hempel) Statistical relevance (1971, Salmon), Causal Mechanical (1984, Salmon), Unificationist (1974, Freidman, 1989, Kitcher) with the hope that there exists ONE OPTIMAL model for explanations.

- Then pragmatic theories (1980, van Fraassen) came out.

  > The discussion of explanation went wrong at the very beginning when explanation was conceived of as a relation like description: a relation between a theory and a fact. Really, it is a three-term relation between theory, fact, and context. No wonder that no single relation between theory and fact ever managed to fit more than a few examples! Being an explanation is essentially relative for an explanation is an *answer*… it is evaluated vis-à-vis a question, which is a request for information. But exactly… what is requested differs from context to context. (1980: 156)

- The importance of "context" is something that ML community also came to (generally) agree.

- I doubt we will over find a single best model of explanation without context.

Carl Gustav Hempel

Wesley C. Salmon

Bas van Fraassen

# What is the best explanation?
## Illuminating example

- Structural explanation by Prof. Sally Haslanger



Prof. Haslanger

- The Invisible Foot (Okin 1989, Cudd 2006): Lisa and Larry, equally intelligent and talented at work, both capable of taking care of a child. But they live in a society where there is wage gap between men and women. They don't have means to pay for childcare. Lisa decides to quit her job.

- What's the "best" explanations for "why did Lisa quit her job?"

  - Why did **Lisa** quit her job and Larry?

  - Why did Lisa **quit** instead of going part time?

  - The society that unconsciously shaped her preference? "I'm not as good as Larry".

  - The society that created the bias and wage gap?

  - …

# What is the best explanation?
# Illuminating example

- Structural explanation by Prof. Sally Haslanger

- The Invisible Foot (Okin 1989, Cudd 2006): Lisa and Larry, equally intelligent and talented at work, both capable of taking care of a child. But they live in a society where there is wage gap between men and women. They don't have means to pay for childcare. Lisa decides to quit her job.

Prof. Haslanger

- What's the "best" explanations for "why did Lisa quit her job?"

  - Why did **Lisa** quit her job and Larry?

  - Why did Lisa **quit** instead of going part time?

  - The society that unconsciously shaped her preference? "I'm not as good as Lar

  - The society that created the bias and wage gap?

  - ...

What's the "best" explanations for "why was this predicted as a dog?"

- This and that pixel?

- Other training data and their delicate interaction during training process?

- The choice of architecture or optimizer?

- How the pictures are taken and when?

- The human history of domesticating wolves into dogs...

# Wait, why are we talking about philosophy?

- Giving "explanations" isn't a new problem. It's century-old one.

- The complexity of "how/what/when" to explain: it's always more complicated than we think.

- We should not take "good" explanation on its face value: we need to be skeptical (as we will see more soon).

# Trying to understand something new isn't new. Neuroscience?



- Understanding human brain: came a long way, but not enough.

  - **"We still don't understand** a worm (Caenorhabditis elegans) with 302 neurons. Humans have **86 billion** of them." – Koch, Allen institute for brain science.

  - "Let's say we could actually record from 1 million neurons in a brain while it's operating. You'd get a lot of data, but **what would we look for**? That is what we have to get some idea of." – Prof. Roland

  - "Throughout, Understanding the Brain reads like a compendium of things we still don't know. **We don't know how many neurons** are in the human brain. [..] **We don't know how alcohol relieves anxiety**, or how dopamine signaling is impaired in schizophrenia[...]" – article



*BRAINS!* —
## *Understanding the Brain* is a catalog of all we don't know about the brain

Updated version of *Creating Mind* mostly tells us what we don't understand.

DIANA GITIG - 11/10/2018, 7:00 AM



## Will It Ever Be Possible to Understand the Human Brain?

Despite technical breakthroughs like Elon Musk's Neuralink, scientists still have no reliable model of how the brain actually works

Brian Bergstein  Aug 21, 2019 · 15 min read ★

# Oh bummer... Are you still giving this talk?

- Yes, neuroscience feels like my future in 40 years... "We still don't understand..."

- But, I'm still optimistic. Because, while we still don't understand human brain, without a doubt **studying human brain helped the world,** because for example, 1) we have ways to help people via psychological treatments 2) we can sometimes cure seizure (e.g., epilepsy surgery) and the list goes on.

- The point is: the goal of interpretability is similar. it's not about understanding **everything** all the time. It's about understanding **enough** so that they are **useful.**

# What's enough?

- "This hammer isn't perfect, but it is good enough!

 [for what I am trying to do = context]"


I'm better off having this tool [for my goal/context]



inf.news

# What's enough in medicine?

- For example:

  - "Solve" medicine (?)

  - Help doctors to be more effective, efficient, and precise.

  - Use less resources, help more patients.

  - ...

  - ...

  - At minimum, do no harm.

Low bar

# What's enough in medicine?

- For example:

  - "Solve" medicine (?)

  - Help doctors to be more effective, efficient, and precise.

  - Use less resources, help more patients.

  - ...

  - ...

  - At minimum, do no harm.

Low bar

# Investigating
# post–training interpretability methods.

Input image



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

prediction

$p(z)$

Junco Bird-ness

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

Why was this a Junco bird?

Sanity Checks for Saliency Maps
Joint work with Adebayo, Gilmer, Goodfellow, Hardt, [NIPS 18]

# Investigating
# post-training interpretability methods.

Input image

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

prediction

$p(z)$

Junco Bird-ness

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

Why was this a Junco bird?

One definition of
explanation:

Tell me how **sensitive**
the prediction is when
we slightly **change**
each input feature
(pixel).

Sanity Checks for Saliency Maps
Joint work with Adebayo, Gilmer, Goodfellow, Hardt, [NIPS 18]

One of the most popular interpretability methods for images:

# Saliency maps

Input image



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

prediction
$p(z)$

Junco Bird-ness

In jargon: take derivative of the prediction wrt each pixel.

a logit $\longrightarrow$ $\dfrac{\partial p(z)}{\partial x_{i,j}}$
pixel i,j $\longrightarrow$

In English: take one pixel in the image, and imagine changing it by a little. See how much prediction changes. Do this for all pixels.

One definition of explanation:

Tell me how **sensitive** the prediction is when we slightly **change** each input feature (pixel).

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

One of the most popular interpretability methods for images:

# Saliency maps

Input image



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

prediction

$p(z)$

Junco Bird-ness

Popular method #1

Popular method #2

My work from 2018 #1

My work from 2018 #2

Popular method #3

Popular method #4

# Sanity check question

Input image



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

prediction

$p(z)$

Junco Bird-ness

So these pixels are the **evidence** of **prediction.**

Sanity Checks for Saliency Maps
Joint work with Adebayo, Gilmer, Goodfellow, Hardt, [NIPS 18]

# Sanity check question

Input image



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

prediction

$p(z)$

Junco Bird-ness

So these pixels are the **evidence** of **prediction.**

When **prediction** changes, the explanations will probably change.

When **prediction** is random, the explanations really should change!

Sanity Checks for Saliency Maps
Joint work with Adebayo, Gilmer, Goodfellow, Hardt, [NIPS 18]

# Some confusing behaviors of saliency maps.

**Original Image**

**Saliency map**

K<sup>th</sup> class

# Some confusing behaviors of saliency maps.

**Original Image**

**Saliency map**

K<sup>th</sup> class

Randomized weights!
Network now makes garbage prediction.

K<sup>th</sup> class

# Some confusing behaviors of saliency maps.



**Original Image**

$K^{th}$ class

**Saliency map**

!!!!!???!?

Randomized weights!
Network now makes garbage prediction.

**Original Image**

$K^{th}$ class

Sanity Checks for Saliency Maps
Joint work with Adebayo, Gilmer, Goodfellow, Hardt, [NIPS 18]

One of the most popular interpretability methods for images:

# Saliency maps

Input image



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

prediction

$p(z)$

Junco Bird-ness

Popular method #1

Popular method #2

My work from 2018 #1

My work from 2018 #2

Popular method #3

Popular method #4

# Sanity check1:
# When prediction changes, do explanations change?
# No!

# Sanity check2:
# Networks trained with random labels,
# Do explanations deliver different messages?
# No!

# Wait, what's so bad about this?

- What's this obsession about prediction? Maybe it's showing "features" that could have been 'used' in prediction. That's still relevant.


Your kidney


Your lung


Your pancreas


Your colon

Explanations: "Dotty" feature **used** to classify cancer.

Oh it's all cancer.

# Wait, what's so bad about this?

- What's this obsession about prediction? Maybe it's showing "features" that could have been 'used' in prediction. That's still relevant.



Cancer

Your kidney

Cancer

Your lung

Not cancer

Your pancreas

Not cancer

Your colon

Explanations: "Dotty" feature **used** to classify cancer.

# Many skeptics followed!  But still long way to go.

25

# But how do some of these methods still helpful for some end-tasks?

## ...

## What are those tasks?



[Adebayo, Muelly, Liccardi, K. Neurips 2020]
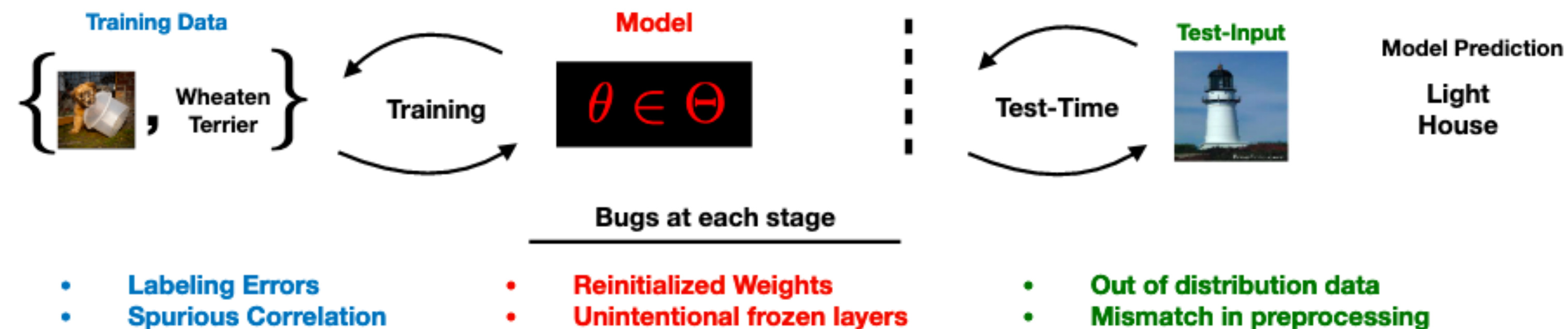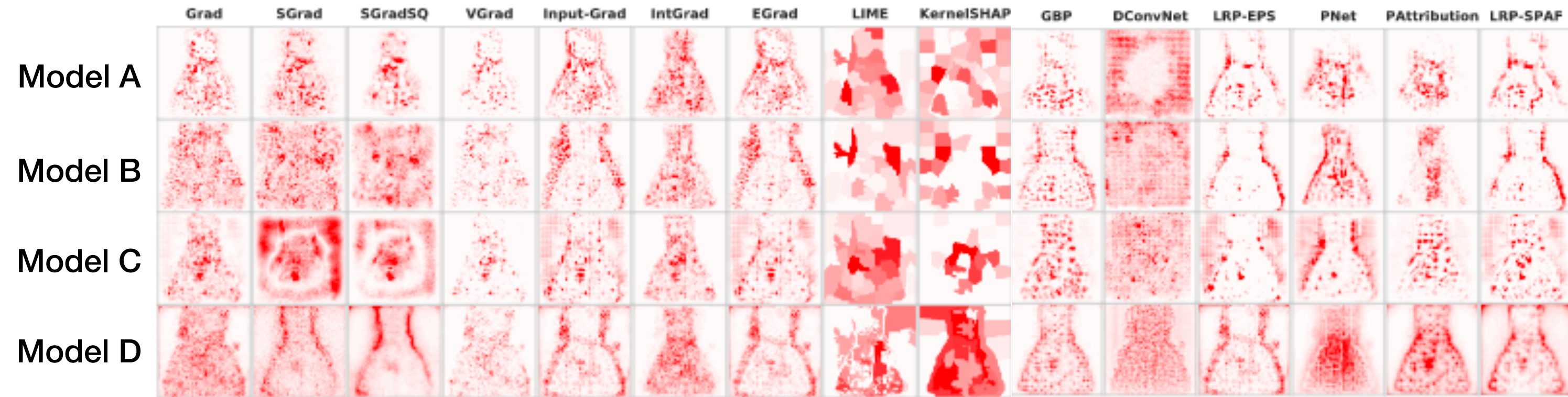
# Testing methods with users and concrete end-tasks
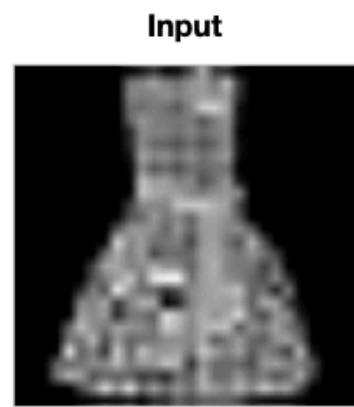


- **Task for subjects**: You work at a start-up selling animal classification ML model. Here are the images, predictions and attribution maps. (We gave users prediction labels as it is unrealistic not to).

- **Questions:** Would you recommend this model? Why? [because the wrong/correct label/explanation]? All in Likert scale.

[Adebayo, Muelly, Liccardi, K. Neurips 2020]

Can these methods tell us about
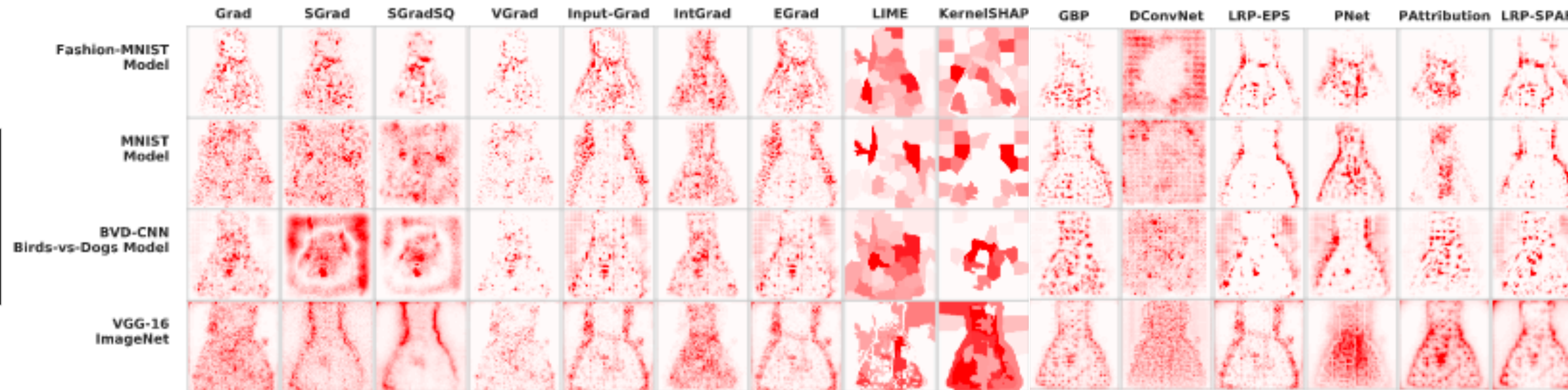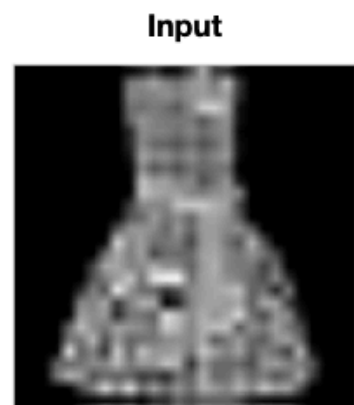
# Out of distribution?

- **Out of distribution data**

Can these methods tell us about

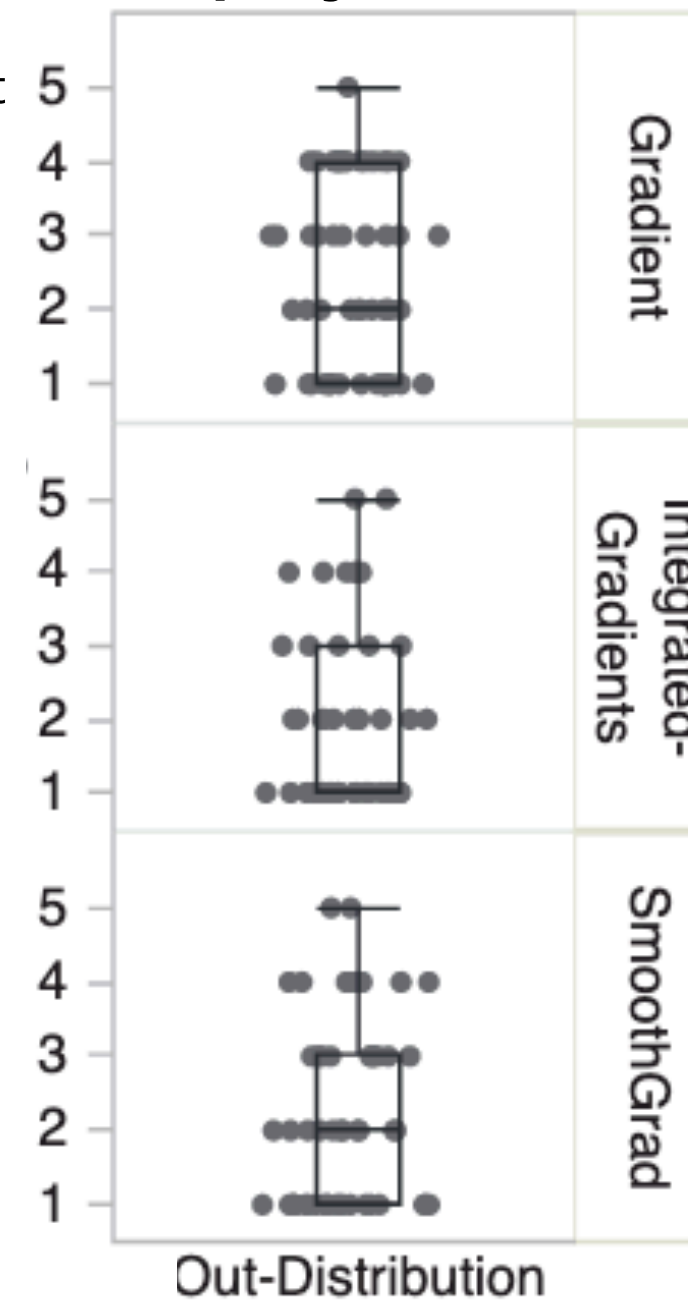# Out of distribution? probably not.
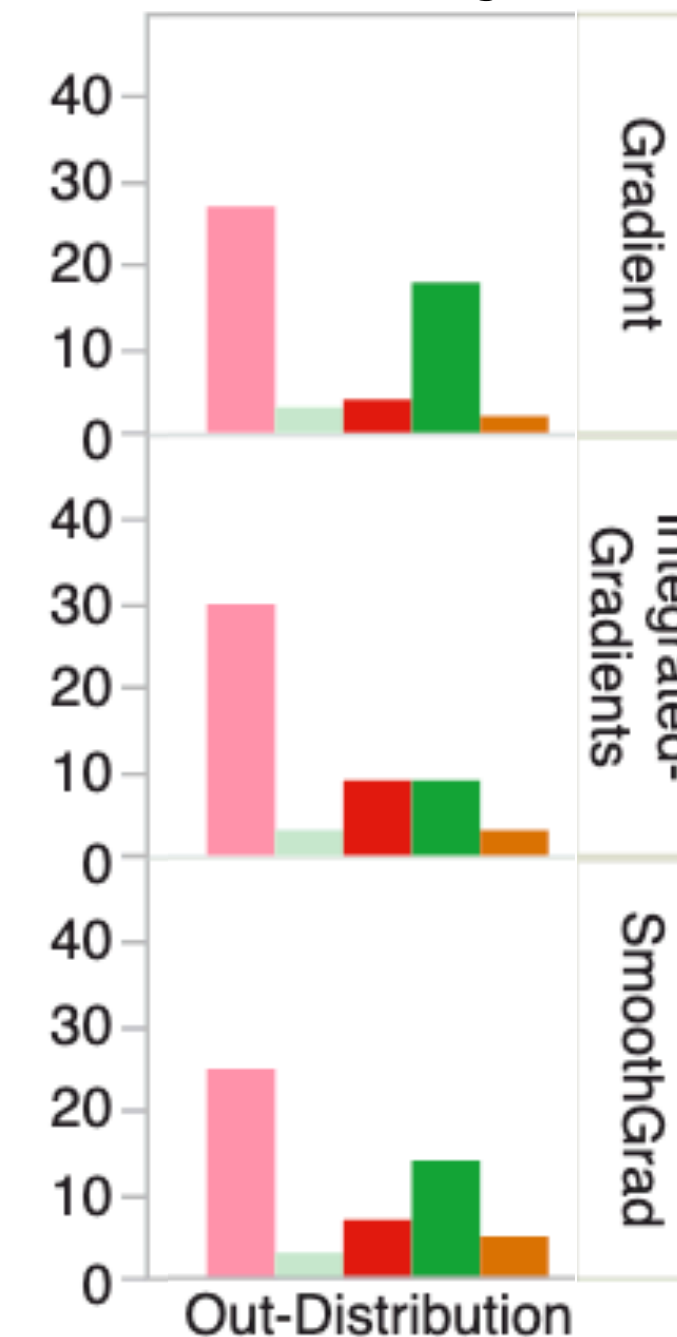
• **Out of distribution data**



Subjects are uncertain, mostly because of wrong label, but some **expected** explanations.

**How confident are you to deploy this model?**

Very confident

Not confident at all

**% why**

Wrong Label
Correct Label
Unexpected Explanation
Expected Explanation
Others

Can these methods tell us about

# Spurious correlation?

Input | Grad | SGrad | SGradSQ | VGrad | Input-Grad | IntGrad | EGrad | LIME | KernelSHAP | GBP | DConvNet | LRP-EPS | PNet | PAttribution | LRP-SPAF

Can these methods tell us about

# Spurious correlation? maybe!

Subjects are uncertain, mostly because of **unexpected** explanations!

How confident are you to deploy this model?

% why

Very confident

Not confident at all

Wrong Label
Correct Label
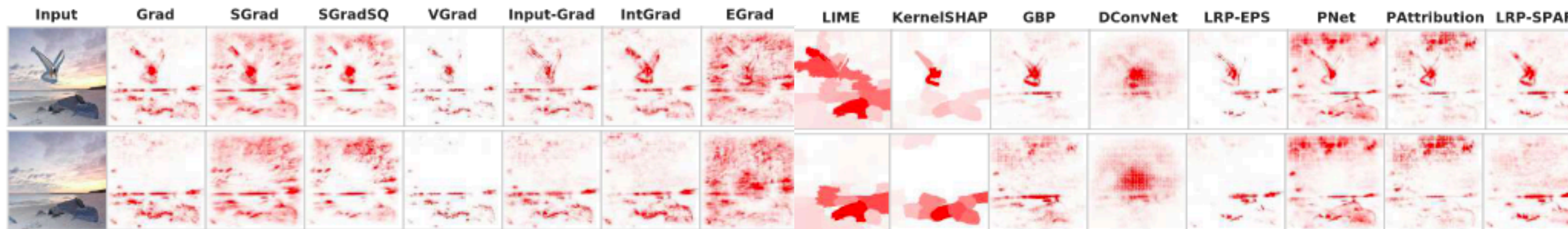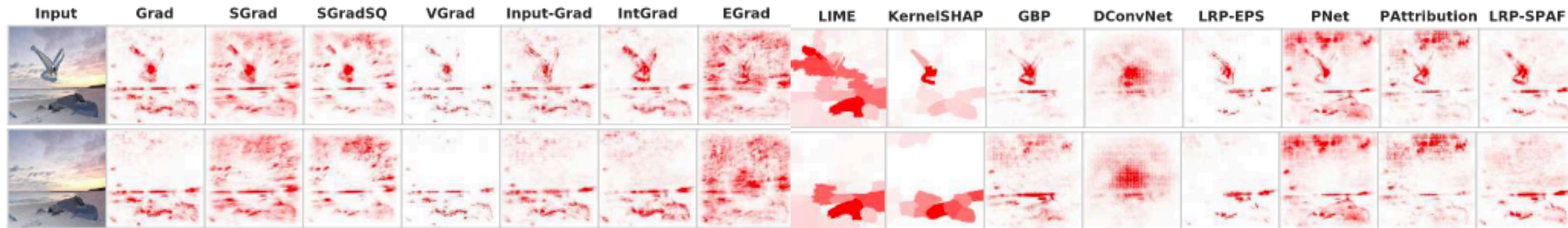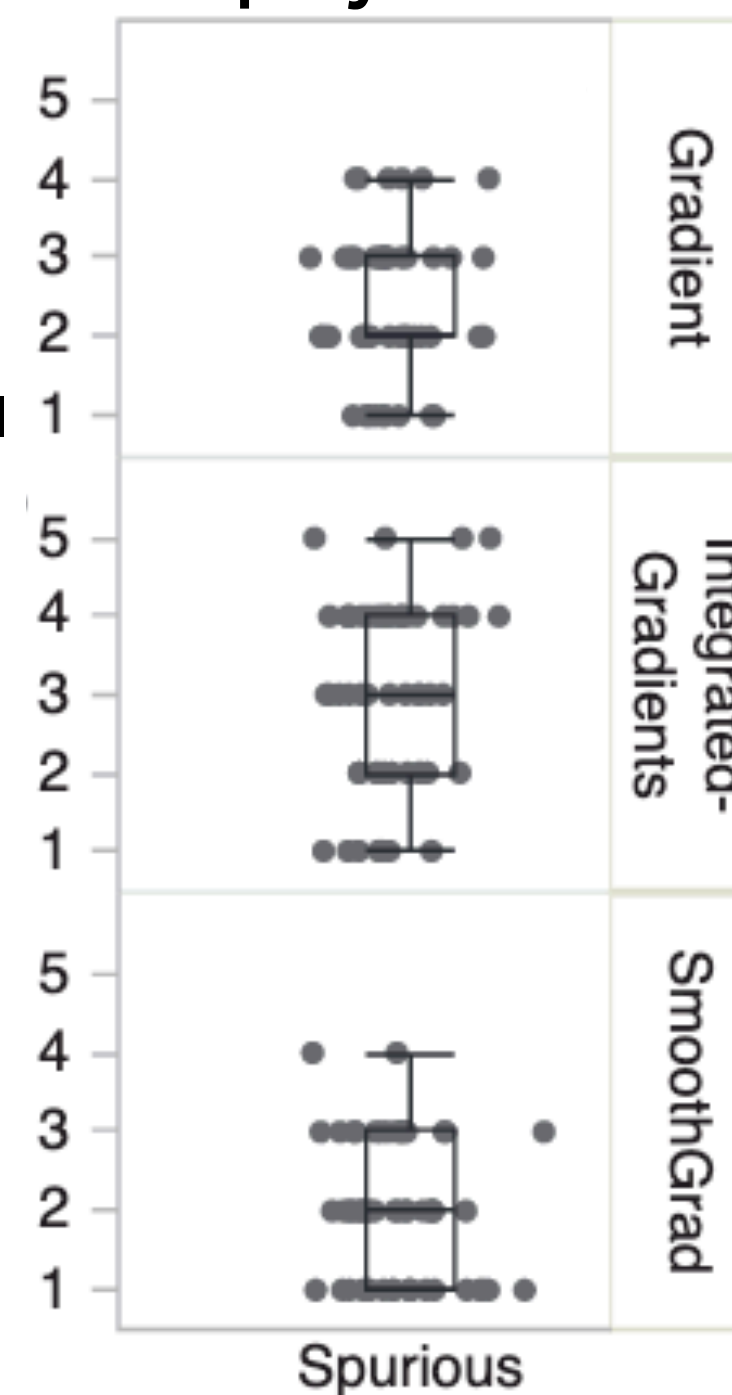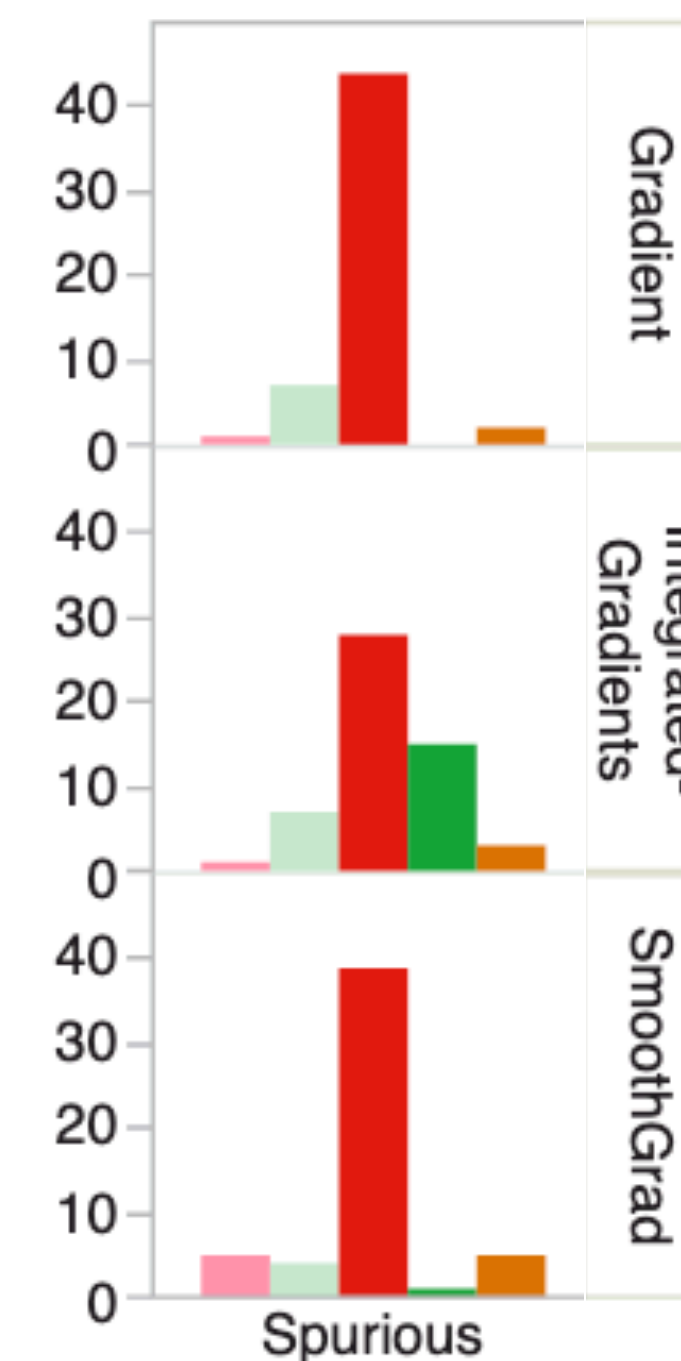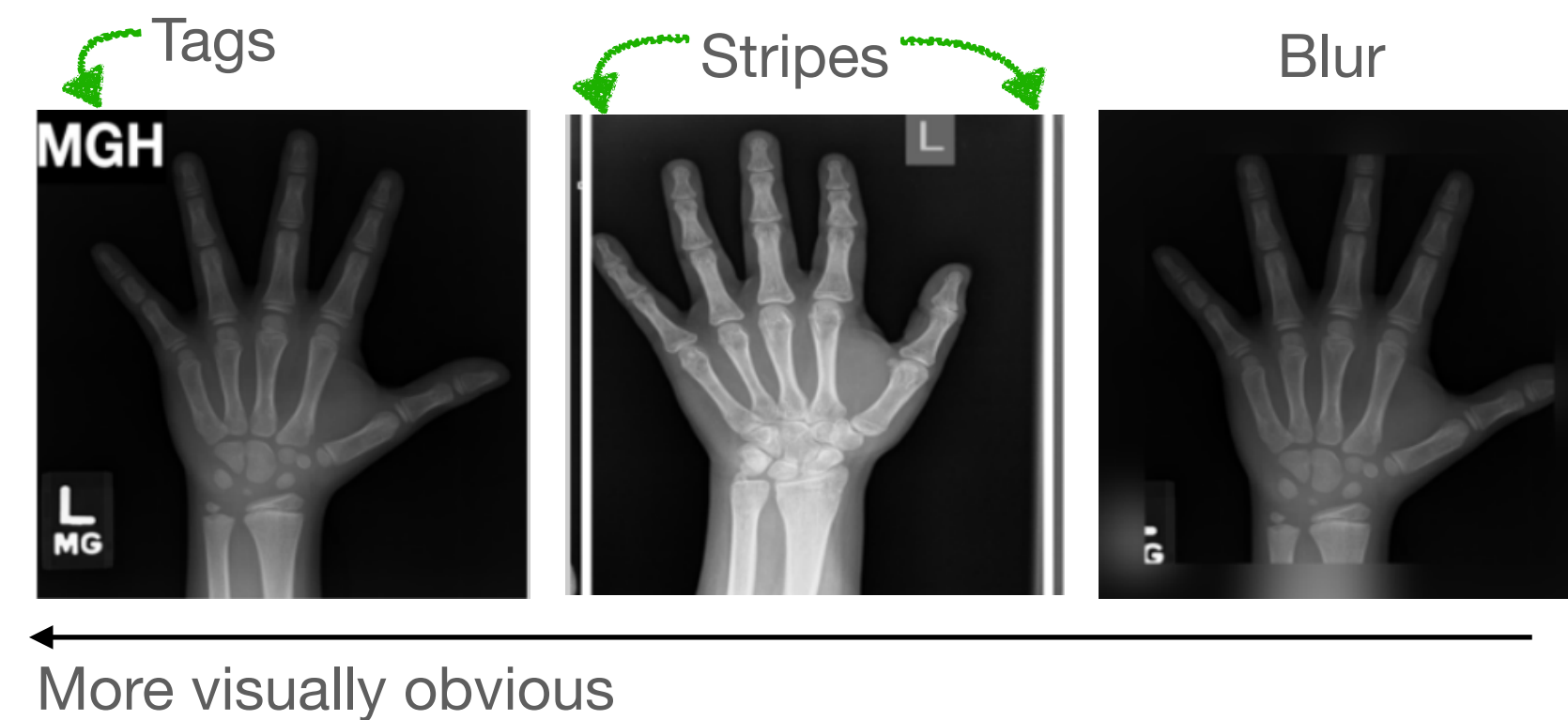Unexpected Explanation
Expected Explanation
Others

[Ongoing work]
What kind of spurious correlation can
we hope to capture?
TL;DR: Not many.
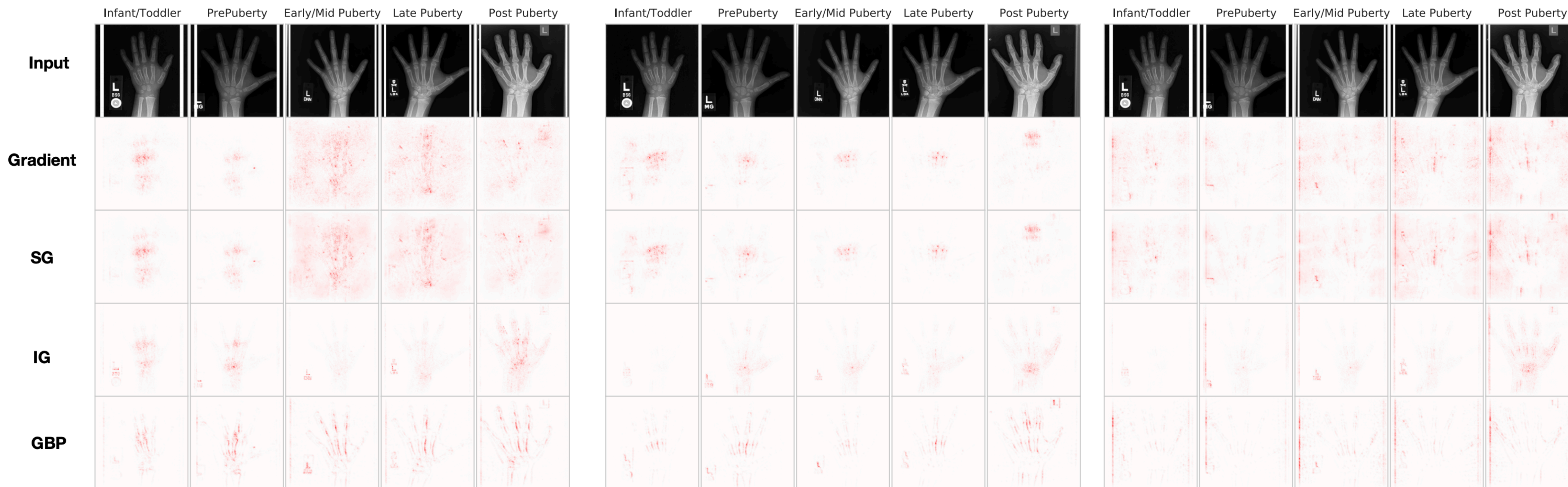


Tags      Stripes      Blur

More visually obvious

**A:** Normal Model Spurious Stripe Inputs     **B:** Spurious Stripe Model on 'Normal' Inputs     **C:** Spurious Stripe Model on Spurious Stripe Inputs

[Adebayo, Muelly, K. In submission]

# Take away

- Please be skeptical! Think of explanations as your (potentially incompetent) colleague. Maybe they are helpful, but maybe not.

- Explanations are complex in nature (we've known this for quite a few centuries); they are powerful, but we need to be careful how we use them.

- Many explanations can give plausible explanations, but we need to be careful (e.g., even explanations from an inherently interpretable model could be misleading in distributional shift)

- Test, test and test.