# Axiomatic Attribution for Deep Networks

Mukund Sundararajan (Google)
Joint work with Ankur Taly, Qiqi Yan

Why did the network label this image as "fireboat"?

# Problem Statement

**Attribute a <u>complex</u> deep network's prediction to its input features**

E.g. Attribute an Object Recognition Network's prediction to its pixels

E.g. Attribute a diagnostic network's prediction to patient's symptoms, measurements, history
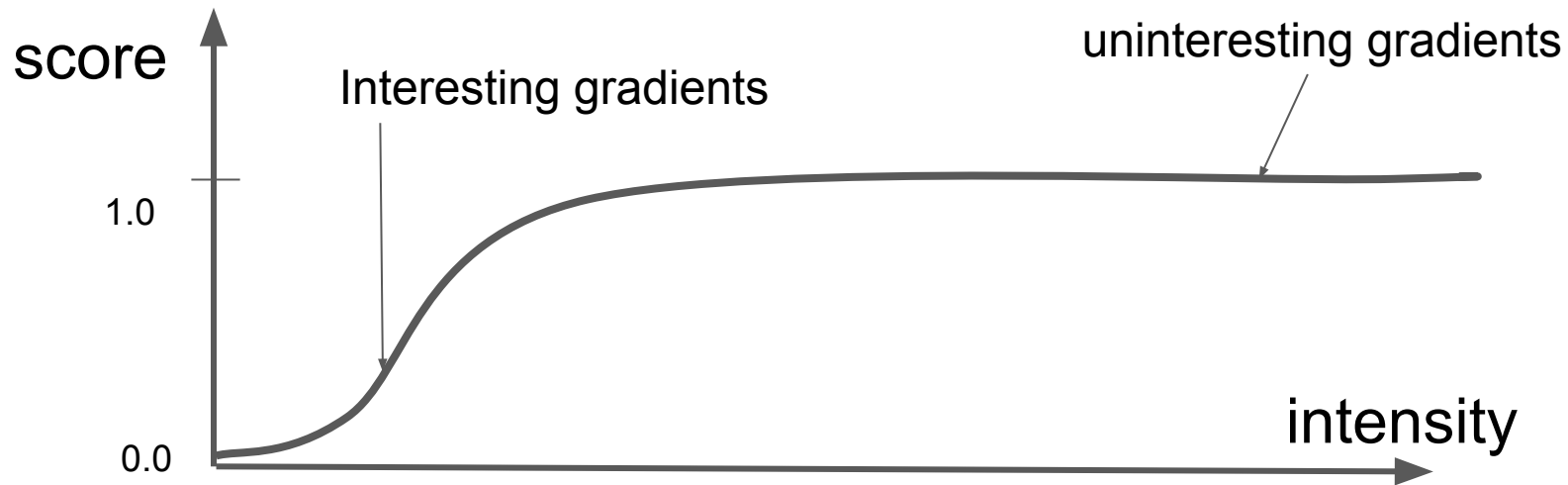
[**DeepLift** '17], [**Layerwise relevance propagation** '16], [**DeConv nets** '14], [**Guided backprop** '14]
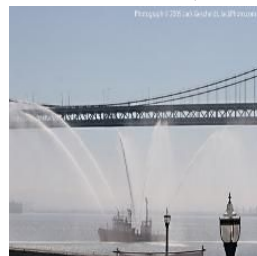
# Prior work

- Prior work is characterized by intuitive design+empirical eval

- Challenges with evaluation

  - Evals from prior work are specific to object recognition
  - **Hard to separate model behavior, attribution errors, eval artifacts**
    - Attributions may "look" incorrect because of unintuitive network behavior
    - Perturbations could change score for artifactual reasons like if they create a "new object" or a "crazy" input

# Our approach: Axioms

- **What makes for a good attribution method?**

  - Inspired by economics literature [Aumann-Shapley 1974]
- Overview

  - List **desirable criteria (axioms)** for an attribution method
  - Establish a uniqueness result: X is the **only** method that satisfies these desirable criteria
    - X = Integrated Gradients (our method)
  - E.g. For f = x*y + z and x,y,z=1,1,1, the "best" attribution is attr(x)=0.5, attr(y)=0.5, attr(z) = 1.0
- Surprisingly, the resulting method is simple

score

uninteresting gradients

Interesting gradients

1.0

0.0

intensity

Scaled images

Scaled gradients

$$IG(input, base) = (input - baseline) *$$

$$\int_{0\ -1} \nabla F(\alpha * input + (1-\alpha) * baseline)\ d\alpha$$

# Easy to implement

```python
def integrated_gradients(inp, baseline, label, steps=range(50)):
    t_input = input_tensor()  # input tensor
    t_prediction = prediction_tensor(label) # output tensor
    t_gradients = tf.gradients(t_prediction, t_input)[0] # gradients
    path_inputs = [baseline + (i/steps)*(inp-baseline) for i in steps]
    grads = run_network(t_gradients, path_inputs)
    return (inp-baseline)*np.average(grads, axis=0)  # integration
```

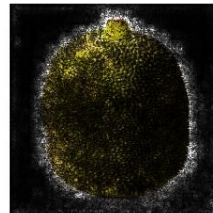Original image — Top label: stopwatch — Score: 0.998507 — Integrated gradients

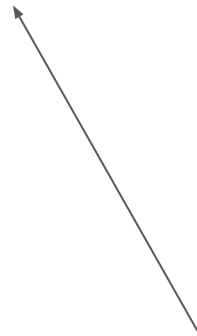Original image — Top label: jackfruit — Score: 0.99591 — Integrated gradients

Original image — Top label: school bus — Score: 0.997033 — Integrated gradients

Attributions overlaid on the image

# Baselines

- **A baseline is an uninformative input** used for comparison
  - E.g. Object recognition: black image, noise image
  - E.g.  For f = x*y + z, and attribution at x,y,z=1,1,1, natural baseline is 0,0,0

- Need for a baseline noted by [Deeplift '16], [Layerwise relevance propagation '16]
- Baselines are necessary for attributions [Kahneman-Miller 86]

# Axiom: Implementation Invariance

**Two networks that compute identical functions <u>for all inputs</u> get identical attributions even if their architecture/parameters differ**

E.g. f1 = x*y + z and f2 = y*x + z should get the same attributions
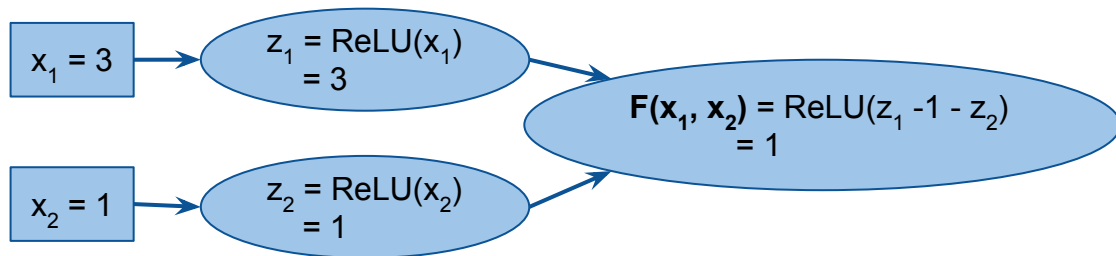
Not satisfied by

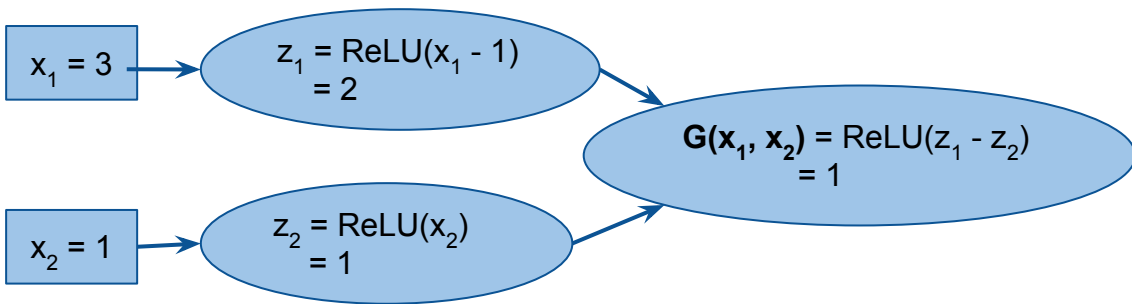- [Deeplift '16], [Layerwise relevance propagation '16]

# Failure of Implementation Invariance

- [**DeepLift** '17], [**Layerwise relevance propagation** '16] fail implementation invariance
- Methods redistribute prediction from output to the input over network structure
  - Chain rule is invalid for discrete gradients

*For all $x_1$ and $x_2$: $F(x_1, x_2) == G(x_1, x_2)$*



$x_1 = 3$ → $z_1 = ReLU(x_1) = 3$

$x_2 = 1$ → $z_2 = ReLU(x_2) = 1$

$F(x_1, x_2) = ReLU(z_1 - 1 - z_2) = 1$

| Integrated gradients | $x_1 = 1.5$, $x_2 = -0.5$ |
|---|---|
| DeepLift | $x_1 = 1.5$, $x_2 = -0.5$ |
| LRP | $x_1 = 1.5$, $x_2 = -0.5$ |

$x_1 = 3$ → $z_1 = ReLU(x_1 - 1) = 2$

$x_2 = 1$ → $z_2 = ReLU(x_2) = 1$

$G(x_1, x_2) = ReLU(z_1 - z_2) = 1$

| Integrated gradients | $x_1 = 1.5$, $x_2 = -0.5$ |
|---|---|
| DeepLift | $x_1 = 2$, $x_2 = -1$ |
| LRP | $x_1 = 2$, $x_2 = -1$ |

# Axiom: Sensitivity

**(a)  If baseline and input have different scores, but differ in a single variable, then that variable gets some attribution.**

**(b)  If a variable has no influence on a function, then it gets no attribution.**

(a) not satisfied by:

- Gradient at output
- [Deconv nets '14]
- [Guided back propagation '14]

# Axiom: Linearity preservation

**Attributions(α*f1 + ß*f2) = α*Attributions(f1) + ß*Attributions(f2)**

Rationale: Attributions have additive semantics. It is best to respect any existing linear structure.

E.g.  For f = x*y + z, the "optimal" attribution should assign blame independently to 'z' and 'x*y'

# Axiom: Completeness

**sum(attributions) = f(input) - f(baseline)**

**Theorem** [Friedman 2004]

Every method that satisfies Linearity preservation, Sensitivity and Implementation invariance, and Completeness is a path integral of a gradient.

# Axiom: Symmetry Preservation

**Symmetric variables with identical values get equal attributions**

E.g. For f = x*y + z, the "optimal" attribution at x,y,z=1,1,2 should be equal for x and y.

**Theorem:** [this work]

Integrated Gradients is the unique path method that satisfies these axioms. (there are other methods that take an average over a symmetric set of paths)
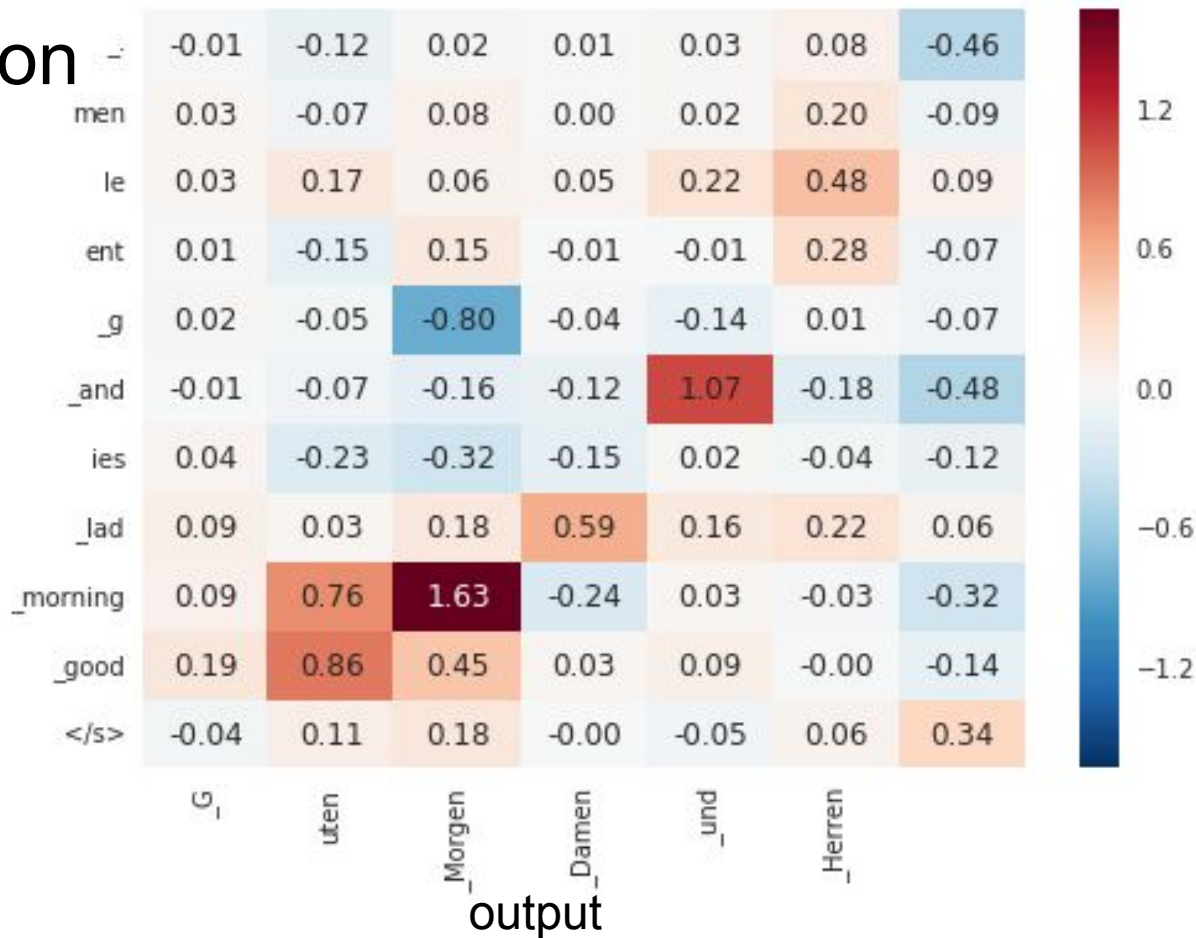
# Virtual Drug Screening



RED: +ve attribution
BLUE: -ve attribtion

**Baseline is zero embedding vector**

# Machine Translation

Attention **does not account** for all the influence of input on output



| input \ output | G¡ | uten | Morgen¡ | Damen¡ | und¡ | Herren¡ |  |
|---|---|---|---|---|---|---|---|
| _ | -0.01 | -0.12 | 0.02 | 0.01 | 0.03 | 0.08 | -0.46 |
| men | 0.03 | -0.07 | 0.08 | 0.00 | 0.02 | 0.20 | -0.09 |
| le | 0.03 | 0.17 | 0.06 | 0.05 | 0.22 | 0.48 | 0.09 |
| ent | 0.01 | -0.15 | 0.15 | -0.01 | -0.01 | 0.28 | -0.07 |
| _g | 0.02 | -0.05 | -0.80 | -0.04 | -0.14 | 0.01 | -0.07 |
| _and | -0.01 | -0.07 | -0.16 | -0.12 | 1.07 | -0.18 | -0.48 |
| ies | 0.04 | -0.23 | -0.32 | -0.15 | 0.02 | -0.04 | -0.12 |
| _lad | 0.09 | 0.03 | 0.18 | 0.59 | 0.16 | 0.22 | 0.06 |
| _morning | 0.09 | 0.76 | 1.63 | -0.24 | 0.03 | -0.03 | -0.32 |
| _good | 0.19 | 0.86 | 0.45 | 0.03 | 0.09 | -0.00 | -0.14 |
| </s> | -0.04 | 0.11 | 0.18 | -0.00 | -0.05 | 0.06 | 0.34 |

# Reductiveness of "Feature Attribution"

- Suppose that the network predicts TRUE if there are **an even number** of black pixels and FALSE otherwise.

- The attributions to the base features will not tell us much.

# Why a vacuum?

# Thanks!

We solve the problem of attributing a network's prediction to its base features

Our technique is called Integrated Gradients, is backed by an axiomatic theory, is <u>simple to implement</u>, and applies to a variety of deep networks

https://github.com/ankurtaly/Integrated-Gradients