



Robert Geirhos*, Roland S. Zimmermann*, Blair Bilodeau*, Wieland Brendel°, Been Kim°

TL/DR How reliable are feature visualizations? We investigate this question through the lens of an **adversary**, **empirically**, and **theoretically**. All three perspectives cast doubt on the reliability of feature visualizations: **They can be manipulated**, **don't reflect how natural input is processed** & are **provably unable to reliably predict even simple function behavior**.

Motivation

Feature visualization is a **foundational interpretability** tool. But are feature visualizations **reliable**, i.e. can we trust & rely on them?

We study this question from **three perspectives**: **Adversarial**, **Empirical**, and **heoretical**.

Summary

- Adversarial Perspective:** Feature visualizations can be fooled by manipulating a model.
- Empirical Perspective:** Even if the model is not manipulated, feature visualizations are processed largely along different paths compared to natural images, which means that they do not explain how neural networks process natural images.
- Theoretical Perspective:** Feature visualizations through activation maximization cannot be used to understand (i.e., predict the behavior of) black-box systems – instead, strong assumptions about the system are necessary.

Conclusion: Feature visualization is best used for exploration / hypothesis generation, not for reliability or confirmation.
Potential way forward: Incorporate more structure & assumptions into networks (instead of seeking black-box explanations). More work needed!

Adversarial Perspective

Through modifications through the network architecture, we can arbitrarily change feature visualizations while maintaining identical behavior on natural input.

Original visualization (last layer, unit #0)

Manipulated visualization (last layer, unit #0)

Identical behavior (e.g., ImageNet accuracy)

A "fooling circuit" manipulates visualizations:

with fooling circuit?

no

yes

unit 0

100

200

300

400

500

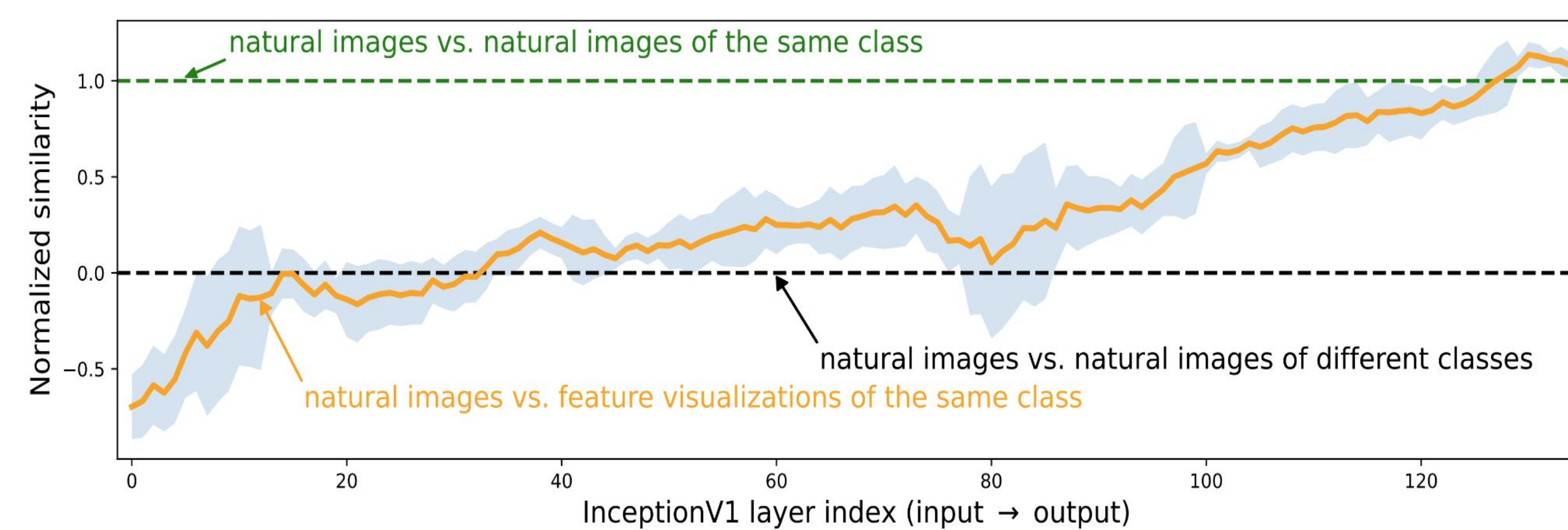
600

700

800

900

Empirical Perspective



Intuition: Feature visualization aims to explain processing of natural input.

This is **only possible if visualizations are processed along similar paths** as natural input.

Empirically, that's not the case: Inception-V1 last-layer visualizations are processed along very different paths compared to natural input throughout the first 2/3rds of network layers.

Theoretical Perspective

| | Given feature visualization for f , can we reliably predict $f(x)$... | | | |
|-----------------------------|--|----------|------------------------|-----------------------|
| | | exactly? | ϵ -approxim.? | closer to min or max? |
| black-box | \mathcal{F} | No | No | No |
| neural network (NN) | \mathcal{F}_{NN} | No | No | No |
| NN trained with ERM | \mathcal{F}_{ERM} | No | No | No |
| L -Lipschitz (known L) | \mathcal{F}_{Lip}^L | No | No | Only for small L |
| piecewise affine | \mathcal{F}_{PAff} | No | No | No |
| monotonic | \mathcal{F}_{Mono} | No | No | No |
| convex | \mathcal{F}_{Conv} | No | No | No |
| affine (input dim. > 1) | $\mathcal{F}_{Aff}^{d>1}$ | No | No | No |
| affine (input dim. = 1) | $\mathcal{F}_{Aff}^{d=1}$ | Yes | Yes | Yes |
| constant | \mathcal{F}_{Const} | Yes | Yes | Yes |

Intuition: How much can you predict about a complex function from knowing its arg max? **Not much ...**

We prove that **feature visualization** based on activation maximization **cannot be used to understand** (i.e., predict meaningful properties of) **a function** unless very strong additional knowledge about the function is available.