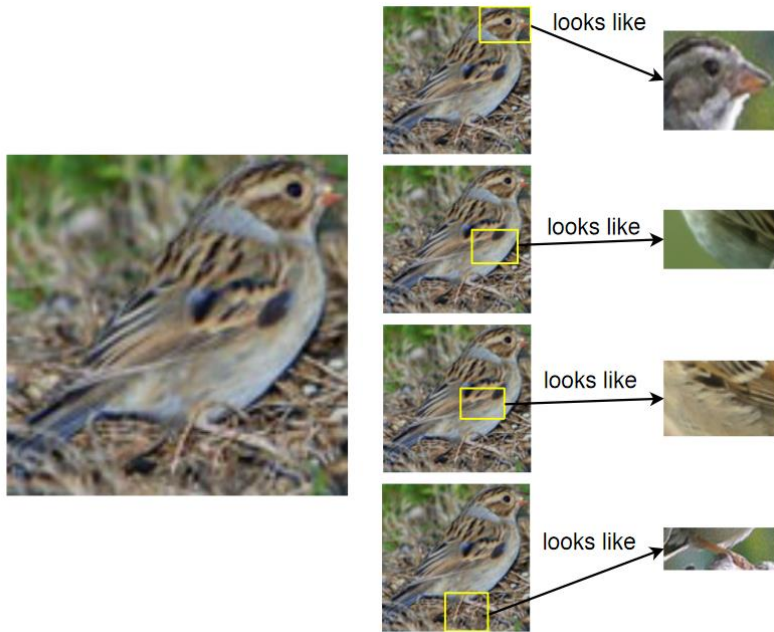# This Looks Like That... Does it?
# Shortcomings of Latent Space Prototype Interpretability in Deep Networks

## Are Deep Interpretable Networks Always Interpretable?
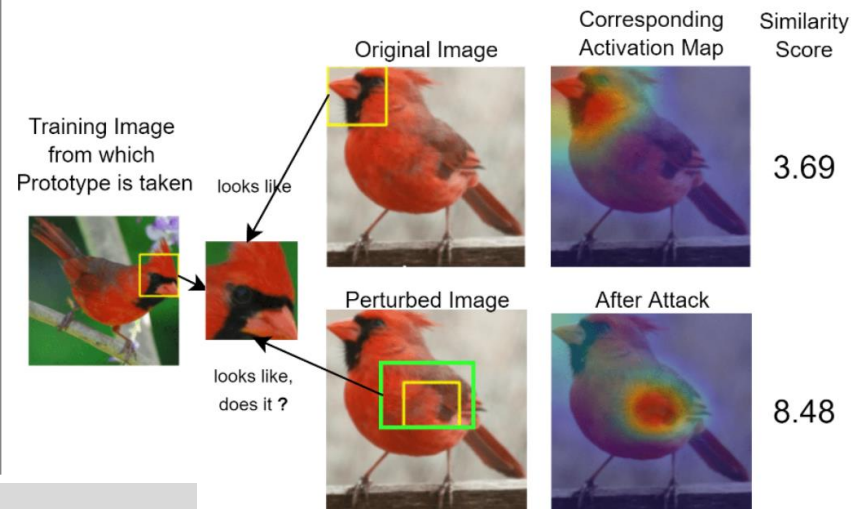
### Prototypical Part Networks



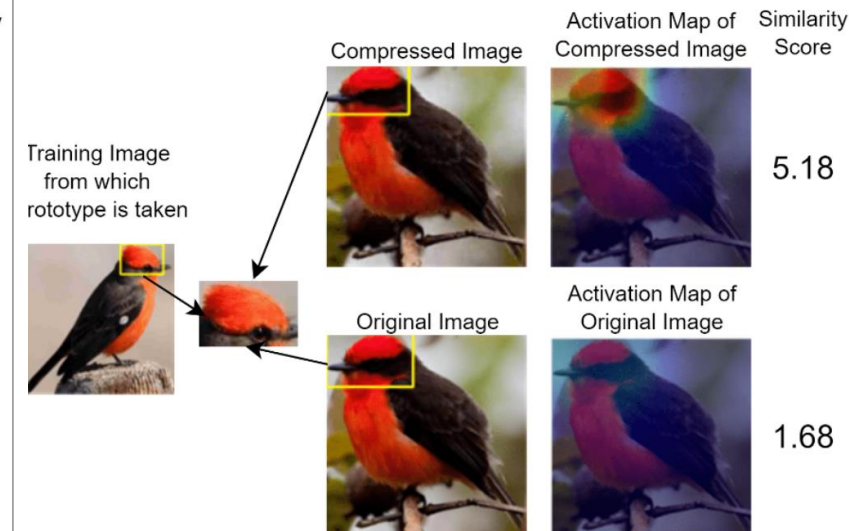## Defining Interpretability for ProtoPNets

Interpretablility holds i.f.f.

*Two image patches look similar to a ProtoPNet* ⇔ *Two image patches look similar to a human*

**not always true**

### Head-on-Stomach Experiment



### JPEG Experiment



## Key Takeaways:

- Use interpretable architectures like ProtoPNets, **with caution**
  - Interpretations might not faithfully convey the network's underlying reasoning process
  - Be aware of different types of noise in dataset when collecting data from different sources

Adrian Hoffmann*, Claudio Fanconi*, Rahul Rade*, Jonas Kohler,  ETH Zurich, Switzerland

**ETH**zürich