# NeuralDivergence

**Georgia Tech**

## Exploring and Understanding Neural Networks by Comparing Activation Distributions

### Haekyu Park    Fred Hohman    Polo Chau
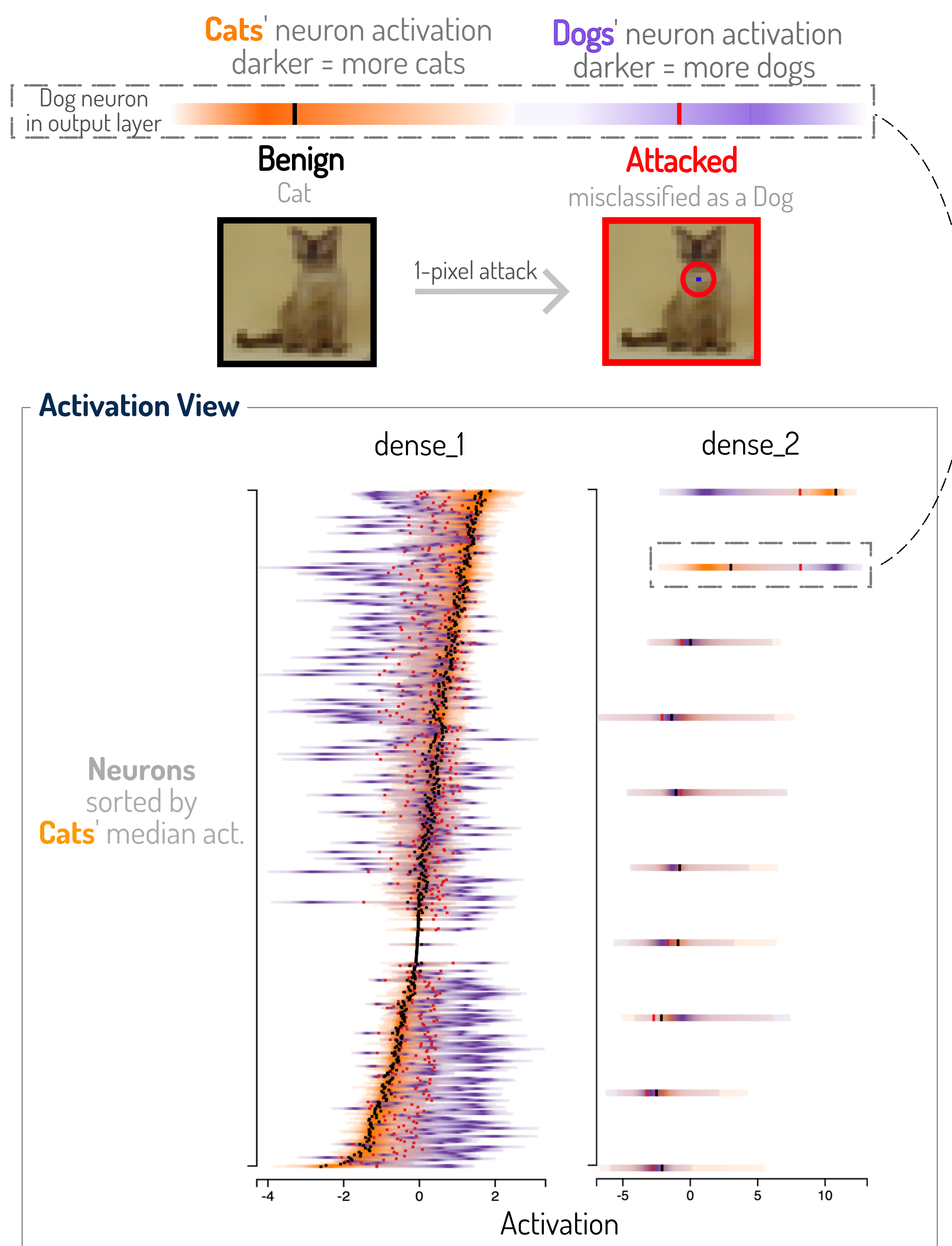
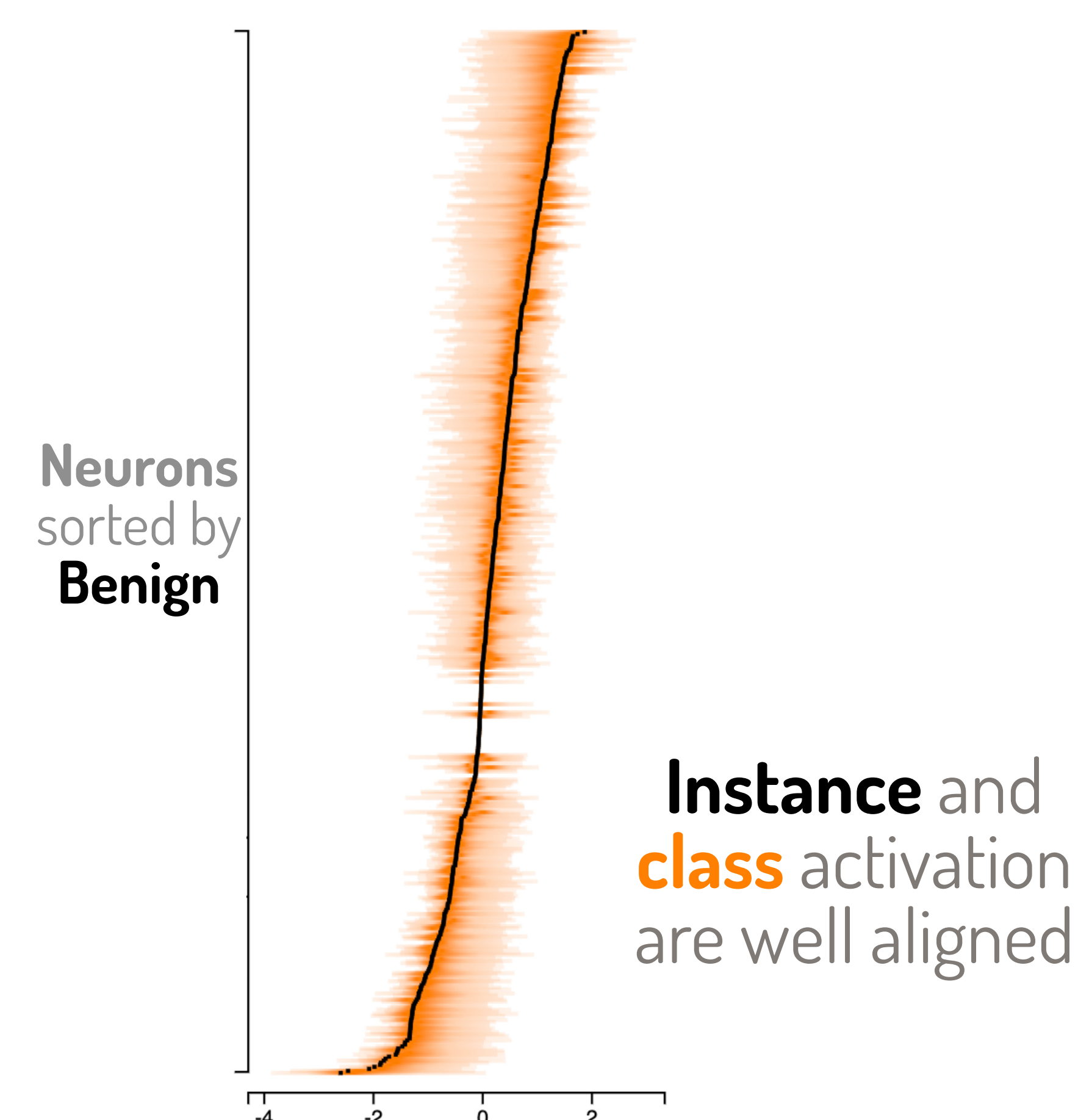haekyu.com / haekyu@gatech.edu

Try at:
**haekyu.com**

## Summary

To **understand** neural networks, NeuralDivergence enables users to explore the models through interactive **summarization** of all neuron activation distribution and **comparison** across layers, classes, and instances (e.g. pairs of adversarial attacked and benign images).

### Scenario: Deciphering Attack on Deep Learning Model

We can use the amount of **"neural divergence"** between an image and its predicted class **to detect one-pixel attack**. The example below shows that an attacked cat image (misclassified as dog) significantly diverges from the "norm" of the real dog class.



**Cats'** neuron activation
darker = more cats

**Dogs'** neuron activation
darker = more dogs

Dog neuron in output layer

**Benign**
Cat

**Attacked**
misclassified as a Dog

1-pixel attack →

**Activation View**

dense_1      dense_2

Neurons sorted by **Cats'** median act.

Activation

**Activation of Cats and Benign**

Neurons sorted by **Benign**

**Instance** and **class** activation are well aligned

**Activation of Dogs and Attacked**

Neurons sorted by **Attacked**

**Instance** and **class** activation **diverge**