# Shared Interest: Large-Scale Visual Analysis of Model Behavior by Measuring Human-AI Alignment
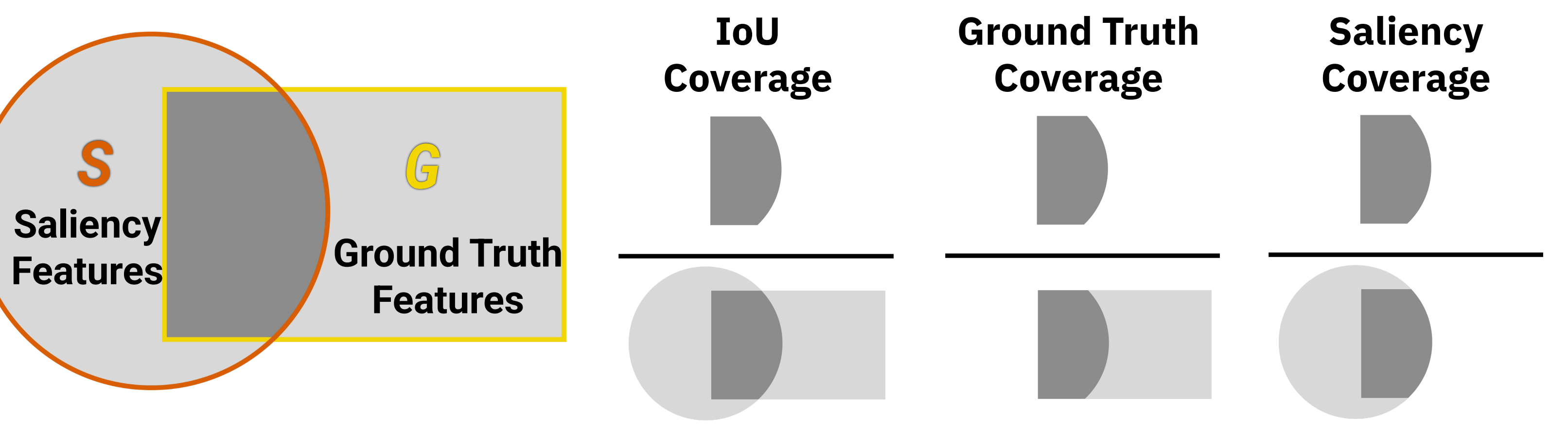
Angie Boggust[1], Benjamin Hoover[2], Arvind Satyanarayan[1], Hendrik Strobelt[2,1]

[1]MIT CSAIL, [2]IBM Research

## Shared Interest Metrics: Computing Human and Model Agreement

To compute human-AI agreement on a data instance, we compute three complementary metrics between the saliency and ground truth features: IoU, Ground Truth, and Saliency Coverage.



**Low Scores:** saliency is disjoint from the ground truth.

**High IoU Coverage:** saliency and ground truth features are identical.

**High Ground Truth Coverage:** all ground truth features are salient to the model.

**High Saliency Coverage:** only ground truth features are salient to the model.

## Identifying Recurring Patterns in Model Behavior

Computing Shared Interest metrics for every instance in a dataset enables us to sort, rank, and aggregate based on model behavior.

We surface common cases that identify dataset limitations and suggest avenues for future model iterations:



## Interactive Probing to Understand Learned Concepts

Shared Interest can also be used as a mechanism to query model behavior.

Given a human annotation, we identify what the model 'knows' about that region by identifying the classes with the highest Shared Interest scores.

Using this procedure, we find our model has not only learned the concept of *dog* (the true label), but also relates the hat annotation to *sombrero*.