# How aligned are different alignment metrics?

Jannis Ahlert, Thomas Klein, Felix A. Wichmann and Robert Geirhos

**1. Many metrics of model-brain alignment yield conflicting results**
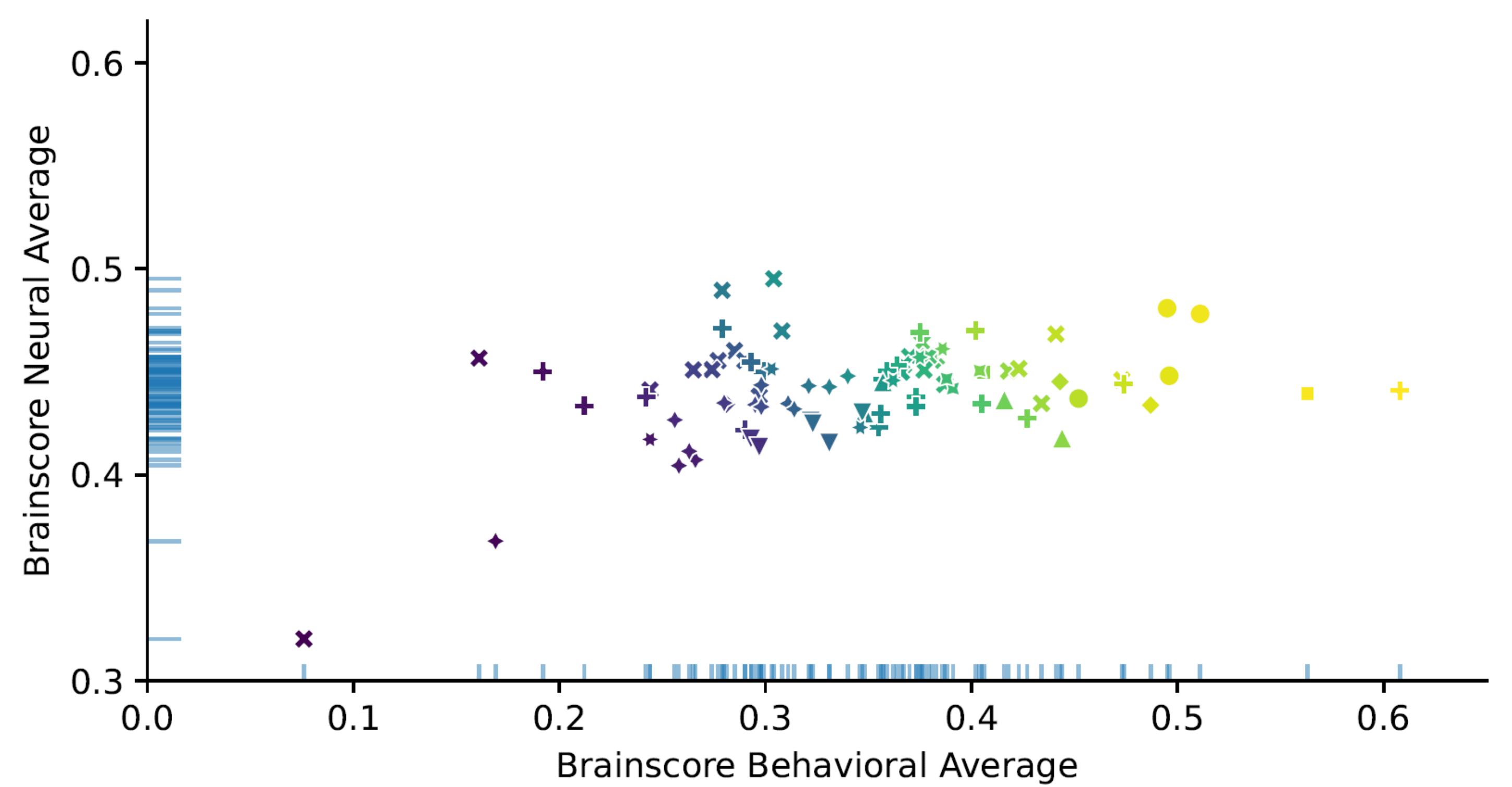
**2. On Brain-Score, behavioral scores have a much larger dynamic range than neural scores**

**3. Integrating metrics should make better use of relationships between benchmarks**

## Motivation

Researchers have put forward many metrics by which the "brain-likeness" of DNNs can be measured.
Benchmark collections such as Brain-Score[1] create overall rankings based on many individual scores.

**Do these metrics agree in their judgements?**
**How many dimensions are there to model-brain alignment?**
**How can we properly integrate multiple metrics into integrated judgements?**



**Comparison of the Neural and Behavioral Average on Brain-Score:** Behavioral scores explain 95% of the variance on the leaderboard. Colour indicates overall rank, with higher ranks being lighter

## Discussion

Which method of aggregating scores is most appropriate?

Desirable properties:
- Independence of irrelevant alternatives
- Accounting for the non-linearity of the measurement scale
- Reflects knowledge about the ventral stream (hierarchy of scores)

## Methods

**Clean scores**
remove one benchmark with >50% missing scores
remove models with any missing scores

**Analysis**
Correlate scores with Spearman's rank correlation for all pairs
Recalculate Brain-Score aggregates for different aggregation techniques

References:
[1] Schrimpf, M., Kubilius, J., Hong, H., Majaj, N. J., Rajalingham, R., Issa, E. B., ... & DiCarlo, J. J. (2018). Brain-score: Which artificial neural network for object recognition is most brain-like?. BioRxiv, 407007.

MAX PLANCK INSTITUTE FOR INTELLIGENT SYSTEMS

imprs-is

EBERHARD KARLS UNIVERSITÄT TÜBINGEN