# Exploring the Platonic Representation Hypothesis Beyond In-Distribution Data

Aryasomayajula Ram Bharadwaj
Independent Researcher
ram.bharadwaj.arya@gmail.com

20th Oct 2024

**Abstract**

The Platonic Representation Hypothesis (PRH) [1] suggests that models trained with different objectives and on various modalities can converge to a shared statistical understanding of reality. This paper explores whether PRH holds when models are exposed to out-of-distribution (OOD) data. We analyze the alignment of model representations using ImageNet-O, a dataset designed with OOD images, and compare these results with those from the original experiment where the evaluation data is within the pretraining data distribution and also with purely random data. Our findings indicate that PRH holds in OOD settings but not for random data, highlighting the importance of underlying structure in data for shared model representations.

## 1    Introduction

The PRH proposes that models converge toward a shared statistical representation of reality, regardless of training specifics such as dataset or training objectives. Previous evidence largely focused on available datasets which are in-distribution with respect to the training data of the models (e.g., Places365's validation dataset). This hypothesis is extended here to include OOD scenarios using ImageNet-O and random noise data. The ImageNet-O dataset, as introduced in [2], contains images that lie outside the distribution of ImageNet, making it a valuable resource for testing generalization.

## 2    Methodology

To test PRH in OOD settings, we used various metrics for measuring model alignment, such as mutual k-NN and CKNN-A. The data was divided into three categories: (1) in-distribution data (ImageNet), (2) OOD data (ImageNet-O), and (3) random noise.

# 3 Results

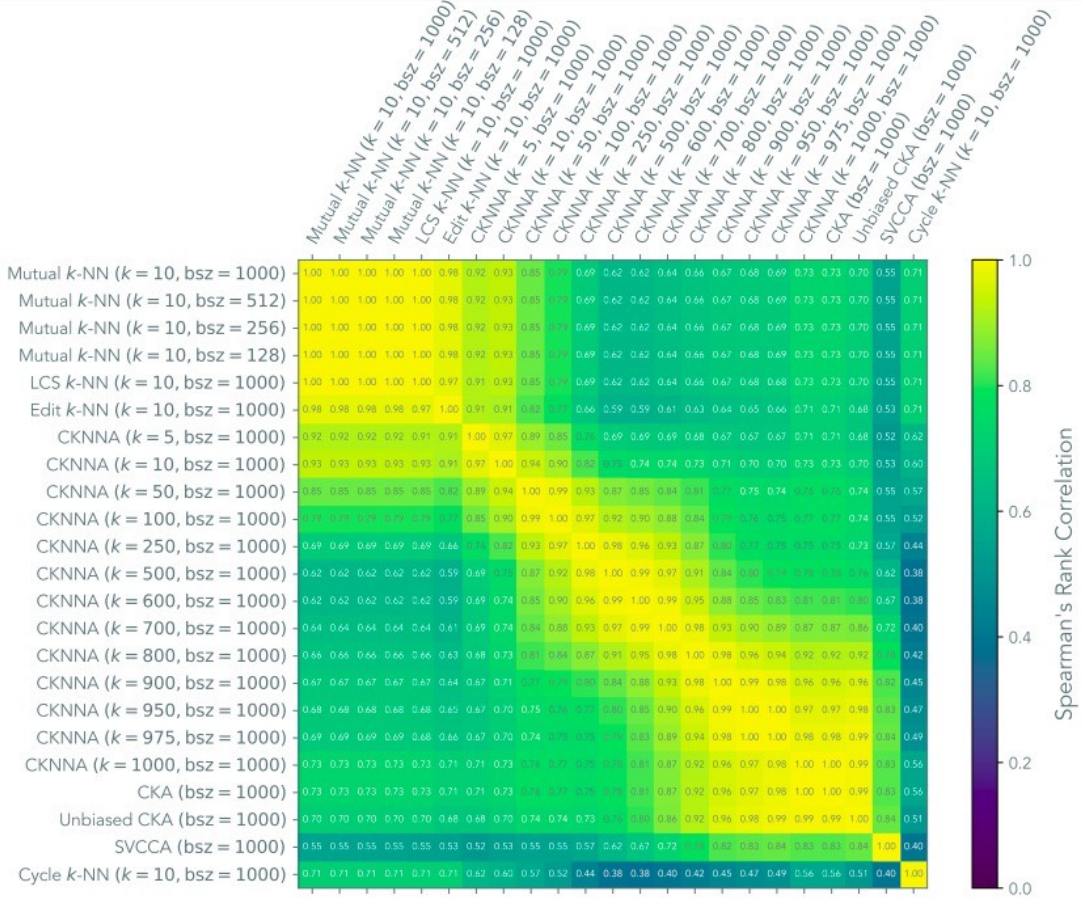## 3.1 Vision-Vision Alignment on Places365's Validation Dataset



*Figure 12.* **Vision-vision alignment measured with various metrics.** Spearman's rank correlation among different metrics and batch sizes (bsz) when used to measure alignment among 78 vision models (see Appendix C.1 for details of these models). All $p$-values are below $2.24 \times 10^{-105}$. Our vision-vision analysis in Figure 2 is based on the first metric (Mutual $k$-NN with $k = 10$ and bsz = 1000).

Figure 1: Vision-vision alignment measured with various metrics. Spearman's rank correlation among different metrics and batch sizes when used to measure alignment among 78 vision models.
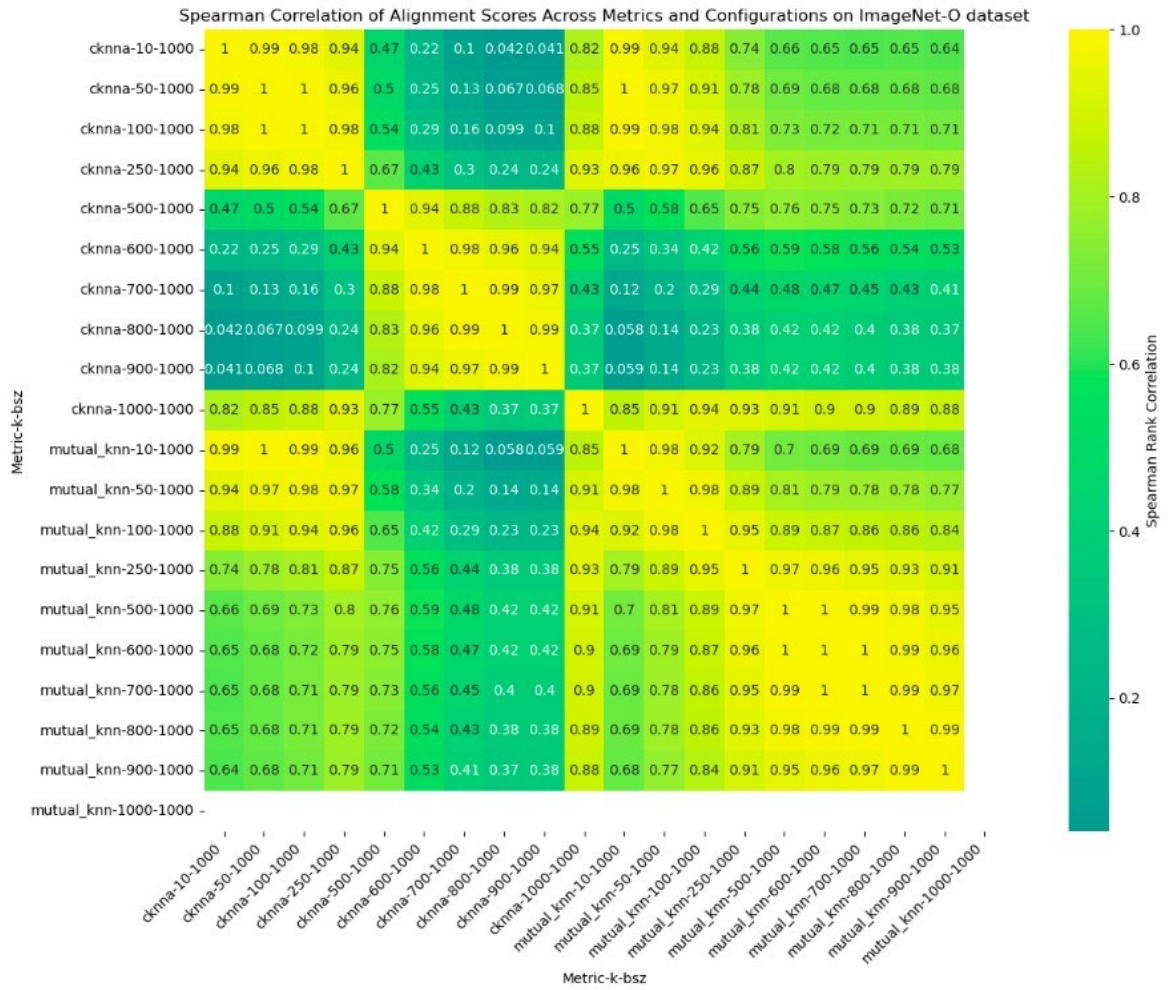
## 3.2 Vision-Vision Alignment on ImageNet-O



Figure 2: Spearman correlation of alignment scores across metrics and configurations on the ImageNet-O dataset.

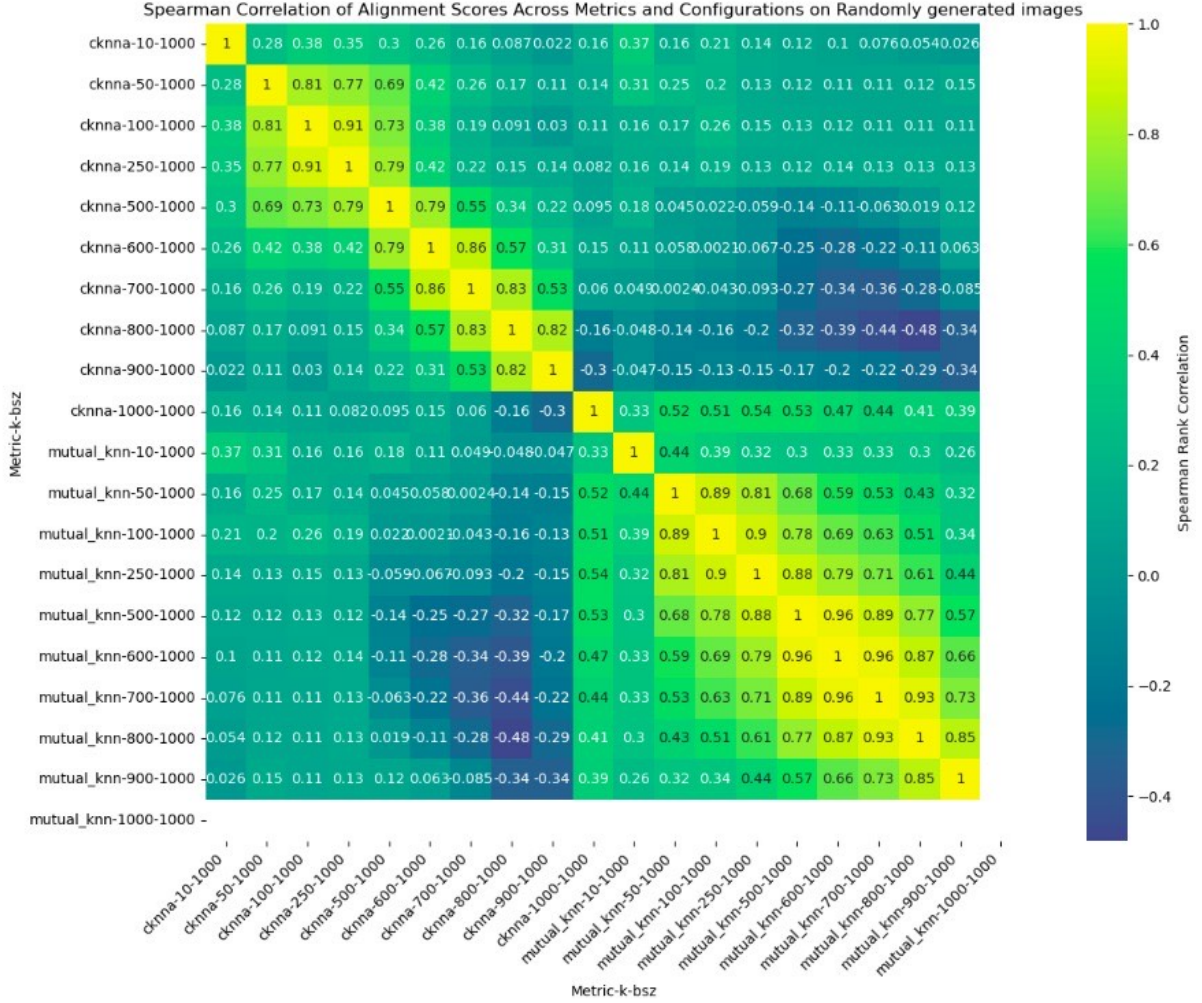## 3.3 Vision-Vision Alignment on Random Data



Figure 3: Spearman correlation of alignment scores across metrics and configurations on randomly generated images.

# 4 Discussion

The results suggest that the Platonic Representation Hypothesis (PRH) holds true even in out-of-distribution (OOD) settings, as demonstrated with the ImageNet-O dataset. Despite being exposed to outlier data, the models maintain high alignment in their representations, indicating that they might develop a shared underlying structure in their perception of reality, even when the training data distribution shifts. However, this alignment breaks down when models are exposed to purely random noise, revealing the limits of the PRH. This distinction is intriguing: on randomly generated data, models do not converge on a shared statistical understanding of reality. Yet, on OOD data like ImageNet-O, the models exhibit strong correlation in their representations, even though their predictions are consistently incorrect but made with high confidence. This suggests that, while the models fail on outlier data, their failure is systematic and predictable.

# 5 Conclusion

We have shown that PRH extends beyond in-distribution data, with models maintaining alignment even in OOD settings. However, the presence of underlying structure in data is crucial for such alignment. Future work should focus on developing better methods to test the generalizability of PRH, finding scenarios where the hypothesis might be falsified, and exploring cross-modality applications of PRH.

The code used for the experiments and analysis in this paper is available on GitHub at https://github.com/rokosbasilisk experiments.

# References

[1] Huh, M., Cheung, B., Wang, T., & Isola, P. (2024). The Platonic Representation Hypothesis. *International Conference on Machine Learning.*

[2] Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., & Song, D. (2021). Natural Adversarial Examples. *CVPR.* arXiv:2405.07987.