# Bridging Critical Gaps in Convergent Learning: How Representational Alignment Evolves Across Layers, Training, and Distribution Shifts

# Chaitanya Kapoor, Sudhanshu Srivastava, Meenakshi Khosla

Department of Cognitive Science, UC San Diego, {chkapoor, sus021, mkhosla}@ucsd.edu

47

77

78

79

80

## Abstract

1

33

Understanding convergent learning—the degree to which 48 2 independently trained neural systems—whether multiple 49 3 artificial networks or brains and models—arrive at sim-<sup>50</sup> 4 ilar internal representations—is crucial for both neuro-<sup>51</sup> 5 science and Al. Yet, the literature remains narrow in 52 6 scope—typically examining just a handful of models with 53 7 one data distribution, relying on one alignment metric,<sup>54</sup> 8 and evaluating networks at a single post-training check-55 9 point. We present a large-scale audit of convergent 56 10 learning, spanning dozens of vision models and thou-57 11 sands of layer-pair comparisons, to close these long-58 12 standing gaps. First, we pit three alignment families 59 13 against one another-linear regression (affine-invariant), 60 14 orthogonal Procrustes (rotation-/reflection-invariant), and 61 15 permutation/soft-matching (unit-order-invariant). We find 62 16 that orthogonal transformations align representations 63 17 nearly as effectively as more flexible linear ones, and al-64 18 though permutation scores are lower, they significantly 65 19 exceed chance, indicating a privileged representational 66 20 Tracking convergence throughout training fur-67 basis. 21 ther shows that nearly all eventual alignment crystallizes 68 22 within the first epoch—well before accuracy plateaus—69 23 suggesting it is largely driven by shared input statistics 70 24 and architectural biases, not by convergence towards the <sup>71</sup> 25 final task solution. Finally, when models are challenged 26 with a battery of out-of-distribution images, early layers 27 remain tightly aligned, whereas deeper layers diverge in 73 28 proportion to the distribution shift. These findings fill crit-74 29 ical gaps in our understanding of representational con-75 30 vergence, with implications for neuroscience and Al. 76 31

# 32 Keywords: Representational Convergence

#### Introduction

Deep Neural Networks (DNNs) are becoming indispens-81 34 able tools in neuroscience, for both predicting neural 82 35 responses (Yamins et al., 2014), (Yamins & DiCarlo, 83 36 2016), (Khaligh-Razavi & Kriegeskorte, 2014) or as mod-84 37 els for reverse-engineering algorithms of neural computa-85 38 tion (Schrimpf et al., 2018), (Schrimpf et al., 2020), (Cichy 86 39 et al., 2016). This congruence invokes the necessity to gain 87 40 a deep understanding of how DNNs learn to represent infor- 88 41 mation. A core question in this domain is whether-and un-89 42 der what conditions-independently trained networks develop 90 43 similar internal representations, and along which dimensions 91 44 this convergence unfolds. Answering question along these 92 45 lines of inquiry can help illuminate how architectural design, 93 46

training objectives and input statistics give rise to emergent features, and whether any aspects of these representations are universal across architectural choices.

Over the past decade, studies have shown that independently trained models often learn highly aligned feature sets-early convolutional layers reliably produce Gabor-like filters (Yosinski et al., 2014), and deeper layers exhibit substantial correspondence when measured via methods such as CCA and its variants (Morcos et al., 2018; Raghu et al., 2017) CKA (Kornblith et al., 2019), RSA (Mehrer et al., 2020) and model stitching (Bansal et al., 2021). More recently, large-scale models appear to converge even across tasks and data modalities (Moschella et al., 2023), suggesting the existence of a modality-agnostic. "Platonic representation" that transcends specific inputs (Huh et al., 2024). Yet, these findings leave open critical gaps: we lack fine-grained metrics to characterize the minimal transformations needed for alignment; we know little about how and when representations align during training; and we do not understand how robust this convergence is under distributional shifts. In this work, we address these questions by evaluating representational alignment along three key axes-across layers, over the course of training, and under out-of-distribution (OOD) inputs-using a suite of metrics with varying degrees of invariances to reveal the geometry and dynamics of convergence in DNNs.

#### Results

Consider a representational pair  $\mathbf{X}_i \in \mathbb{R}^{M \times N_x}$  and  $\mathbf{X}_j \in \mathbb{R}^{M \times N_y}$  as responses to M unique stimuli, where  $N_x$  and  $N_y$  denote the number of units in each representation. We compare these under three mappings—permutation (and its extension when  $N_x \neq N_y$ , Soft-Matching), Procrustes and Linear regression, by finding the optimal mapping  $\mathbf{M}$  under each transformation. We then report the alignment as a pairwise correlation  $\operatorname{corr}(\mathbf{X}_j, \mathbf{MX}_i)$ . These metrics, ordered from strict to flexible, isolate alignment in tuning functions, geometric shape and information content respectively.

**Convergence with Network Depth:** Across independent seeds of the same architecture, alignment is highest in the early layers—known to extract universal, low-frequency features (*e.g.*, edges, corners) (Rahaman et al., 2019; Bau et al., 2017; Zeiler & Fergus, 2014)—and gradually decreases with network depth, consistent across all three metrics (Fig. 1-A). Allowing more flexible mappings (Permutation  $\rightarrow$  Procrustes  $\rightarrow$  Linear) yields marginal gains beyond Procrustes, indicating that simple rotations capture most of the shared structure and additional degrees of freedom (*e.g.*, scaling, shearing) contribute little. Since Procrustes is symmetric, these results



Figure 1: Representational Alignment Across Network Depth, Training and Distribution Shifts. (A) Alignment score vs. layer depth, showing stricter metrics yield lower scores (Linear > Procrustes > Permutation). (B) Layer-layer alignment using Procrustes (Top) and Soft-Matching (Bottom) with maximums over rows (top) / columns (right) denoted in gray lines. (C) Procrustes alignment over the first 10 ImageNet training epochs (lighter = earlier), with task performance. (D) Layer-wise Procrustes alignment for within-distribution (WD) and 17 OOD datasets. Error bars denote standard error across (n = 17) OOD sets. (E) Correlation between Procrustes alignment and task performance over normalized network depth.

130

highlight that convergent learning reflects a deep geometric<sub>127</sub> 94 similarity in feature organization across networks, not merely128 95 129

the ability of one network to predict another. 96

97

Hierarchical Correspondence holds Across Metrics: Pre-131 98 vious studies have shown that, for architecturally identical net-132 99 works trained from different initializations, the most similar133 100 layer in one network to a given layer in another is the corre-134 101 sponding architectural layer (Kornblith et al., 2019). However,135 102 this finding has primarily been supported using affine-invariant<sup>136</sup> 103 metrics (e.g., CCA, SVCCA). We extend this result by showing137 104 that stricter metrics-such as Procrustes and soft-matching138 105 scores—also reveal the same hierarchical correspondence139 106 (Figs. 1-B), even when comparing networks with different ar-140 107 chitectures. This suggests that the hierarchical alignment of 141 108 representations is a fundamental property of neural networks,142 109 robust to architectural differences. Moreover, our results show143 110 that both representational shape (Procrustes) and neuron-144 111 level tuning (Soft-Matching) follow similar alignment patterns,145 112 reinforcing the consistency of this hierarchical organization<sup>146</sup> 113 across different levels of representational analysis. 147 114 148 115

Evolution of Convergence Over Training: We measure Pro-149 116 crustes alignment between independently trained ImageNet 117 networks over the first ten epochs and found that the bulk of<sup>150</sup> 118 convergence occurs within the first epoch—well before appre-151 119 ciable task accuracy (Fig. 1-C). This rapid early convergence152 120 suggests that shared input statistics, architectural inductive bi-153 121 ases, and initial training dynamics are primary drivers of align-154 122 ment, rather than the final task-specific solution. Such find-155 123 ings stand in contrast to hypotheses such as the contravari-156 124 ance principle (Cao & Yamins, 2021) and task generality hy-157 125 pothesis (Huh et al., 2024), which attribute convergence to158 126

constraints imposed by the final, high-performance solution. These results are consistent with studies on the early training phase that identify a rapid representational reorganization before meaningful task learning begins (Frankle et al., 2020).

Convergence Across Distribution Shifts: We investigate the extent of within-distribution (WD) representational alignment of ImageNet-trained DCNNs under 17 OOD variants of ImageNet (Geirhos et al., 2018) sharing ImageNet's 16 coarse labels by computing layerwise Procrustes scores. Early convolutional layers maintained alignment levels nearly identical to WD stimulus, whereas later layers exhibit amplified divergence under OOD conditions (Fig. 1-D). We attribute this to early layers encoding universal features (e.g., edges, corners) that generalize across distributions, while deeper layers represent task-specific abstractions that are more sensitive to distributional shifts. Further, we find that representational alignment in deeper (but not earlier or middle) layers correlates with OOD classification accuracy (Fig. 1-E). These findings suggest that OOD stimuli can serve as an effective probe to differentiate model architectures in model-brain comparisons, and that selective fine-tuning of later layers can serve as a promising strategy for improving OOD generalization.

#### Discussion

This study fills critical gaps in our understanding of convergent learning, offering a comprehensive analysis of how representational alignment between independently trained networks varies across network depth, training, and distribution shifts. We systematically explored how different alignment metricswith varying levels of transformation invariance-capture representational similarities, providing a more nuanced view of convergent learning than previous work.

# References

212

Bansal, Y., Nakkiran, P., & Barak, B. (2021). Revisiting model<sup>213</sup>
 stitching to compare neural representations. *Advances in*<sup>214</sup>
 *neural information processing systems*, *34*, 225–236.

159

- <sup>163</sup> Bau, D., Zhou, B., Khosla, A., Oliva, A., & Torralba, A. (2017).<sup>216</sup>
- <sup>164</sup> Network dissection: Quantifying interpretability of deep vi-<sup>217</sup>
- sual representations. In *Proceedings of the ieee confer-<sup>218</sup> ence on computer vision and pattern recognition* (pp. 6541–<sup>219</sup>
- ence on computer vision and pattern recognition (pp. 6541–<sup>219</sup>
   6549).
- Cao, R., & Yamins, D. (2021). Explanatory models in<sup>221</sup>
   neuroscience: Part 2–constraint-based intelligibility. arXiv<sup>222</sup>
   preprint arXiv:2104.01489.
- 171 Cichy, R. M., Khosla, A., Pantazis, D., Torralba, A., & Oliva,224
- A. (2016). Comparison of deep neural networks to spatio-<sup>225</sup>
- temporal cortical dynamics of human visual object recogni-<sup>226</sup>
- tion reveals hierarchical correspondence. *Scientific reports*,
   6(1), 27755.
- Frankle, J., Schwab, D. J., & Morcos, A. S. (2020). The<sup>229</sup>
   early phase of neural network training. *arXiv preprint*<sup>230</sup>
   *arXiv:2002.10365*.
- Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wich-<sup>232</sup>
   mann, F. A., & Brendel, W. (2018). Imagenet-trained cnns<sup>233</sup>
   are biased towards texture; increasing shape bias improves<sup>234</sup>
   accuracy and robustness. *arXiv preprint arXiv:1811.12231*.
- Huh, M., Cheung, B., Wang, T., & Isola, P. (2024).
  The platonic representation hypothesis. *arXiv preprint arXiv:2405.07987*.
- Khaligh-Razavi, S.-M., & Kriegeskorte, N. (2014). Deep supervised, but not unsupervised, models may explain it cortical representation. *PLoS computational biology*, *10*(11), e1003915.
- Kornblith, S., Norouzi, M., Lee, H., & Hinton, G. (2019). Sim ilarity of neural network representations revisited. In *Inter- national conference on machine learning* (pp. 3519–3529).
- Mehrer, J., Spoerer, C. J., Kriegeskorte, N., & Kietzmann,
   T. C. (2020). Individual differences among deep neural
   network models. *Nature communications*, *11*(1), 5725.
- Morcos, A., Raghu, M., & Bengio, S. (2018). Insights on
   representational similarity in neural networks with canoni cal correlation. Advances in neural information processing
   systems, 31.
- Moschella, L., Maiorca, V., Fumero, M., Norelli, A., Locatello, F., & Rodolà, E. (2023). Relative representations en-
- able zero-shot latent space communication (2023). *arXiv* preprint arXiv:2209.15430.
- Raghu, M., Gilmer, J., Yosinski, J., & Sohl-Dickstein, J. (2017).
   Svcca: Singular vector canonical correlation analysis for
   deep learning dynamics and interpretability. *Advances in neural information processing systems*, *30*.
- Rahaman, N., Baratin, A., Arpit, D., Draxler, F., Lin, M., Hamprecht, F., ... Courville, A. (2019). On the spectral bias
- of neural networks. In International conference on machine
- <sup>211</sup> *learning* (pp. 5301–5310).

- Schrimpf, M., Kubilius, J., Hong, H., Majaj, N. J., Rajalingham, R., Issa, E. B., ... others (2018). Brain-score: Which artificial neural network for object recognition is most brain-like? *BioRxiv*, 407007.
- Schrimpf, M., Kubilius, J., Lee, M. J., Murty, N. A. R., Ajemian, R., & DiCarlo, J. J. (2020). Integrative benchmarking to advance neurally mechanistic models of human intelligence. *Neuron*, 108(3), 413–423.
- Yamins, D. L., & DiCarlo, J. J. (2016). Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience*, 19(3), 356–365.
- Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., & DiCarlo, J. J. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the national academy of sciences*, 111(23), 8619–8624.
- Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in deep neural networks? *Advances in neural information processing systems*, *27*.
- Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer vision–eccv 2014:* 13th european conference, zurich, switzerland, september 6-12, 2014, proceedings, part i 13 (pp. 818–833).