

Vision Sciences Lab

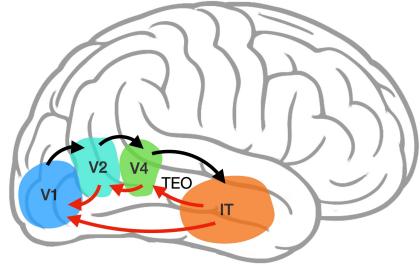
# Representational Geometry Dynamics in Networks After Long-Range Modulatory Feedback

Cindy Kexin Luo, George A. Alvarez, Talia Konkle

Department of Psychology &amp; Kempner Institute, Harvard University

## Introduction

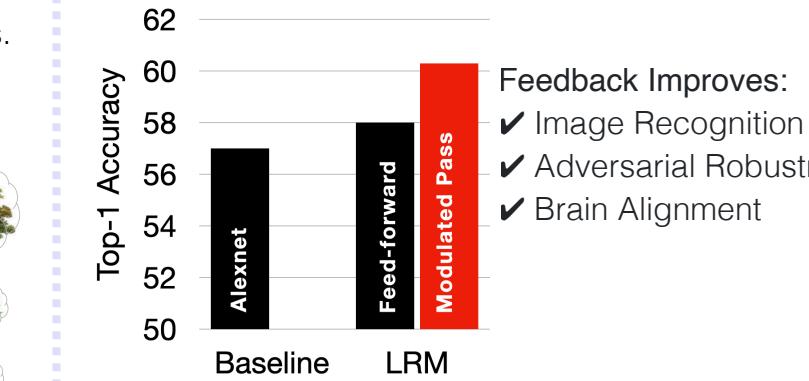
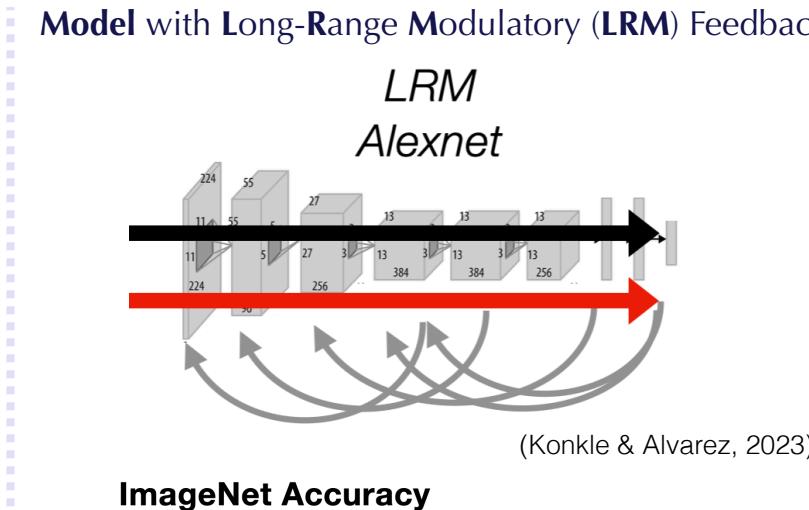
### Feedback in the Visual System



- Extensive feedback connections in the neural system (pathways from IT to V4, V4 to V1).
- Feedforward and recurrent processing iteratively refine perceptual interpretations.

### Roles of Feedback:

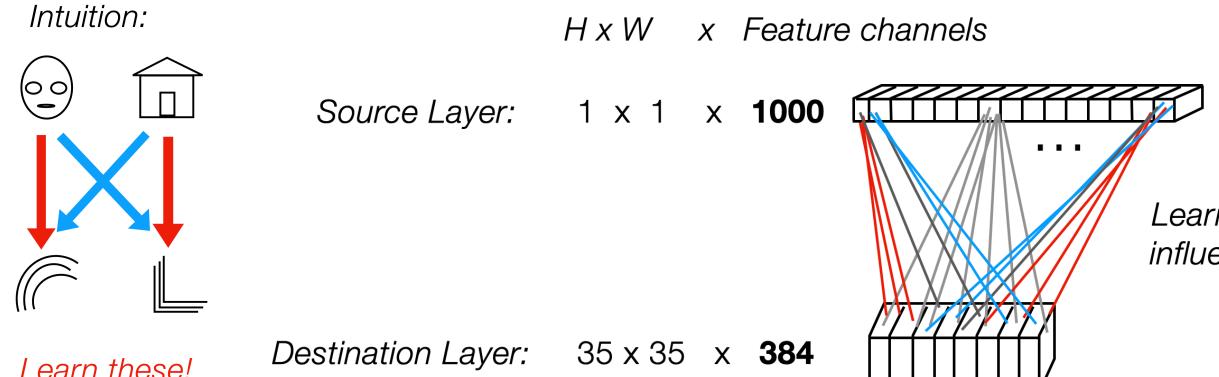
- Emergence of coherent visual objects (Di Lollo, 2012)
- Sharpening representation (Jehee et al., 2007)
- Predictive coding (Rao & Ballard, 1999)
- Functional equivalence to increased depth of processing (Kar et al., 2019)
- Figure-ground segregation (Heinen, Jolij, & Lamme, 2005)
- Construct new representations (Williams et al., 2008)



**Q: How Does Feedback Intrinsically Reshape Representational Geometry?**

## How the LRM Model Integrates Feedback

### Key Design #1: Channel-to-Channel Influences



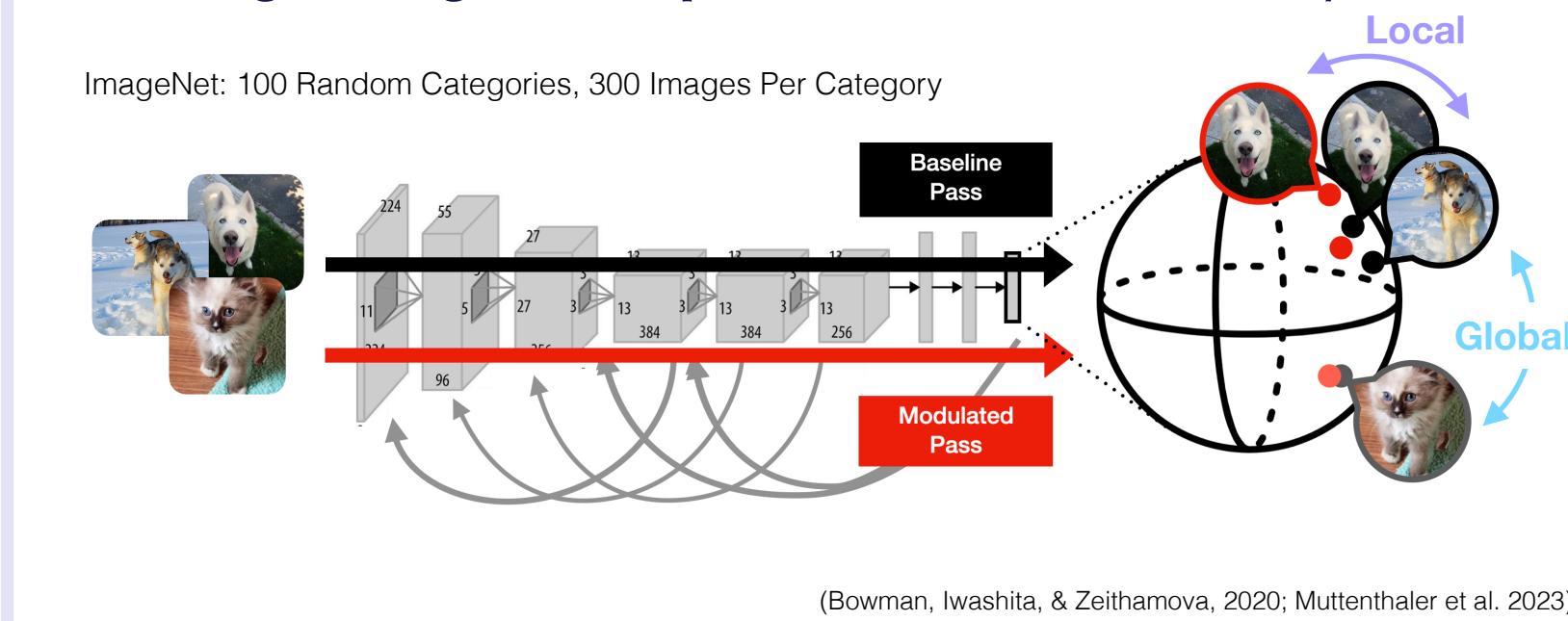
### Key Design #2: Representation Learning

ImageNet: 1000-way categorization

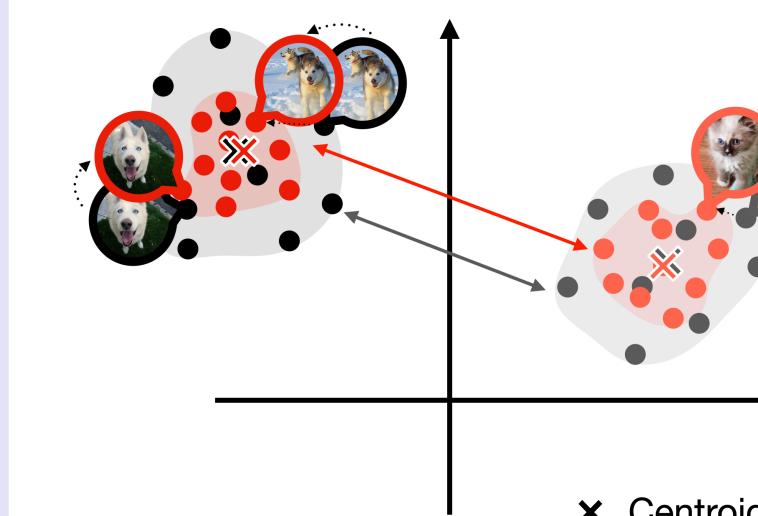
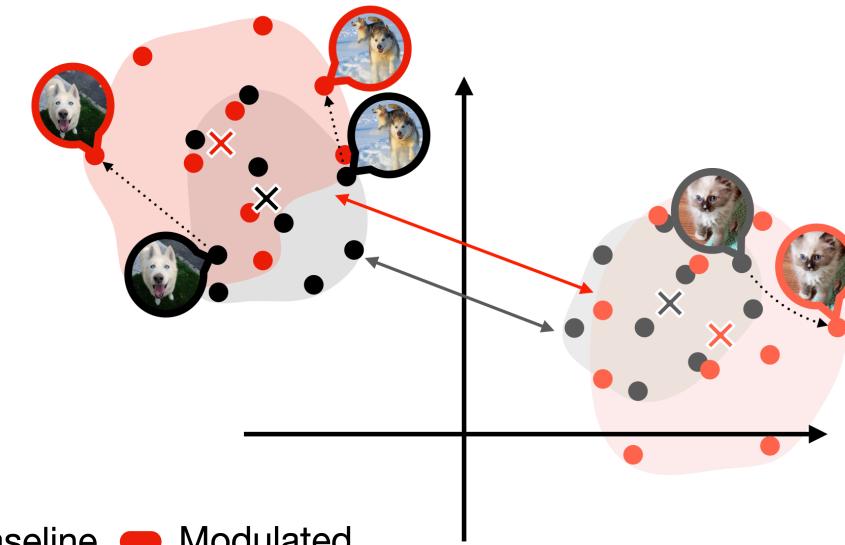
$$\mathcal{L}_{\text{total}} = \frac{\mathcal{L}_{\text{feedforward}} + \mathcal{L}_{\text{Modulated pass}}}{2}$$

## Probing Changes in Representational Geometry

ImageNet: 100 Random Categories, 300 Images Per Category



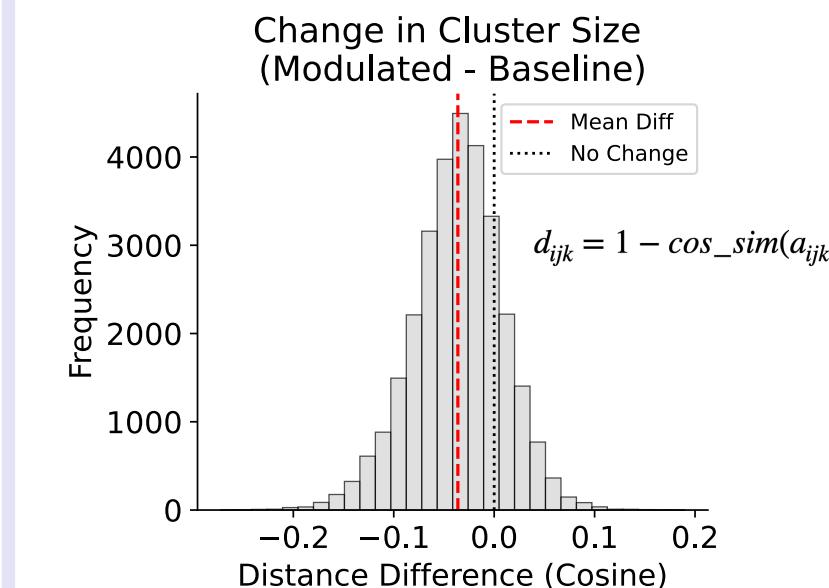
## Hypotheses: How Feedback Reshapes Representations

(a) **Shrink + Stay**(b) **Expand + Shift**

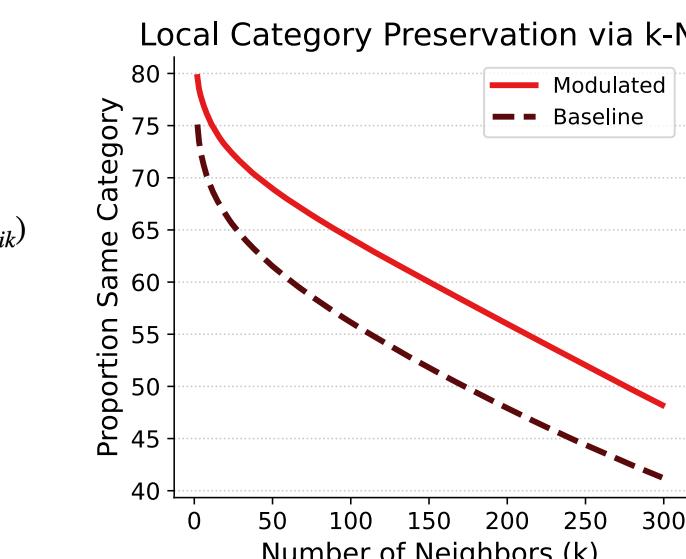
✗ Centroid    ● Baseline    ● Modulated

## Local Effect

### Feedback Compacts Category Clusters



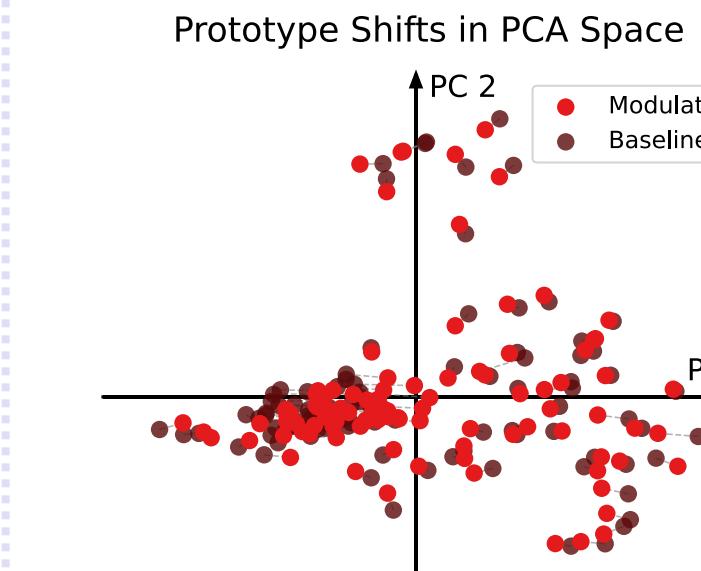
→ Exemplars move significantly closer to their category centroids.



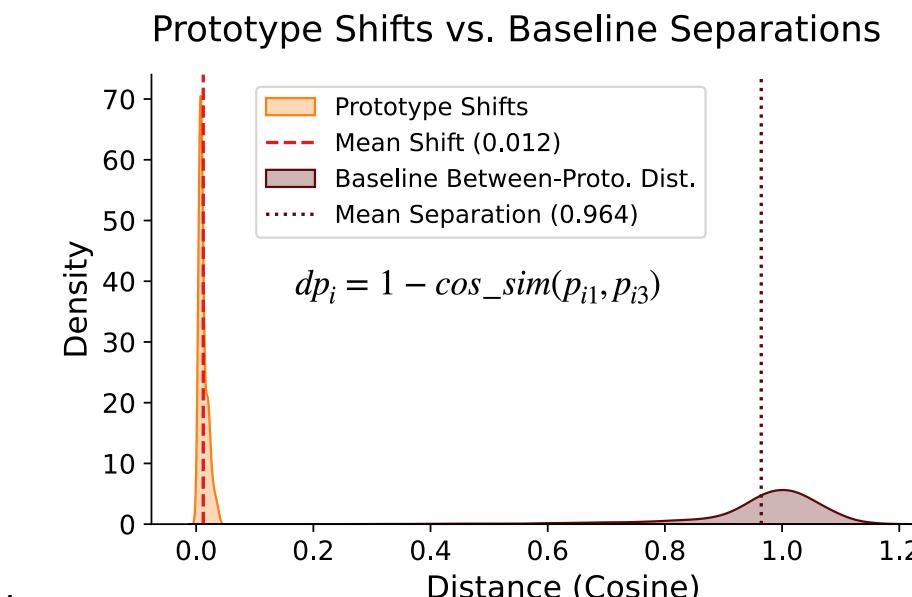
→ Local neighbors have more of the same category after modulation.

## Global Effect

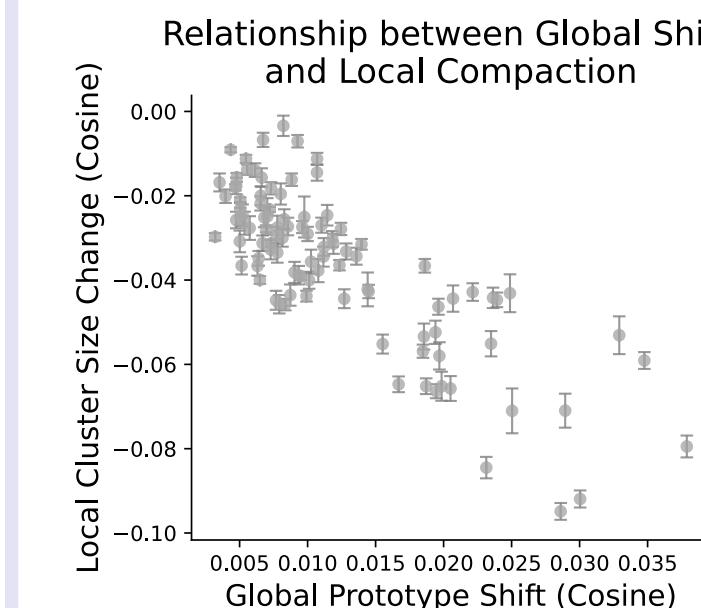
### Overall Category Arrangement Remains Stable

RSA:  $p = 0.95$  between baseline and modulated RDMS.

→ The global structure of inter-category relationships remains largely intact.



## Exploratory Analysis



**Are Categories with Larger Global Shifts Also More Locally Refined?**  
→ Yes!  
A robust correlation ( $r = -.80, p < .001$ ) between local size changes and global prototype shifts.

## Conclusions

Feedback modulation induces an automatic “prototype effect” on representations (Shrink + Stay):

- It compacts and refines local representations within a category.
  - It largely preserves global structure.
- These adjustments may enhance downstream category-based task efficiency.

>>> But what role does the training objective play?  
Extending the analysis to the Self-Supervised LRM—stay tuned!

>>> Does this reflect what happens in the brain?  
We also observed that the LRM model achieved higher brain alignment on the Brain-Score benchmark.  
- Extra-Classical Receptive Field?  
- Changes across layers → Different stages of processing