# Understanding Representational Alignment in Neural Nets
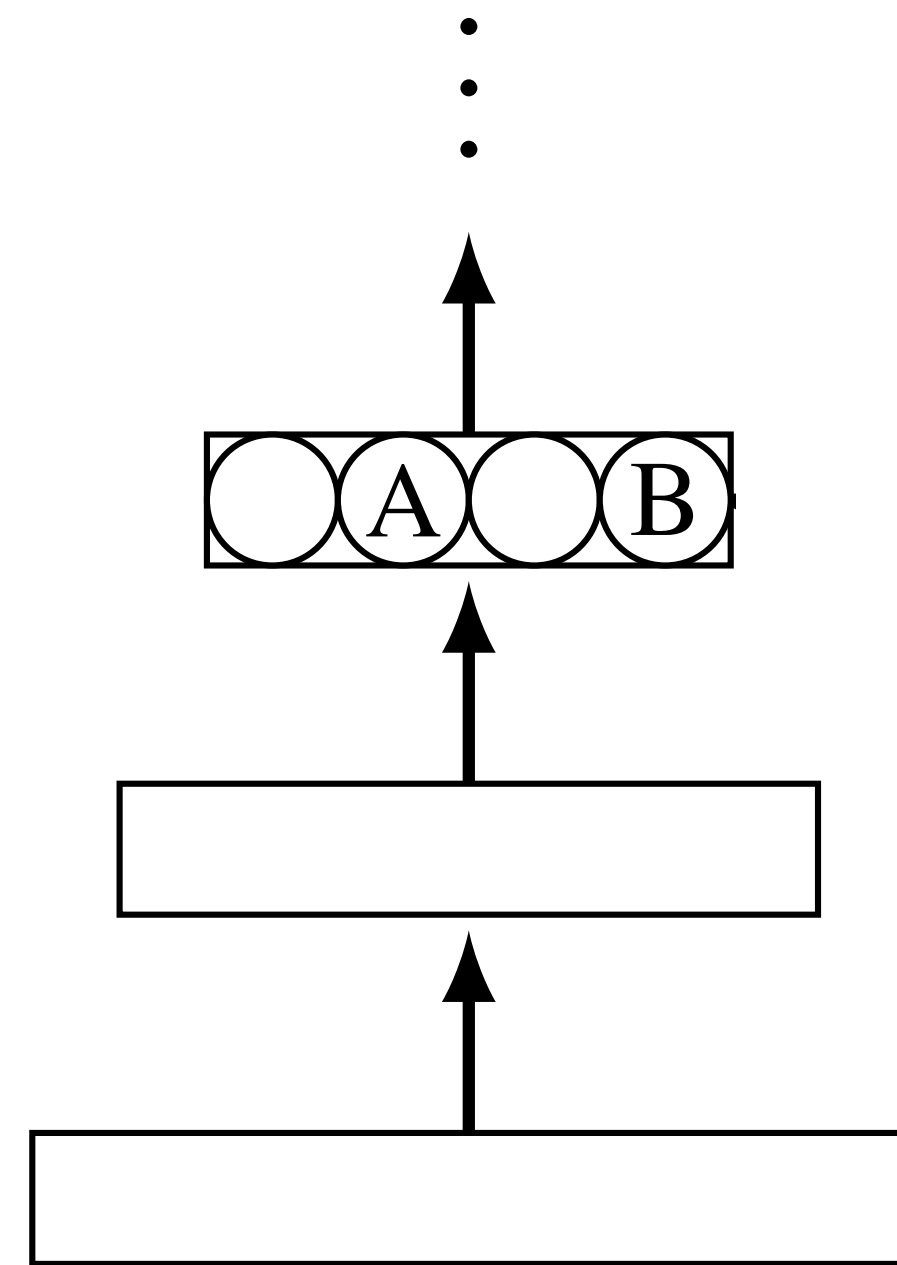
Phillip Isola, MIT
IN2346, TUM
July 15th , 2025

# Object Detectors Emerge in Deep Scene CNNs
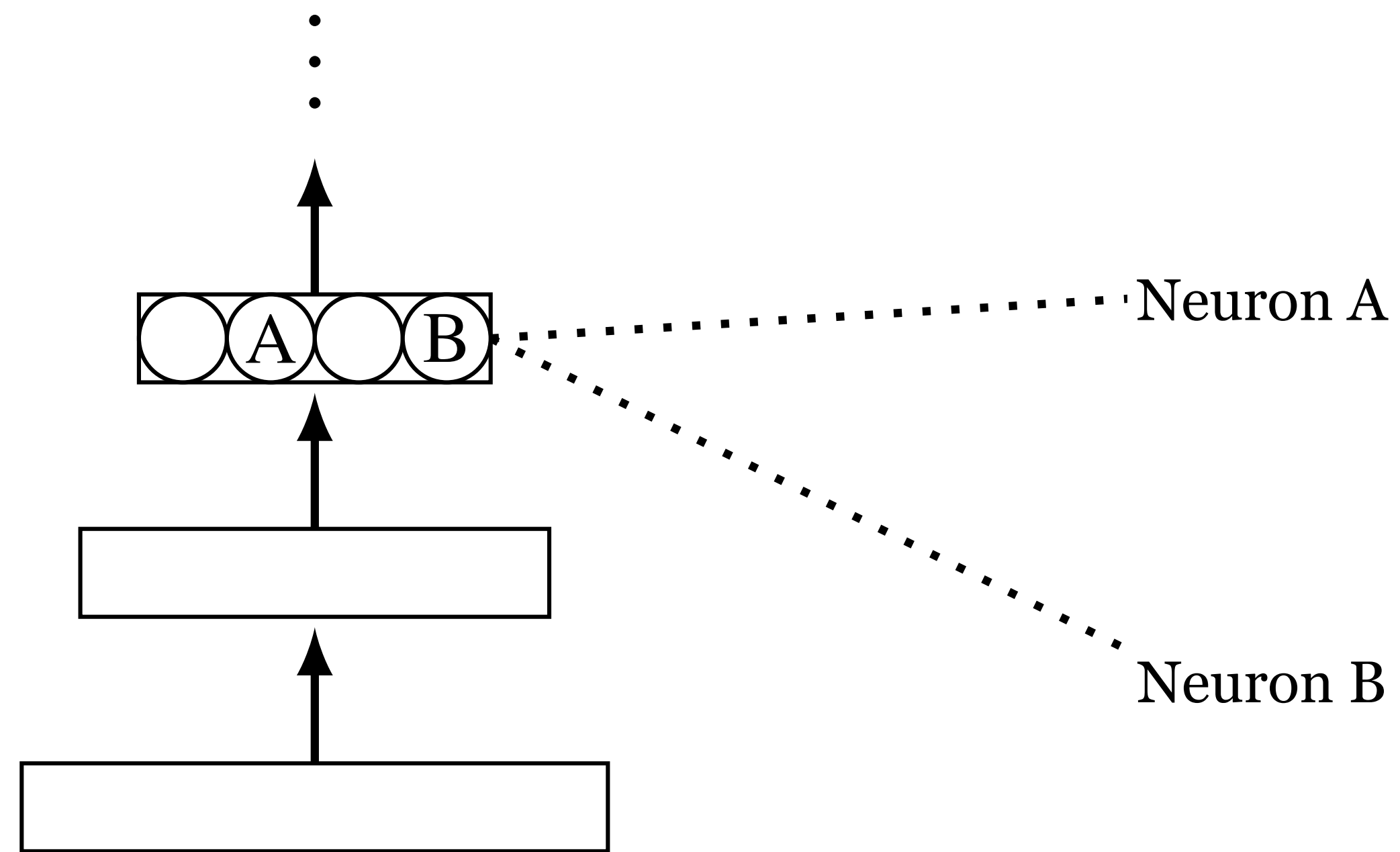
scene



Neuron A

Neuron B

[Zhou, Khosla, Lapedriza, Oliva, Torralba 2015]

[fig modified from: Torralba, Isola, Freeman 2024]

`scene label`

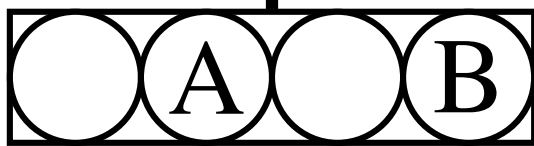Images that maximally activate these neurons



Neuron A

Neuron B

[Zhou, Khosla, Lapedriza, Oliva, Torralba 2015]
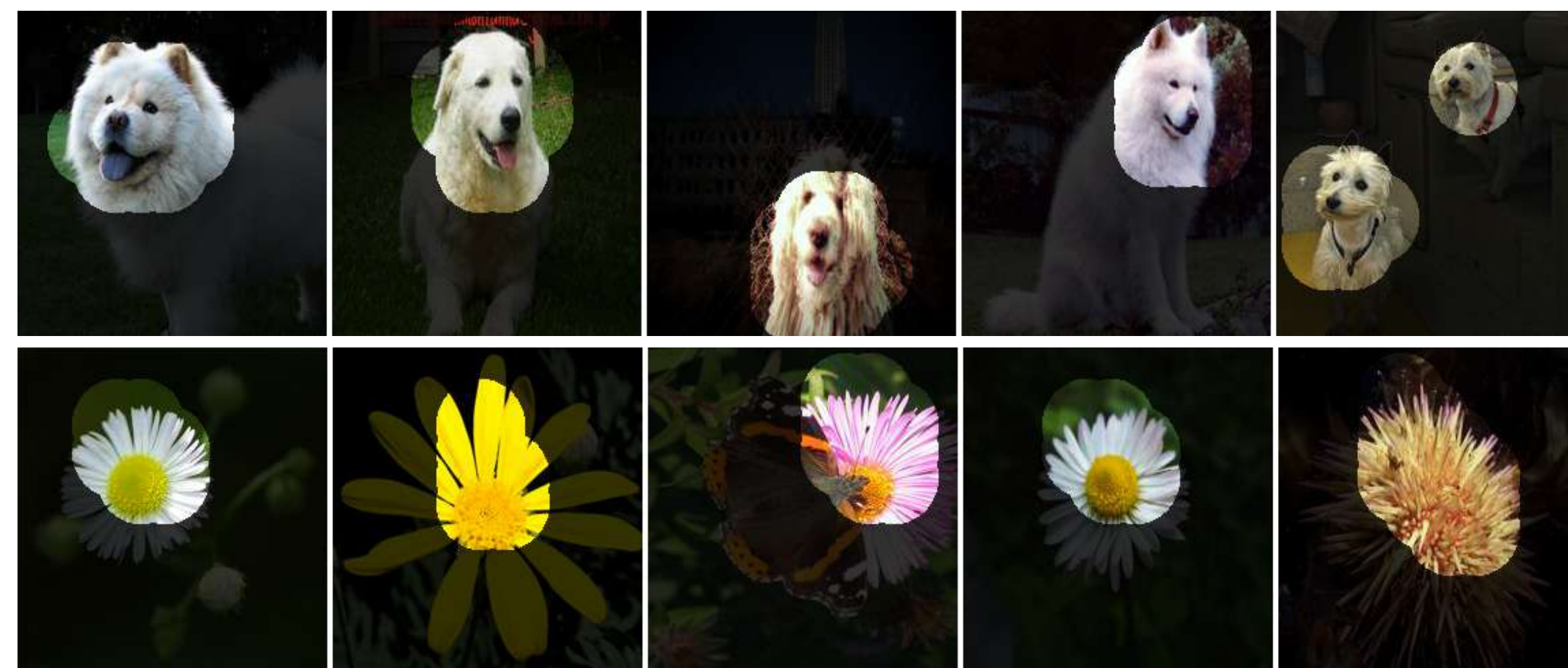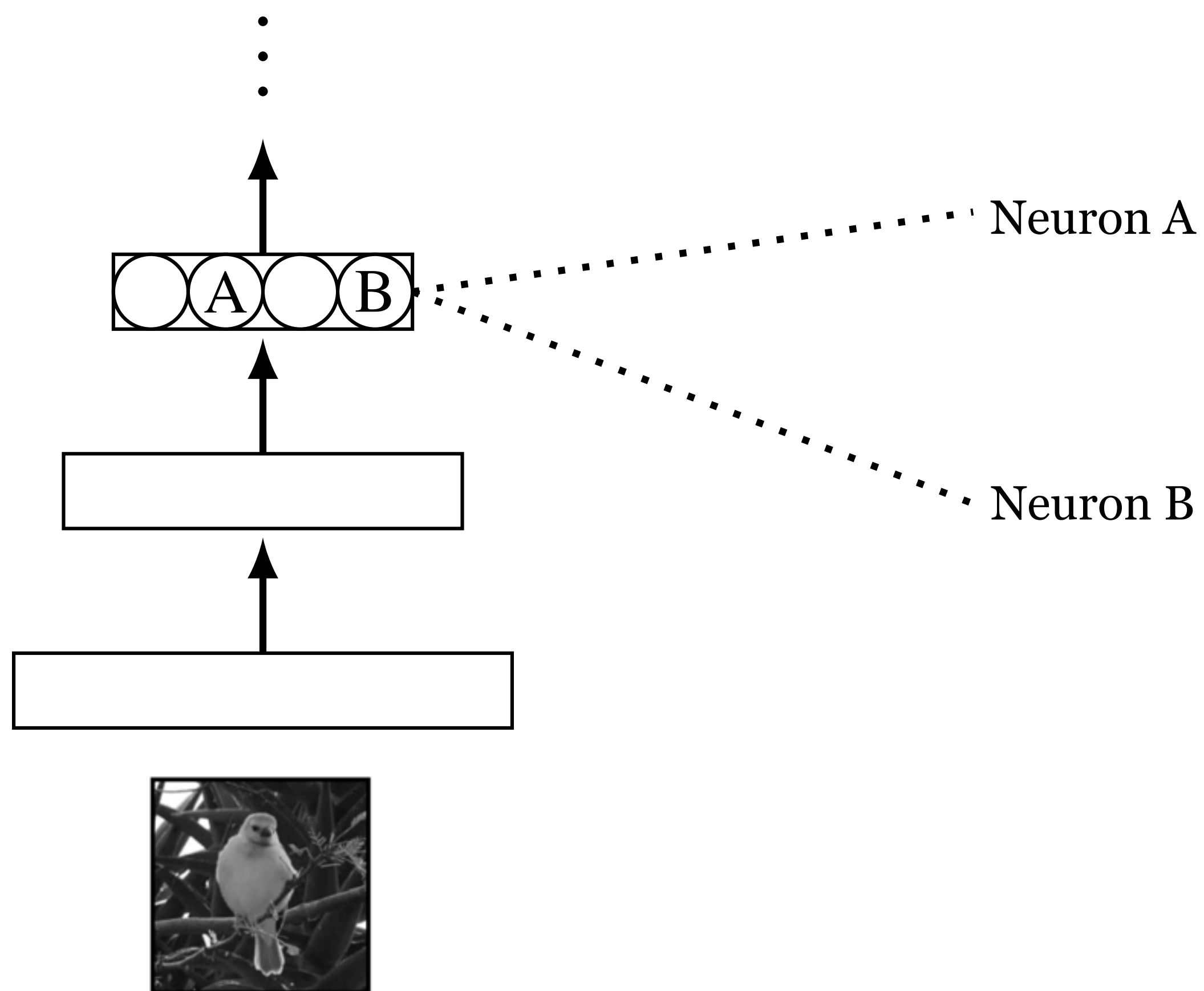
[fig modified from: Torralba, Isola, Freeman 2024]

Neuron A

Neuron B

Images that maximally activate these neurons

Neuron A

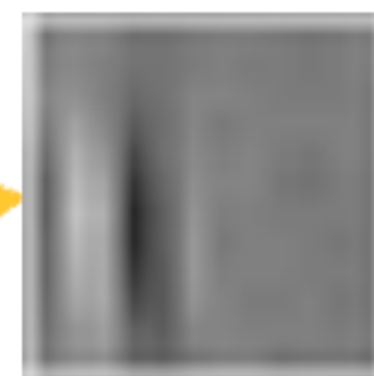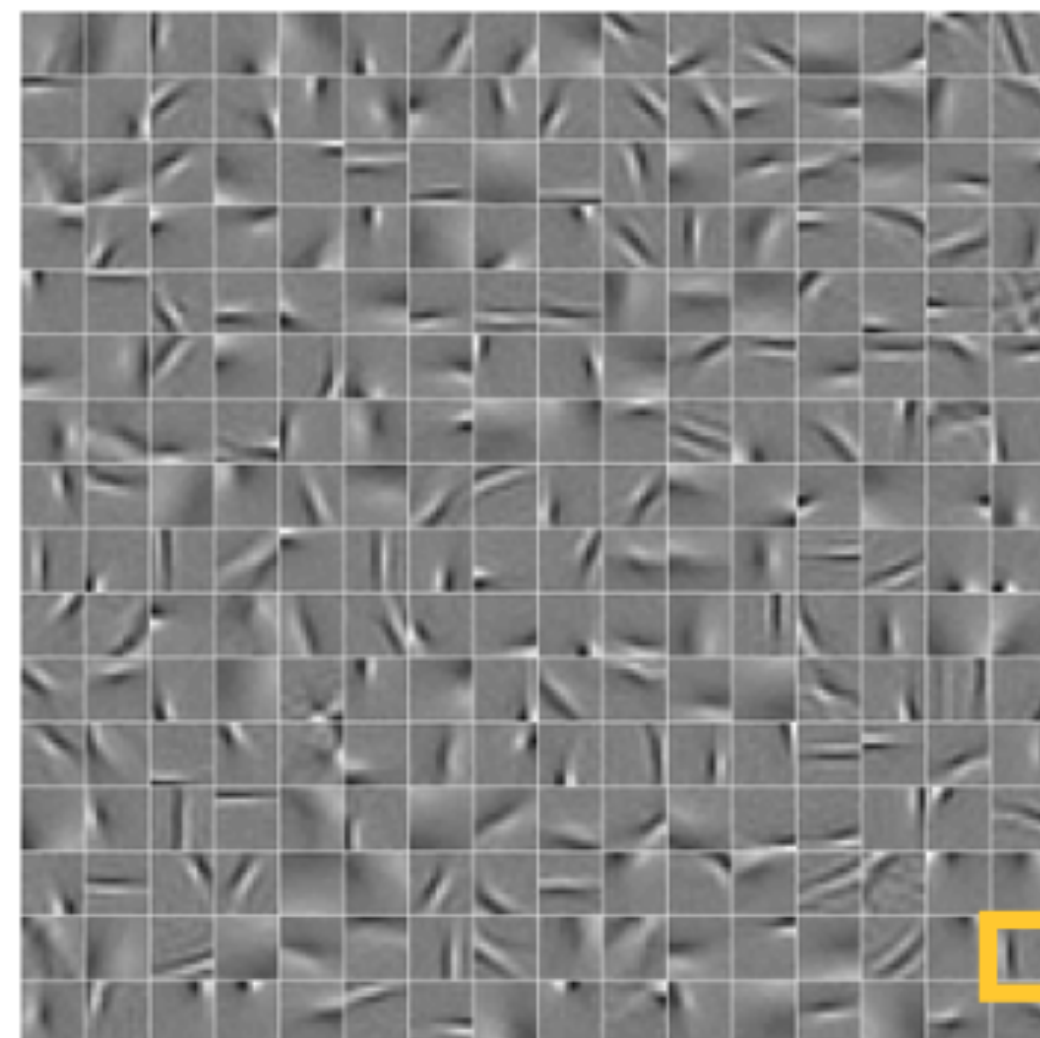Neuron B

[Zhang, Isola, Efros 2016]

[fig from: Torralba, Isola, Freeman 2024]

# Common features

[Hubel and Wiesel 59]



Electrical signal from brain

Recording electrode

Visual area of brain

Stimulus

oriented filter

Filters in AlexNet

[fig from Andrea Vedaldi]

# Rosetta Neurons



StyleGAN2  ResNet50  CLIP-RN  DINO-RN  DINO-ViT  MAE

Example Image

[Dravid*, Gandelsman*, et al. 2023]

Outline:

1. What's a representation?

2. How to measure representational similarity?

3. Which representations are similar and which are different?

4. What drives representational alignment?

5. Making representations more aligned

Outline:

Embedding

Representation learning

Generative modeling

Data

# What does training a deep net classifier look like?

Input data

x

Series of geometric transformations

(i.e., a neural net)

Target output

y       "Cat"

"Dog"

linear

$$\mathbf{x}_{\mathbf{out}} = \mathbf{W}\mathbf{x}_{\mathbf{in}} + \mathbf{b}$$

Embedding

Representation learning

Generative modeling

Data

$\mathbf{x}_{\mathbf{out}}$

$\mathbf{x}_{\mathbf{in}}$

**linear**

$$\mathbf{x_{out}} = \mathbf{W}\mathbf{x_{in}} + \mathbf{b}$$

**relu**

$$x_{\mathtt{out}}[i] = \max(x_{\mathtt{in}}[i], 0)$$

**L2-norm**

$$x_{\mathtt{out}}[i] = \frac{x_{\mathtt{in}}[i]}{\|\mathbf{x_{in}}\|_2}$$

**softmax**

$$x_{\mathtt{out}}[i] = \frac{e^{-\tau x_{\mathtt{in}}[i]}}{\sum_{k=1}^{K} e^{-\tau x_{\mathtt{in}}[k]}}$$

Loss: 0.65

softmax

linear

relu

linear

relu

linear

Embedding

Representation learning

Generative modeling

Data

**Each layer is a representation**

A "representation" is an assignment datapoints to locations in some space

i.e. a labeled point cloud:

# Definitions and notation

- A **representation** is a mapping $f : \mathcal{X} \to \mathcal{Z}$, where $x \in \mathcal{X}$ is data and $z \in \mathcal{Z}$ is some transformation of the data.

  - Typically we have $\mathcal{Z} = \mathbb{R}^d$, i.e. the representation maps data to vector embeddings.

## Summary #1:

All layers are a representation, and so are the input data and the output beliefs.

Representations can be understood in terms of their geometry.

Outline:

# Definitions and notation

- **Representational similarity** is a measure $d : f_1 \times f_2 \to \mathbb{R}$

  - It takes two representations as input and outputs a number that is higher if the two representations are to be considered more alike.

  - Often we will measure $d$ over a finite set of datapoints, $\mathbf{Z}_1 = \{f_1(x^{(i)})\}_{i=1}^{n}$, $\mathbf{Z}_2 = \{f_2(x^{(i)})\}_{i=1}^{n}$, with $d^z : \mathbf{Z}_1 \times \mathbf{Z}_2 \to \mathbb{R}$

# The main question

Neural net 1's embeddings ($\mathbf{Z}_1$)

Neural net 2's embeddings ($\mathbf{Z}_2$)



**How similar are these two point clouds?**

# Regression-based metrics

Neural net 1's embeddings ($\mathbf{Z}_1$) $\Rightarrow h$ Neural net 2's embeddings ($\mathbf{Z}_2$)

$$d(\mathbf{Z}_1, \mathbf{Z}_2) = \frac{1}{n} \sum_{i=1}^{n} ||h^*(z_1^{(i)}) - z_2^{(i)}||$$

$$h^* = \arg\min_{h} \frac{1}{n} \sum_{i=1}^{n} ||h(z_1^{(i)}) - z_2^{(i)}||$$

# Two equivalent representations under linear regression

Neural net 1's embeddings ($\mathbf{Z}_1$)

Neural net 2's embeddings ($\mathbf{Z}_2$)

# Two equivalent representations under linear regression

Neural net 1's embeddings ($\mathbf{Z}_1$)

Neural net 2's embeddings ($\mathbf{Z}_2$)

# Kernel-alignment metrics

$$\mathbf{K}_{\text{vision}}$$



similar

dissimilar

Restrict our attention to **vector embeddings**

$$f : \mathcal{X} \to \mathbb{R}^n$$

Characterize a representation in terms of its **kernel**

$$K : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$$

$$K(x_i, x_j) = \langle\, f(\text{🍎}), f(\text{🍊})\, \rangle$$

# Kernel-alignment metrics

# Two representations with equivalent kernels

Neural net 1's embeddings ($\mathbf{Z}_1$)

Neural net 2's embeddings ($\mathbf{Z}_2$)



Rigid transformations don't change distances

# Centered Kernel Alignment (CKA)

- Kernel alignment metrics are invariant to isometries (i.e. rotation, translation, mirror flips, "glide reflections")

- CKA says: Let's also be invariant to isotropic scaling

kernel similarity

subtracts mean similarity

normalize by scale

$$\text{CKA}(\mathbf{K}_1, \mathbf{K}_2) = \frac{\text{tr}(\mathbf{K}_1 H \mathbf{K}_2 H)}{\sqrt{\text{tr}(\mathbf{K}_1 H \mathbf{K}_1 H)\text{tr}(\mathbf{K}_2 H \mathbf{K}_2 H)}}$$

[Kornblith, Norouzi, Lee, Hinton, ICML 2019]

# Two equivalent representations under CKA

Neural net 1's embeddings ($\mathbf{Z}_1$)

Neural net 2's embeddings ($\mathbf{Z}_2$)

# Nearest-neighbor kernel-alignment metric

What percent of my nearest-neighbors under representation f are also my nearest neighbors under representation g?



$f_1$ $f_2$

[Huh*, Cheung*, Wang*, Isola, ICML 2024]
[Park et al. (2024), Klabunde et al. (2023) Oron et al. (2017)]

# Metrics measure sameness up to a transformation T



Bijection

**Cycle-consistency loss**

Linear, MLP

**Regression-based metrics**

Isometry

**Kernel-alignment metrics**

Identity

**Representation,** $z_0$

$$\{z \mid \exists \theta \text{ s.t. } z = T(z_0, \theta)\}$$

Space of all points that are a transformation of z, up to some class of transformations $T$.

# Which way of measuring is best?

- My opinion: **kernel alignment metrics**

- Why? Because *distance* is the thing that matters for most downstream tasks

  - Two representations that are related by an isometry are the same for most practical purposes

  - Linear isometry —> equivalence in: retrieval, k-NN classifier, min-norm linear regression, MLPs in the NTK regime, …

- (We could make this definitional: *a representation is a specificiation of* $d : \mathscr{X} \times \mathscr{X} \to \mathbb{R}$)

[See more: "Getting Aligned on Representational Alignment," Sucholutsky*, Muttenthaler*, et al. arXiv 2024]

Summary #2:

Representations can be compared via distance functions.

Each distance yields different inferences you can make about how a representation will behave, and what you can do with it.

Outline:

1. What's a representation?

2. How to measure representational similarity?

**3. Which representations are similar and which are different?**

4. What drives representational alignment?

5. Making representations more aligned

# How different are these images?

$$D(\quad\quad,\quad\quad)$$

# Which image is more similar to the middle?



< Clap >

< DINO >

# Which image is more similar to the middle?



< Clap >

DINO

# Which image is more similar to the middle?



< Clap >

< Clap >

# Which image is more similar to the middle?



< Clap >

< Clap >

**Metric-Human Alignment**

Fu*, Tamir*, Sundaram*, Chai, Zhang, Dekel, Isola. *DreamSim.* NeurIPS 2023.

# Investigating representations in the brain

How similar are these two images?



How about these two?



[Kriegeskorte, Mur, Ruff, et al. 2008]

# Investigating a representation via similarity analysis

**Representational Dissimilarity Matrix**



$$||z^{(i)} - z^{(j)}||$$

Neural activation vector

[Kriegeskorte, Mur, Ruff, et al. 2008]

# Investigating a representation via similarity analysis

IT Neuronal Units

Deep net (in paricular, HMO)



[Yamins, Hong, Cadieu, Solomon, Seibert, DiCarlo, PNAS 2014]

# What's the color space in which a language (model) sees?

*Color space*: a mapping from a spectral power distribution to 3 numbers

- Camera CCD: RGB color space
- Human vision: Lab color space

How did we determine this for humans?

- Ask them which colors are similar and which are different
- Find a 3D projection that best preserves distances

# What's the color space in which a language (model) sees?

- Ask an LLM which colors are similar and which are different
- Find a 3D projection that best preserves distances

How similar is red to orange. Output a single number between 0 and 1.

0.8

# What's the color space in which a language (model) sees?

- Ask an LLM which colors are similar and which are different
- Find a 3D projection that best preserves distances



GPT-4.5's perceptual color space

but maybe it lied…

# What's the color space in which a language (model) sees?

- Measure distance between LLM embeddings of different color words
- Find a 3D projection that best preserves distances



CIELAB

BERT, controlled context

[Abdou et al. 2021]

# Brains vs Machines

Deep nets and the human/primate brain both learn similar metric spaces.

Deep nets organize visual information similarly to how our brains do!

[Yamins, Hong, Cadieu, Solomon, Seibert, DiCarlo, PNAS 2014]

# Alignment between different computer vision systems

# Experiment: Is alignment between vision models increasing as vision systems become stronger?

**Hypothesis 1:**

There are many different ways one can represent the visual world, and each can be highly effective.

**Hypothesis 2:**

All strong visual representations are alike.

["Anna Karenina scenario," Bansal et al. 2021]

# Experiment: Is alignment between vision models increasing as vision systems become stronger?

- 78 vision models: different architectures, objectives, training data distributions.

- Group models by performance on VTAB, and measure representational similarity within each group.

# Experiment: Is alignment between vision models increasing as vision systems become stronger?

All strong representations are alike, each weak representation is weak in its own way.



UMAP of model representations

# VTAB tasks solved

- ▲ Random Initialization
- ◆ Classification
- ✚ MAE
- ● Contrastive
- ★ CLIP

# Alignment between different modalities

# Experiment: Is language-vision alignment increasing?

**Hypothesis 1:**
As language models get better and better, they will become more and more specific to language, and start being less generally useful for vision.

**Hypothesis 2:**
Better language models are better vision models.

**Hypothesis 2+:**
The best language model is the best vision model. They converge to the *same* representation.

$$\text{sim}\left(\ \boxed{\begin{array}{c} \mathbf{K}_{\text{vision}} \end{array}}\ ,\ \boxed{\begin{array}{c} \mathbf{K}_{\text{text}} \end{array}}\ \right)$$

$z_1^{\text{vision}}$

$f_{\text{vision}}$

"Sunset over Glacier Point"

$f_{\text{text}}$

$z_1^{\text{text}}$

$z_2^{\text{vision}}$

"Yosemite valley"

$z_2^{\text{text}}$

$z_3^{\text{vision}}$

"San Francisco during the California Gold Rush"

$z_3^{\text{text}}$

$z_4^{\text{vision}}$

"A Boston Red Sox game at Fenway Park"

$z_4^{\text{text}}$

**Wikipedia Image Text Dataset**

[Srinivasan, Raman, Chen, Bendersky, Najork 2021]

# Strong models converge in representation

Summary #3:

Humans and deep nets both measure distances between images in similar ways.

Different vision models, and language models, seem to be converging in how they measure distanes.

Outline:

1. What's a representation?

2. How to measure representational similarity?

3. Which representations are similar and which are different?

**4. What drives representational alignment?**

5. Making representations more aligned

# The Multitask Scaling Hypothesis



### Hypothesis space

task gradient

Solves task 1

Solves task 2

task gradient

---

**The Multitask Scaling Hypothesis**

There are fewer representations that are competent for $N$ tasks than there are for $M < N$ tasks. As we train more general models that solve more tasks at once, we should expect fewer possible solutions.

"Anna Karenina principle"

"Contravariance Principal"



Boaz Barak @boazbaraktcs

Yet another work demonstrating "Anna Karenina" principle of deep learning - successful deep nets seem to learn the same internal representations, up to the "right" notion of symmetry. Supports bold conjecture of arxiv.org/abs/2110.06296 @rahiment @HanieSedghi @osaukh @bneyshabur



Daniel Yamins @dyamins

13/ We thus argue for a "contravariance" principle: the harder the constraint, the smaller the set of mechanisms that can solve the constraint, and thus the more likely any two solving mechanisms (whether biological or artificial) are to be similar in key ways.

# The Multitask Scaling Hypothesis



Hypothesis space

# Corollary: more data —> more convergence



Conwell et al. 2024 found that, of the factors they tested, **data diet** plays the greatest role in determining brain-machine alignment.

Models trained on more data are more aligned with the brain.

[Conwell, Prince, Kay, Alvarez, Konkle, Nature Communications 2024]

# The Capacity Hypothesis

**The Capacity Hypothesis**

Bigger models are more likely to converge to a shared representation than smaller models.

**The Simplicity Bias Hypothesis**

Deep networks are biased toward finding simple fits to the data, and the bigger the model, the stronger the bias. Therefore, as models get bigger, we should expect convergence to a smaller solution space.

[Gunasekar et al. 2018, Arora et al. 2019, Valle-Perez et al. 2019., Huh et al. 2023, etc]

$$f^* = \underset{f \in \mathcal{F}}{\arg\min} \; \mathbb{E}_{x \sim \text{dataset}} [\mathcal{L}(f, x)] + \mathcal{R}(f)$$

trained model

training objective

function class

regularization

## Task/data pressures

## Regularization

## Model size

Hypothesis space

task gradient

Solves task 1

Solves task 2

task gradient

Hypothesis space

simplicity bias

Functions that solve the tasks

Simple functions

simplicity bias

Hypothesis space 2

optimum

Hypothesis space 1

Scale up architectures

"Contravariance Principal" [Cao & Yamins 2024]

LVX VENIT IN MVNDVN ET DILEXERVNT HOMINES MAGIS TENEBRAS QVAM LVCEM. Io.3.19

ANTRVM PLATONICVM

Maxima pars hominum cecis immersa tenebris
Voluitur assiduè, et s tulto letatur inani:
Adspice ut obiec tis obtutis hereat umbris,
Vt VERI simulacra omnes mirentur amentq̃,

Et s tolidi vanâ ludantur imagine rerum.
Quàm pauci meliore luto, qûi in lumine puro
Secreti à s tolida turba, ludibria cernunt
Rerum umbras rectaq̃, expendunt omnia lance:

Hi posteà erroris nebulâ dignoscere possunt
Vera bona, atque alios cecâ sub nocte latentes
Extrahere in claram lucem conantur, at illis
Nullus amor lucis, tanta es t rationis egestas.

C.C. Harlemensis In.
Janredam Sculpsit.
Henr. Hondius excudit.
1604.

H L. SPIEGEL FIGVRARI ET SCVLPI CVRAVIT. AC DOCTISS. ORNATISSq̃ D PEI. PAAW IN LVGDVN. ACAD. PROFESSORI MEDICO D.D.

that drives it: different models are all trying to arrive at a *representation of reality*, meaning a representation of the joint distribution over events in the world that generate the

IT. Correspondence to:  <FIXME>.

ons that solve
the tasks

← ← ← simplicity bias

- All data is mediated via observations

- In this world, we will model **cooccurrences**

$$P_{\mathsf{coor}}(x_a, x_b) \quad \propto \quad \sum_{(t,t'):\, |t-t'| \leq T_{\mathsf{window}}} \mathbb{P}(X_t = x_a, X_{t'} = x_b)$$

- Positives: two observations that cooccur; Negatives: two samples from marginals



$\Big\{$  ,  $\Big\}$

I parked the **car** in a nearby **street.**

- Contrastive learner, with NCE objective converges to PMI :

$$\langle f_X(x_a), f_X(x_b) \rangle \approx \log \frac{P_{\text{coor}}(x_a, x_b)}{P_{\text{coor}}(x_a)P_{\text{coor}}(x_b)}$$

- An embedding in which similarly = (normalized) cooccurrence rate .

"orange"

"apple"

"elephant"

- For bijective, discrete **obs** functions, PMI over **obs** equals PMI over events, which implies that different **obs** converge to same kernel.

Summary #4:

Scaling up task/data/model can drive convergence.

Certain contrastive learners converge to kernel = rate at which events co-occur in nature.

Outline:

1. What's a representation?

2. How to measure representational similarity?

3. Which representations are similar and which are different?

4. What drives representational alignment?

**5. Making representations more aligned**

# Benefits of alignment

- Can share data/supervision between modalities
- A common representation can serve as a bridge/translation
- Can scaffold new models onto existing representations

# Detriments of alignment

- Lack of diversity in the population of models.

- Sometimes one modality has access to qualitatively different information than another, and this information can be useful; alignment will remove this information.

- There might not be a single best representation for all problems. (And in theory there isn't.)

# Aligning and translating representations



**Aligning:**
Train a representation that is aligned with a supervised target representation, up to T

**Translating:**
Find the transformation T that relates two representations

**Identity**

# Training to align representations up to identity transformation

$$\text{"green triangle"} \quad \text{"red square"} \quad \text{"blue circle"}$$

$f_t \qquad f_t \qquad f_t$

$\mathbf{z}_t^{(1)} \qquad \mathbf{z}_t^{(2)} \qquad \mathbf{z}_t^{(3)}$

| | $\mathbf{z}_t^{(1)}$ | $\mathbf{z}_t^{(2)}$ | $\mathbf{z}_t^{(3)}$ |
|---|---|---|---|
| $f_\ell$, $\mathbf{z}_\ell^{(1)}$ | $\mathbf{z}_\ell^{(1)} \cdot \mathbf{z}_t^{(1)}$ | $\mathbf{z}_\ell^{(1)} \cdot \mathbf{z}_t^{(2)}$ | $\mathbf{z}_\ell^{(1)} \cdot \mathbf{z}_t^{(3)}$ |
| $f_\ell$, $\mathbf{z}_\ell^{(2)}$ | $\mathbf{z}_\ell^{(2)} \cdot \mathbf{z}_t^{(1)}$ | $\mathbf{z}_\ell^{(2)} \cdot \mathbf{z}_t^{(2)}$ | $\mathbf{z}_\ell^{(2)} \cdot \mathbf{z}_t^{(3)}$ |
| $f_\ell$, $\mathbf{z}_\ell^{(3)}$ | $\mathbf{z}_\ell^{(3)} \cdot \mathbf{z}_t^{(1)}$ | $\mathbf{z}_\ell^{(3)} \cdot \mathbf{z}_t^{(2)}$ | $\mathbf{z}_\ell^{(3)} \cdot \mathbf{z}_t^{(3)}$ |

## Contrastive Language-Image Pre-training (CLIP)

- Tries to find a representation in which an image and its caption are assigned identical embeddings.

[Radford*, Kim* et al., ICML 2021]

**MLP**

# Training to align representation up to MLP transformation

## Representation Alignment for Generation (REPA)

[Yu, Kwak, Jang, Jeong, Huang, Shin*, Xie*, ICLR 2025]

that drives it: different models are all trying to arrive at a *representation of reality*, meaning a representation of the joint distribution over events in the world that generate the data we observe. Figure 1 conveys this hypothesis: there exists a real world out there (labeled $Z$), which we measure

With sufficiently non-degenerate data, relates $Z$ to $Z$. Can translate between two representations related by an isometry with zero paired examples.

Unique rigid transformation

Learn

Isometry

Translating between representations

Learn

World

airplane

cat

dog

$f_l$

$f_v$

Factorized Hahn-Grant

cat
airplane
dog

[Schnaus, Araslanov, Cremers, arXiv 2025]

see also: Sorscher, Ganguli, Sompolinsky, PNAS 2022;
Lazaridou, Bruni, Baroni, ACL 2014

Bijection

CycleGAN        [Zhu*, Park*, Isola, Efros, ICCV 2017]

$\mathbf{Z}_1$                    $\mathbf{Z}_2$

,

[Zhu*, Park* et al. 2017], [Yi et al. 2017], [Kim et al. 2017]

Bijection

CycleGAN [Zhu*, Park*, Isola, Efros, ICCV 2017]

$\mathbf{Z}_1$

$\mathbf{Z}_2$

$\mathbf{Z}_1$

$\mathbf{Z}_2$

[Zhu*, Park* et al. 2017], [Yi et al. 2017], [Kim et al. 2017]

Bijection

# CycleGAN   [Zhu*, Park*, Isola, Efros, ICCV 2017]

$\mathcal{Z}_1$

$\mathcal{Z}_2$

reconstruction error

Bijection   Isometry

# vec2vec   [Jha, Zhang, Shmatikov, Morris, arXiv 2025]

Embeddings [Original]          Latent Representations [vec2vec]

Method: GAN + cycle consistency loss + kernel matching loss

See also: [Conneau, Lample, Ranzato, Denoyer, Jégou, ICLR 2018]

Summary #5:

Many important problems involve aligning or translating between representations.

You don't necessarily need paired data to do so.

# Cognitive Science

# Neuroscience

# Machine Learning
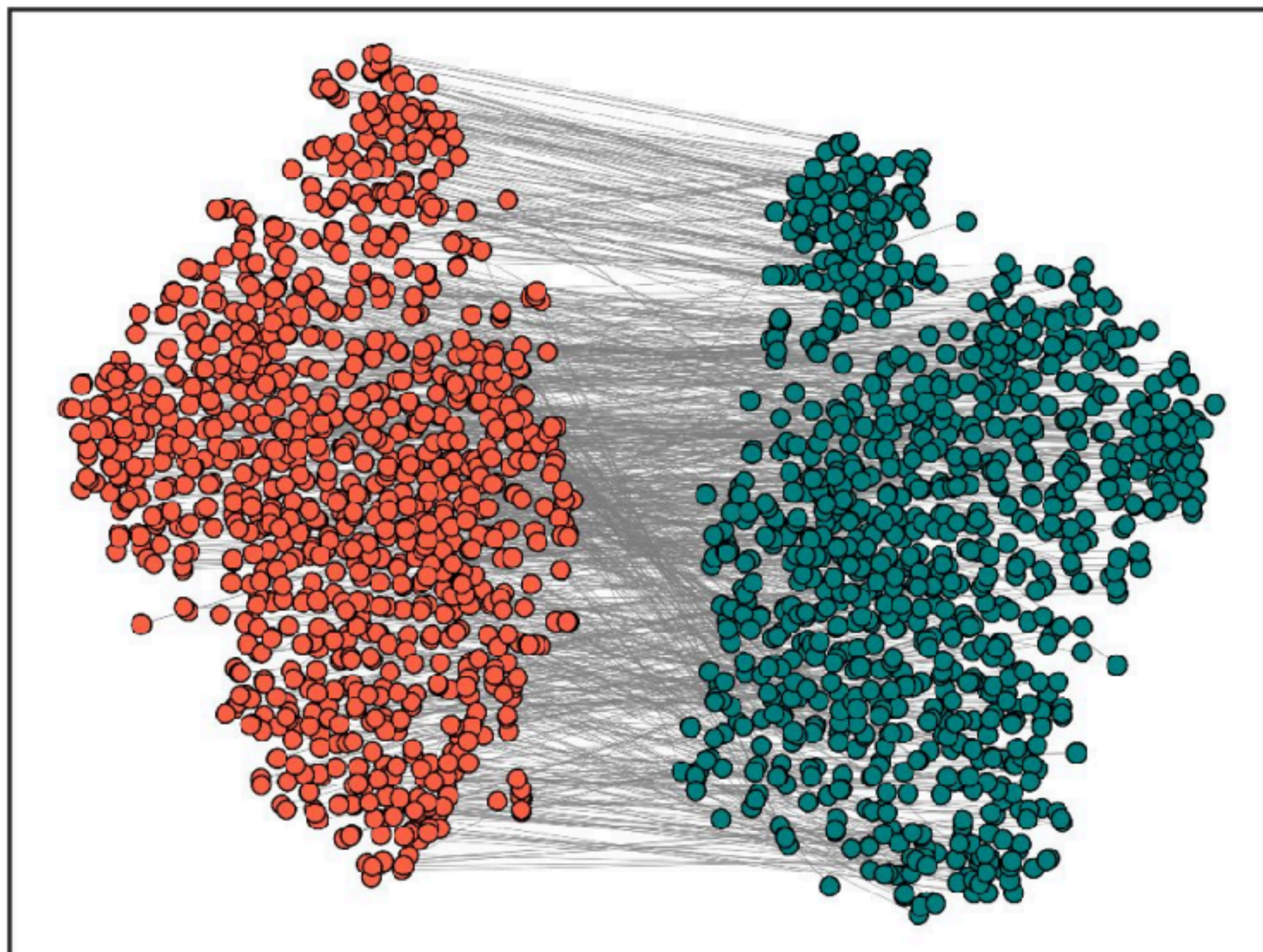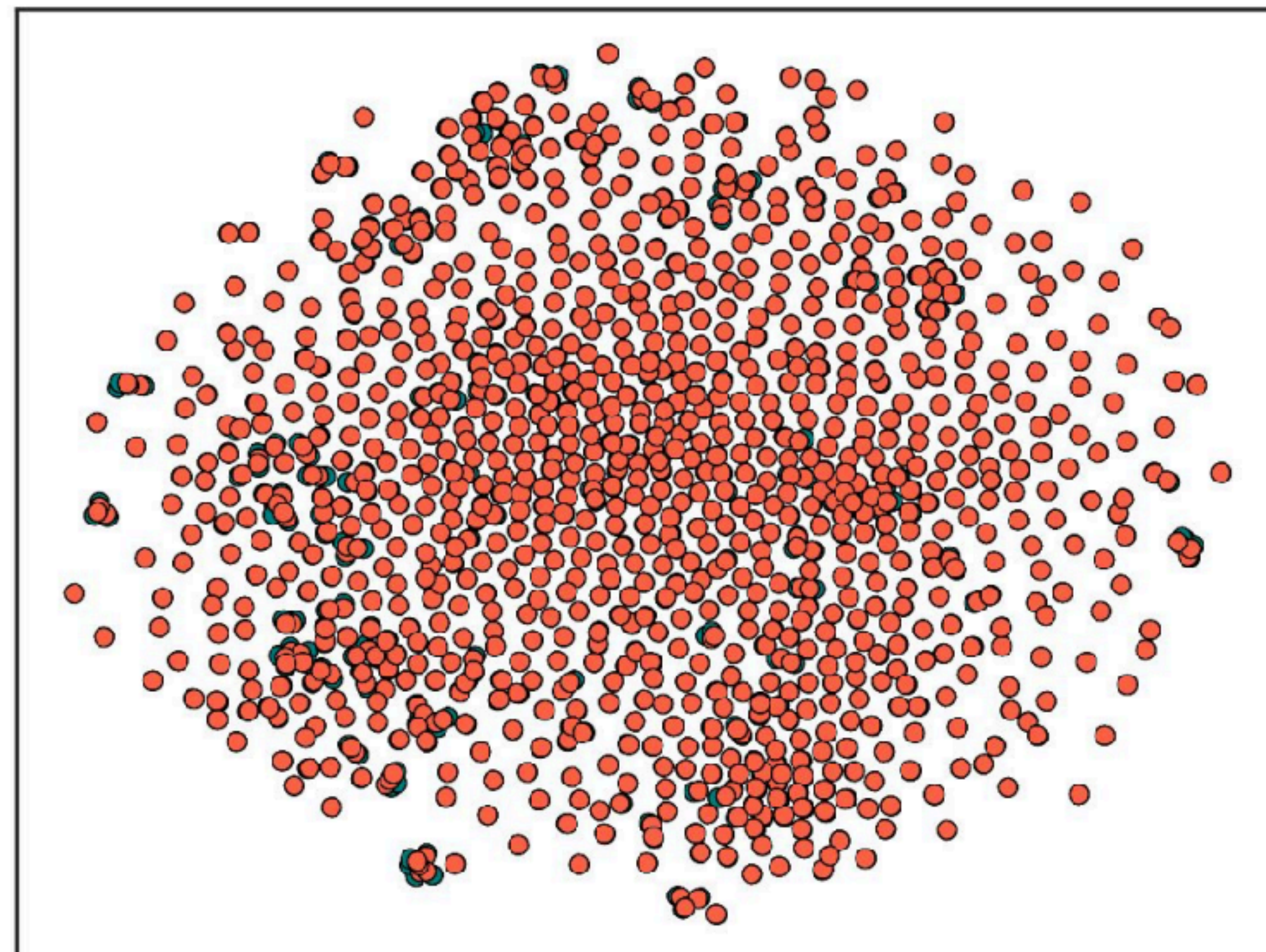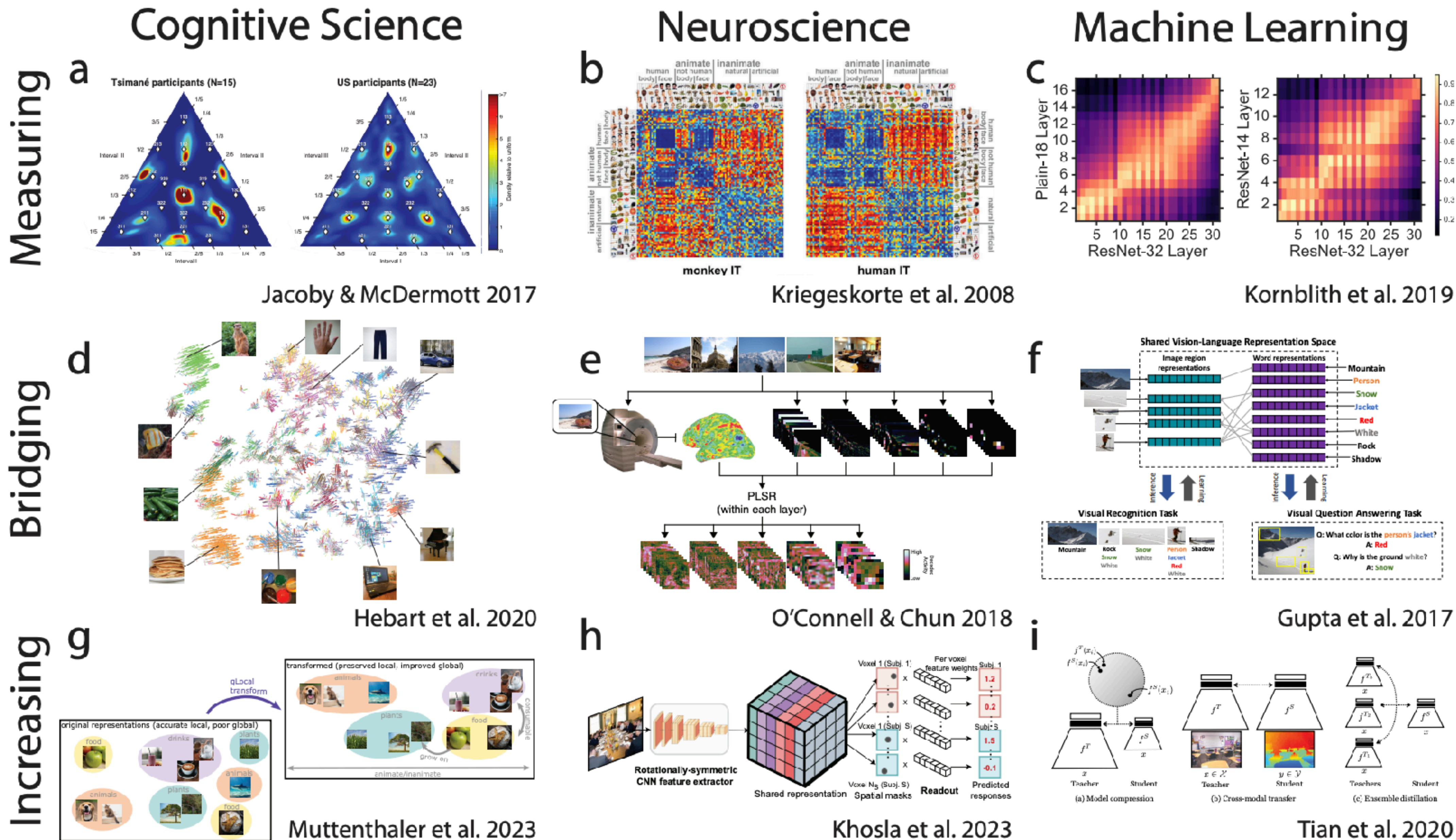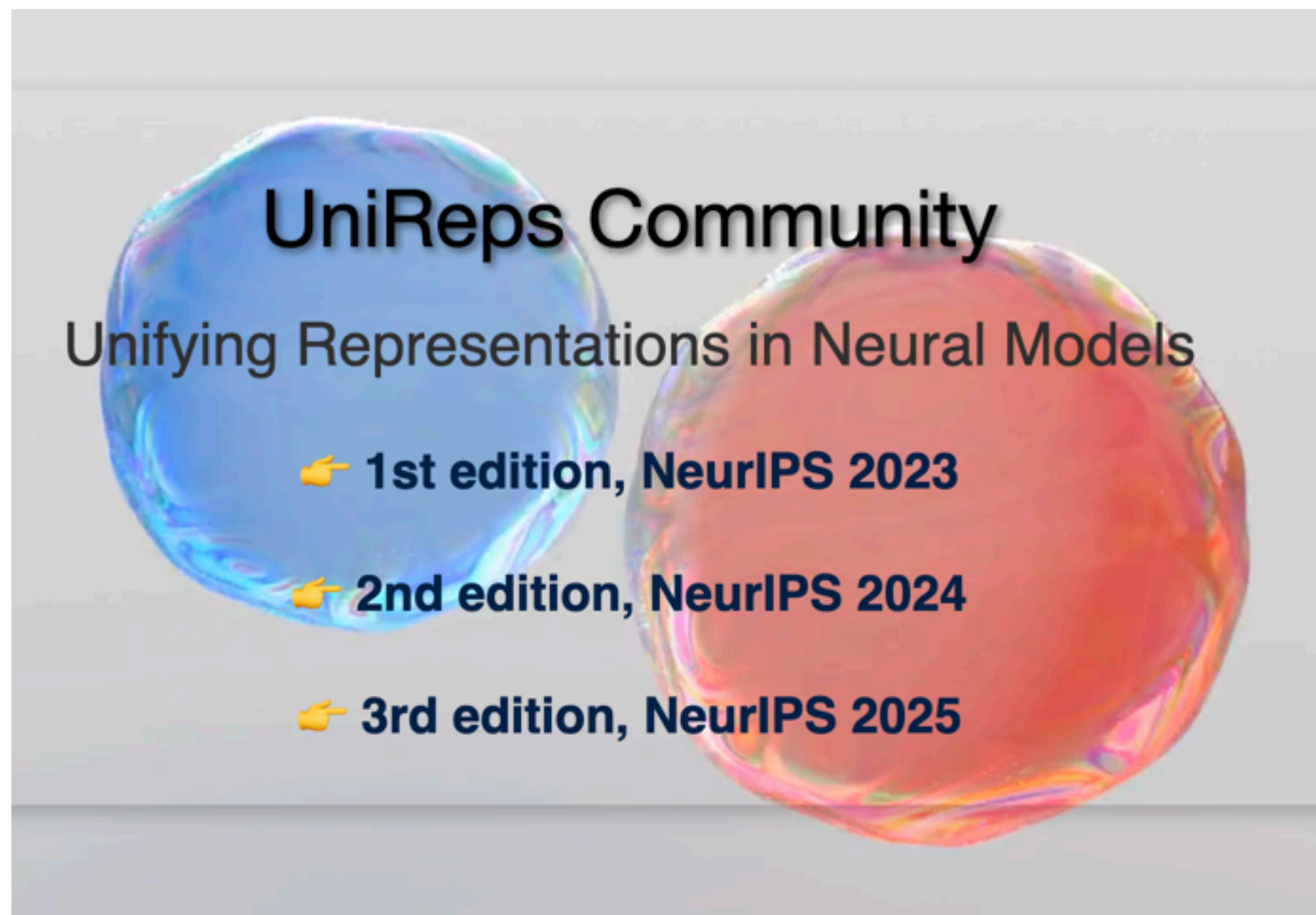
## Measuring



a

Tsimané participants (N=15)    US participants (N=23)

Jacoby & McDermott 2017

b

monkey IT    human IT

Kriegeskorte et al. 2008

c

Kornblith et al. 2019

## Bridging

d

Hebart et al. 2020

e

PLSR
(within each layer)

O'Connell & Chun 2018

f

Shared Vision-Language Representation Space

Visual Recognition Task    Visual Question Answering Task

Gupta et al. 2017

## Increasing

g

Muttenthaler et al. 2023

h

Rotationally-symmetric
CNN feature extractor

Shared representation    Spatial masks    Readout    Predicted responses

Khosla et al. 2023

i

(a) Model compression    (b) Cross-modal transfer    (c) Ensemble distillation

Tian et al. 2020

[See more: "Getting Aligned on Representational Alignment," Sucholutsky*, Muttenthaler*, et al. arXiv 2024]

## UniReps Community

**Unifying Representations in Neural Models**

👉 1st edition, NeurIPS 2023

👉 2nd edition, NeurIPS 2024

👉 3rd edition, NeurIPS 2025

## ICLR 2025 Workshop on Representational Alignment (Re-Align)

## Community Event

*Wednesday, August 13, 10:00 am – 12:00 pm, Room C1.03*

**Universality and Idiosyncrasy of Perceptual Representations**

# Thanks!