

Cognitive computational neuroscience of vision

Nikolaus Kriegeskorte

Department of Psychology, Department of Neuroscience

Zuckerman Mind Brain Behavior Institute

Affiliated member, Electrical Engineering, Columbia University

Kriegeskorte & Douglas 2018

drawing by
Matteo Farinella

COGNITION

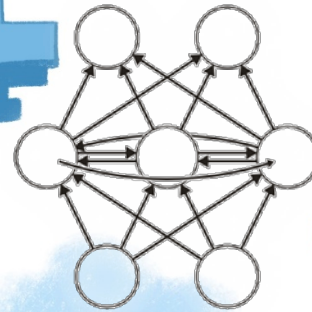
BRAIN

Level of
Description

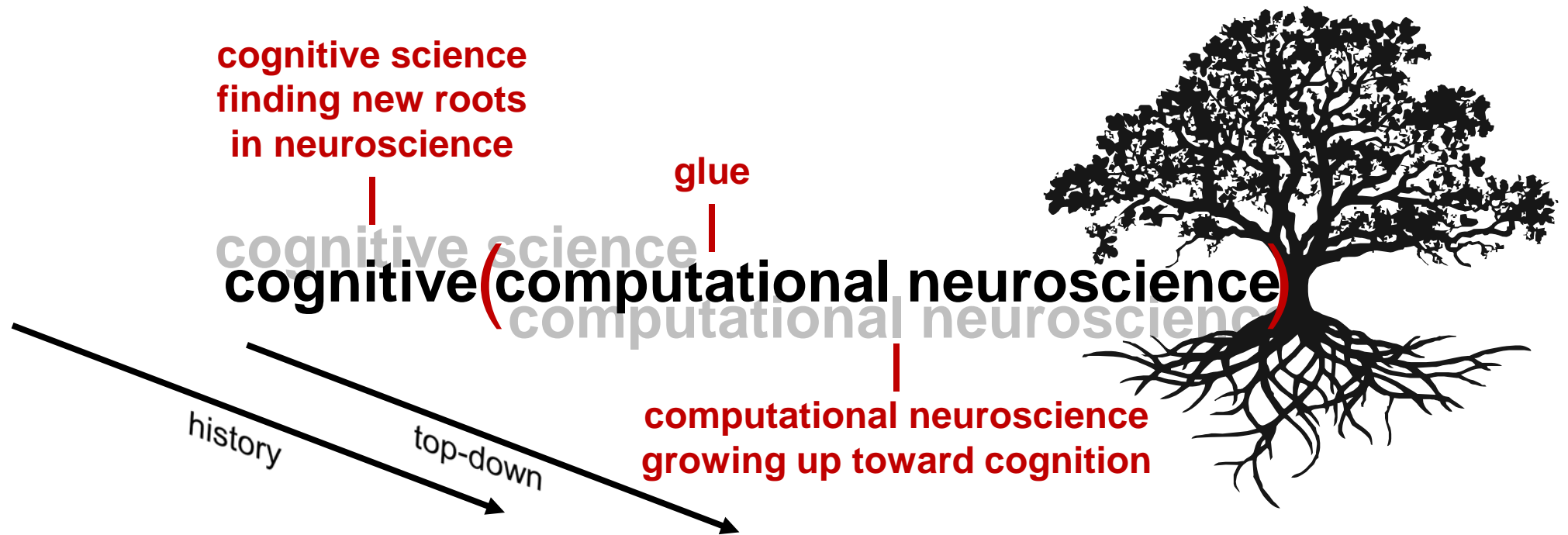
Cognitive
Science

Computational
Neuroscience

Spatial &
Temporal Scale

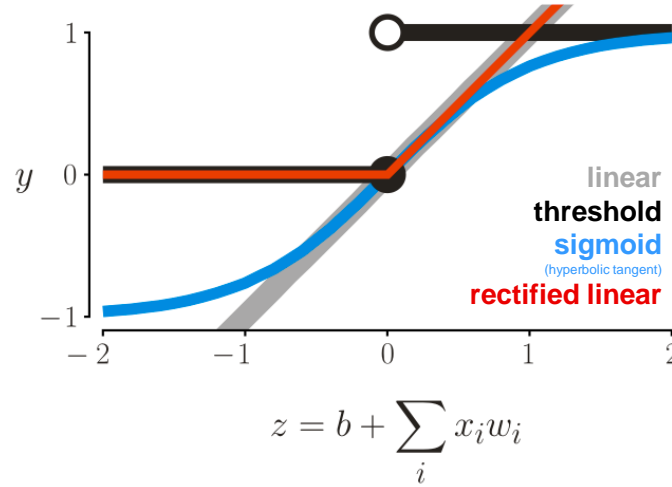
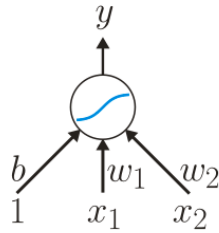


Neural network models
a *language* for expressing theories
of how cognition might be implemented in brains

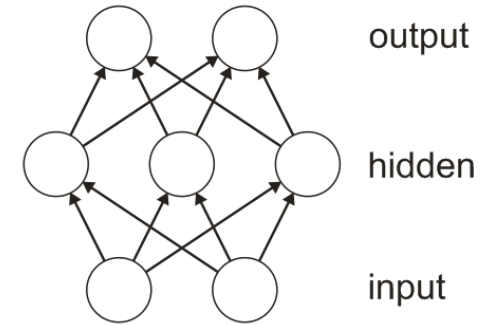


Neural network models

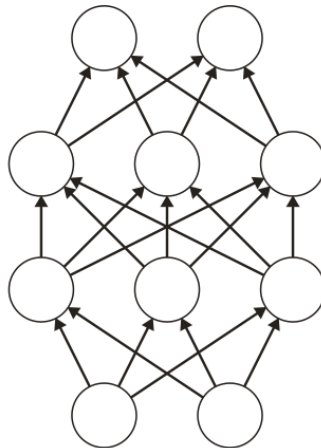
Neural network models



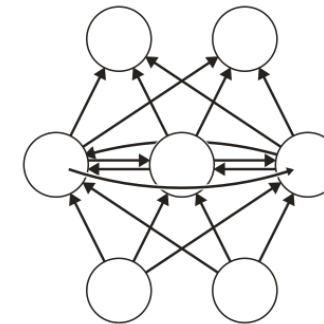
shallow feedforward
(1 hidden layer)



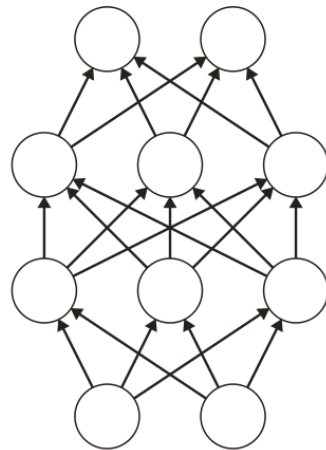
deep feedforward
(>1 hidden layer)



recurrent

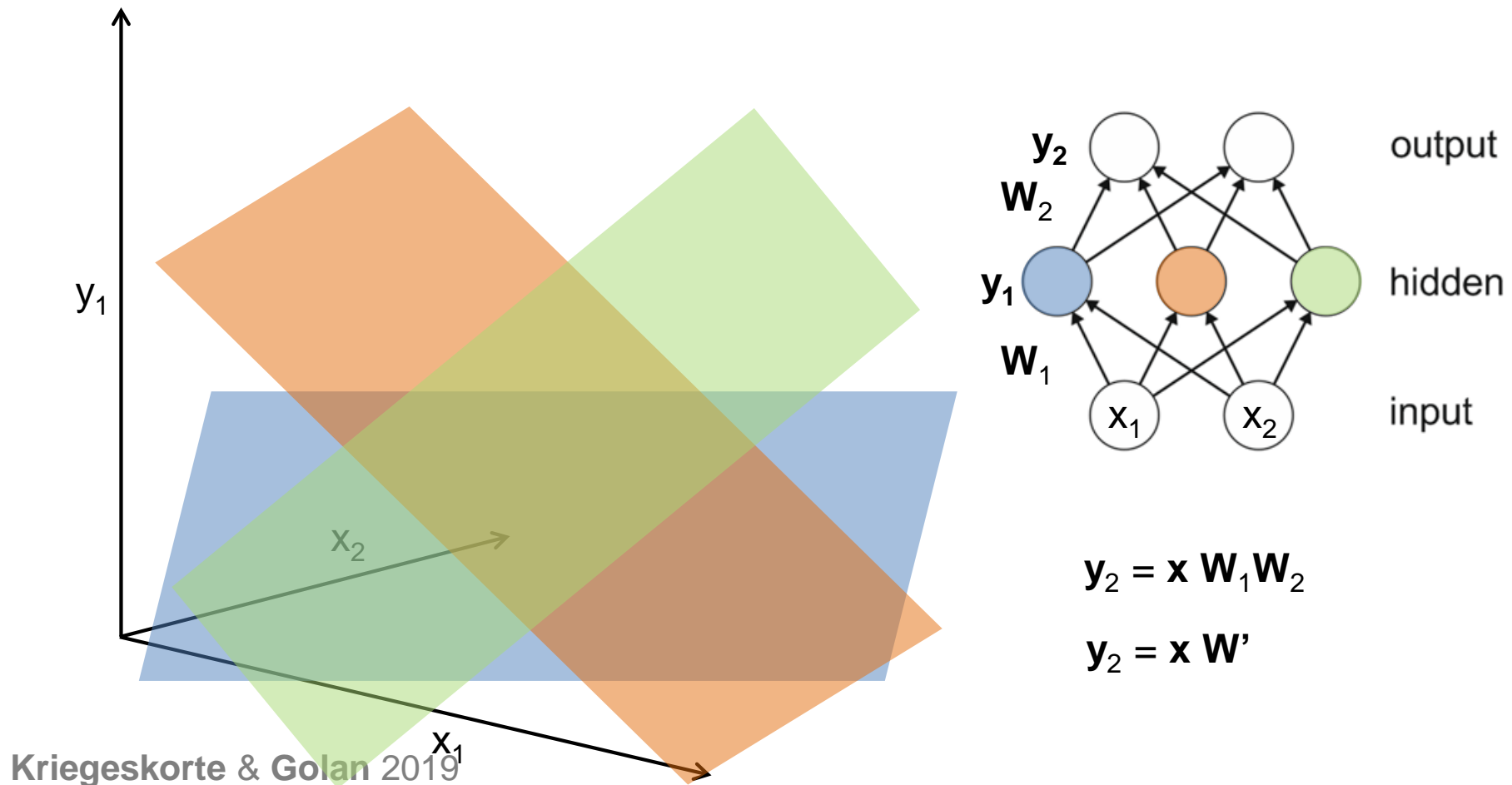


Neural network models



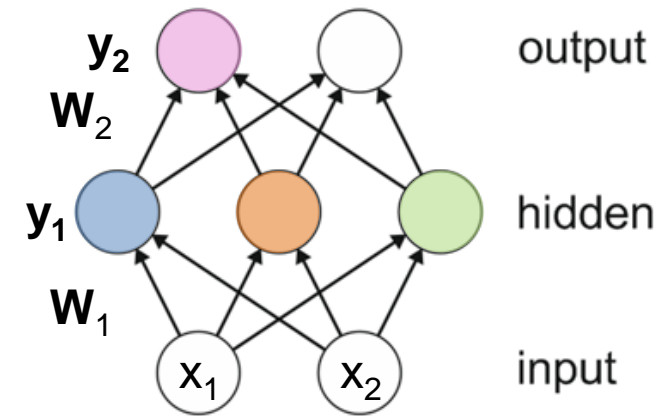
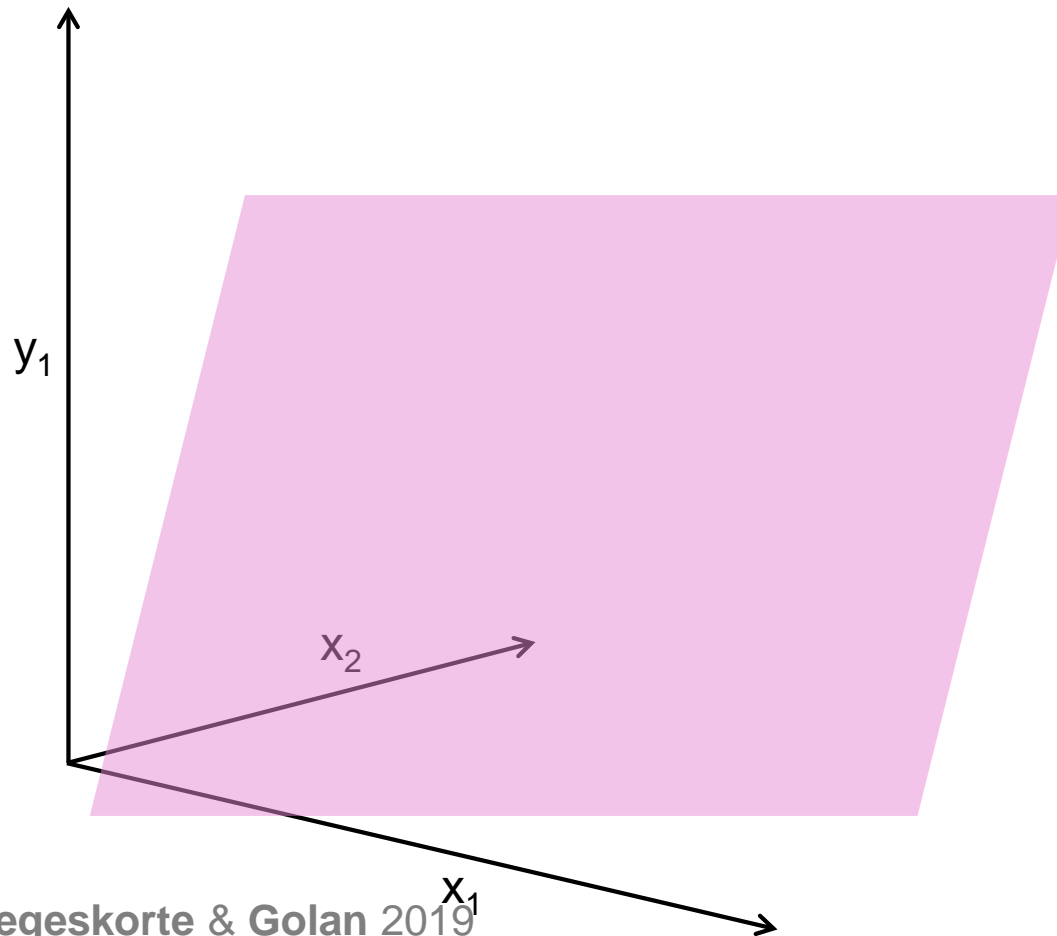
Nonlinear activation function needed to make a hidden layer useful

linear activation functions



Nonlinear activation function needed to make a hidden layer useful

linear activation functions

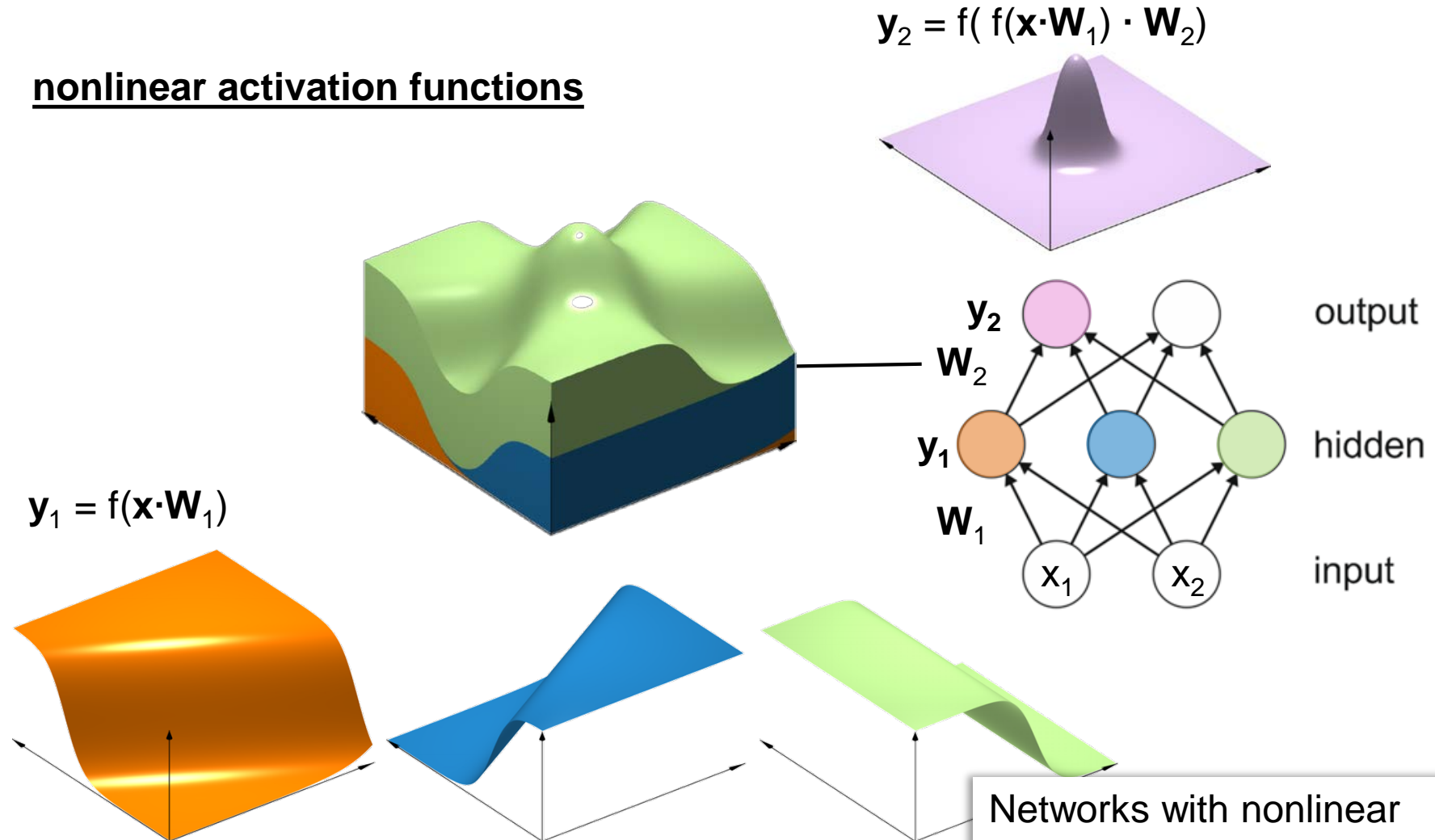


$$y_2 = x W_1 W_2$$

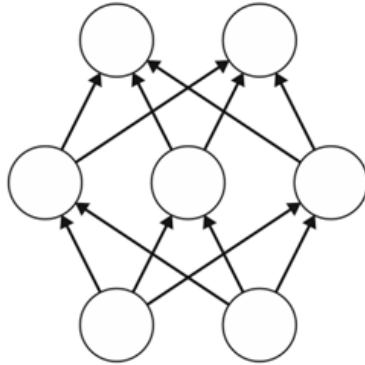
$$y_2 = x W'$$

Nonlinear activation function needed to make a hidden layer useful

nonlinear activation functions



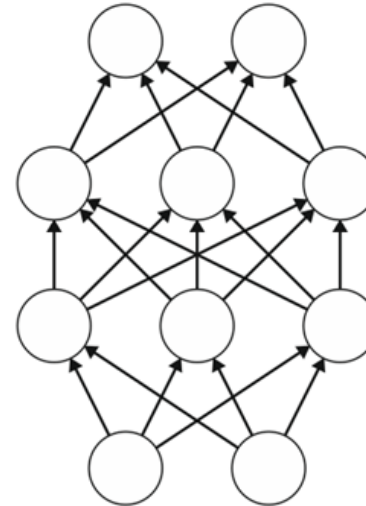
Why deep?



shallow

1 hidden layer

Networks with nonlinear hidden units are *universal function approximators*.



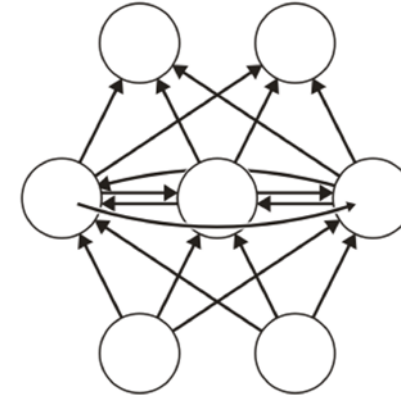
deep

>1 hidden layer

Deep nets can

- *reuse features* downstream
- represent many complex functions more concisely (fewer units and weights).

Why recurrent?



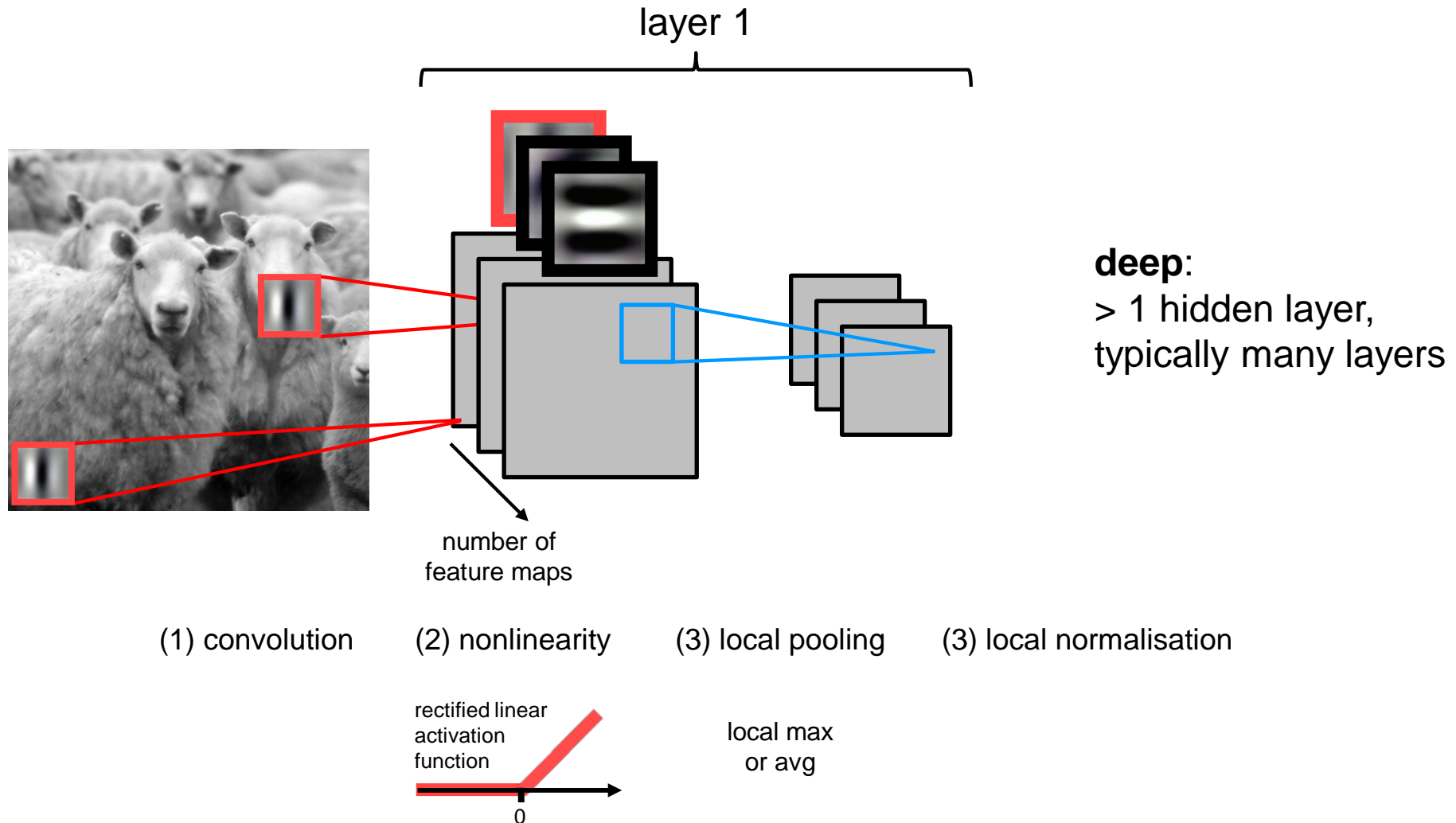
recurrent

Recurrent networks

- can *recycle weights and units* over time
- are *universal approximators of dynamical systems*.

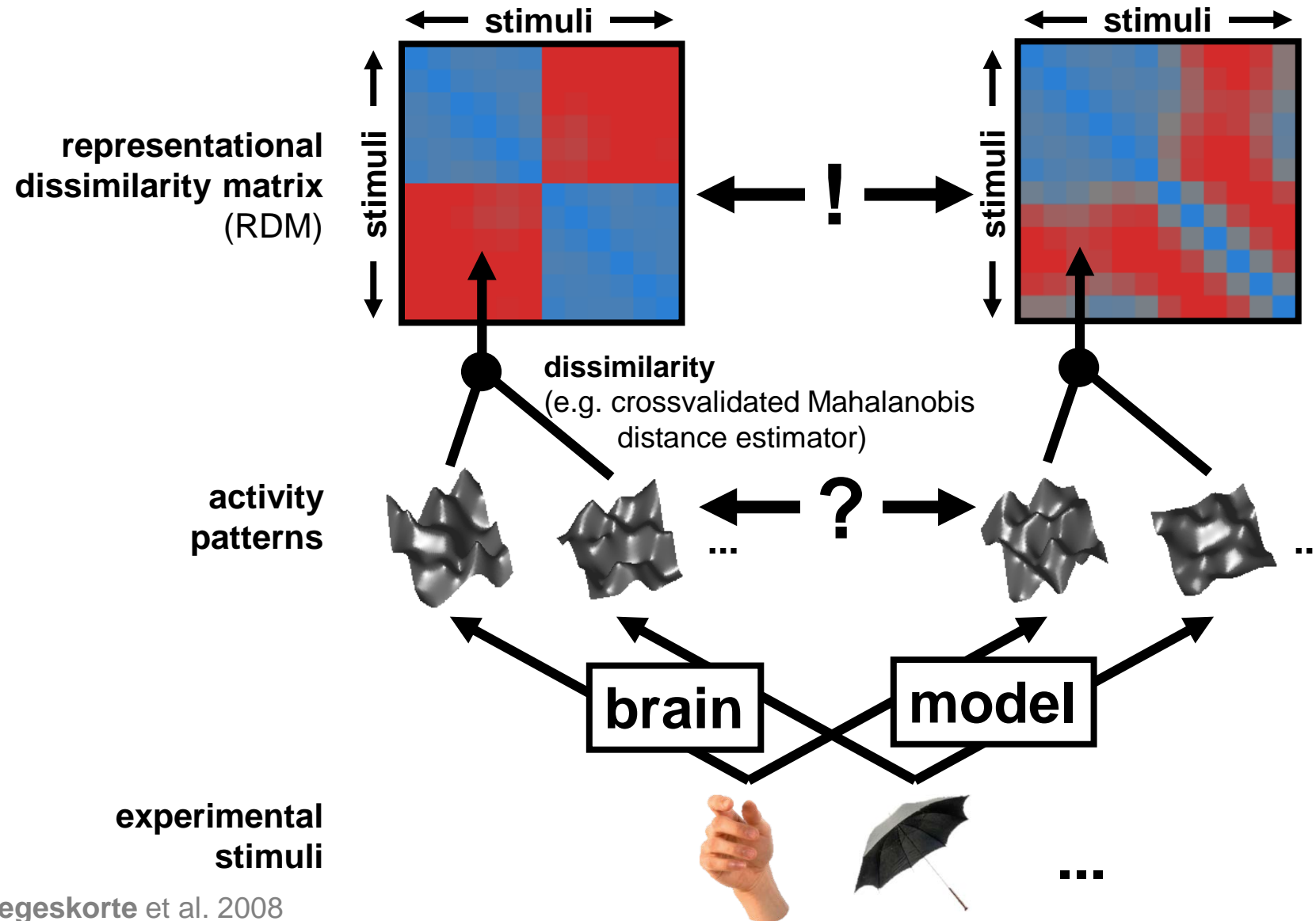
hidden units are *universal function approximators*.

Deep *convolutional* feedforward neural networks

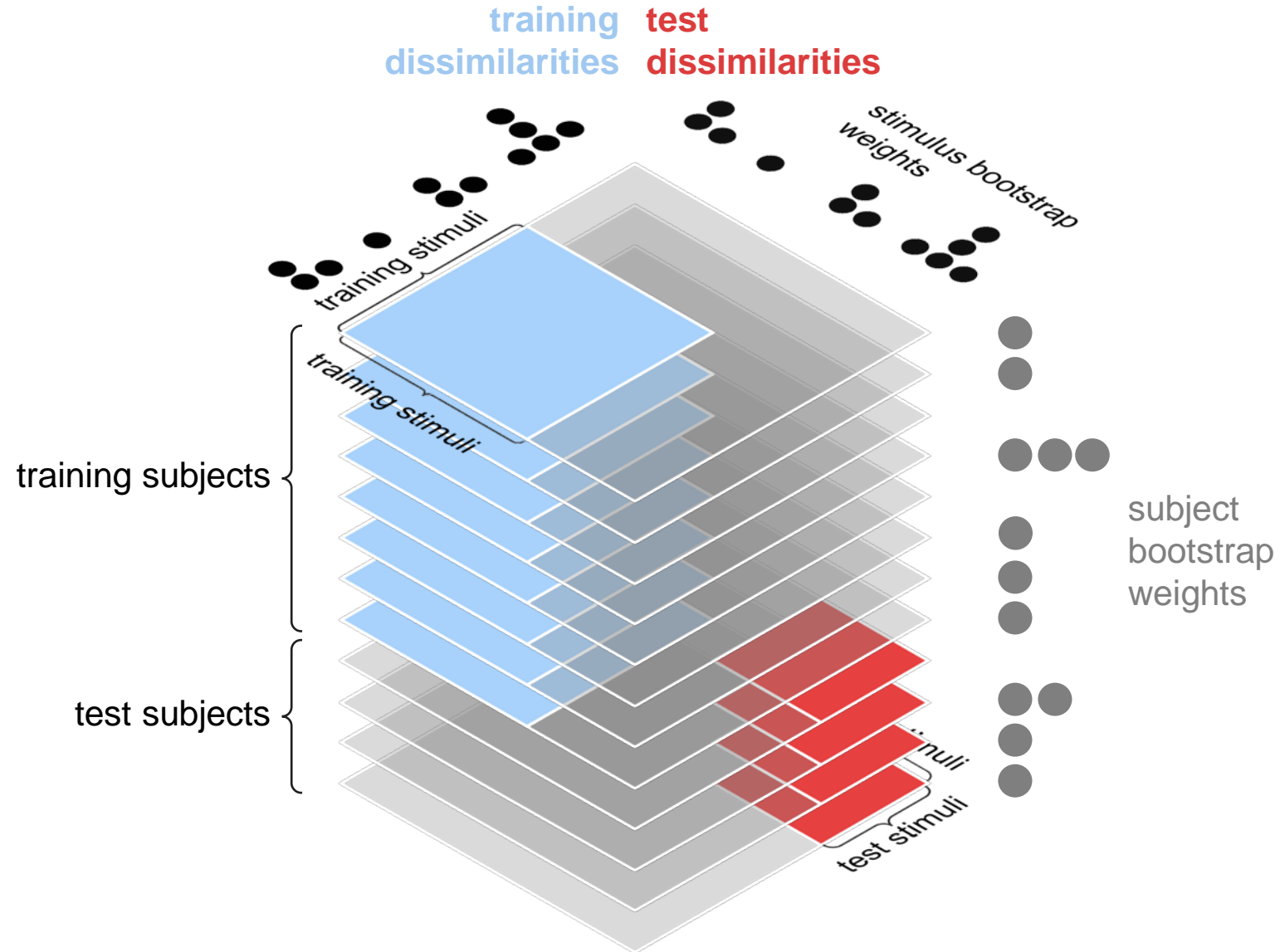


Testing neural network models with brain-activity data

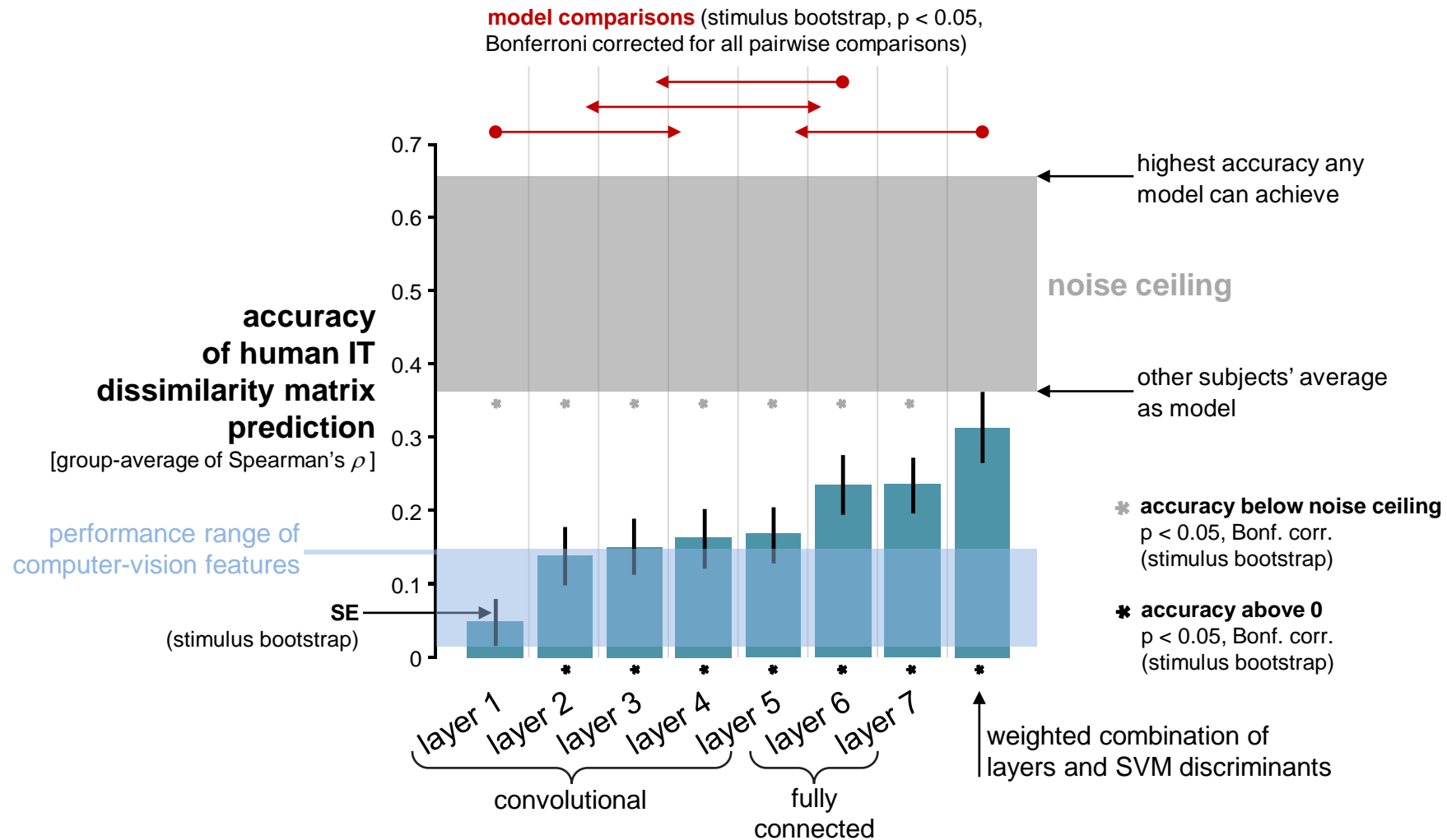
Representational similarity analysis



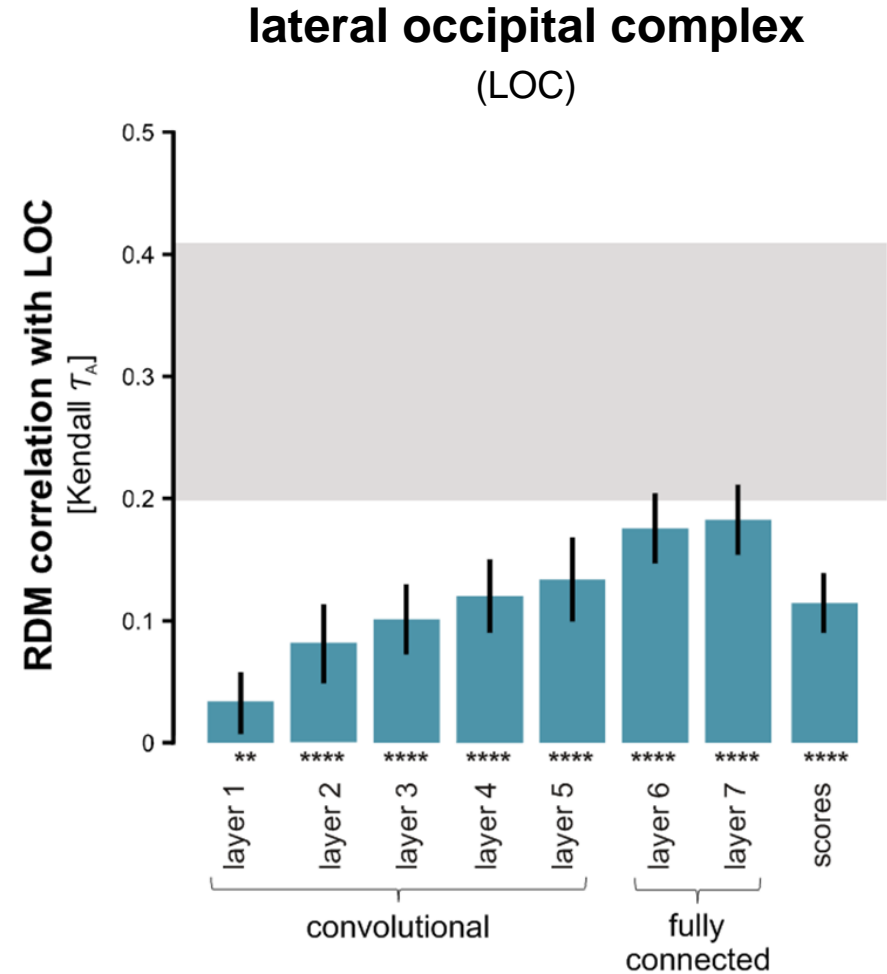
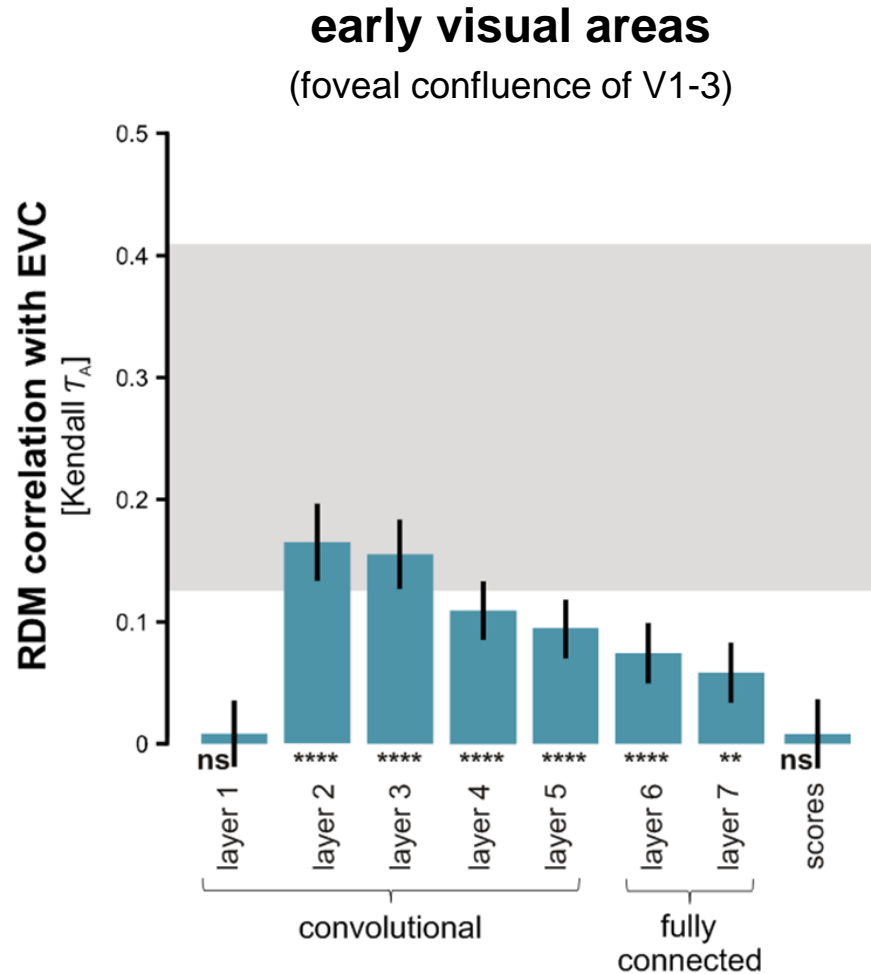
Training and testing in crossvalidation



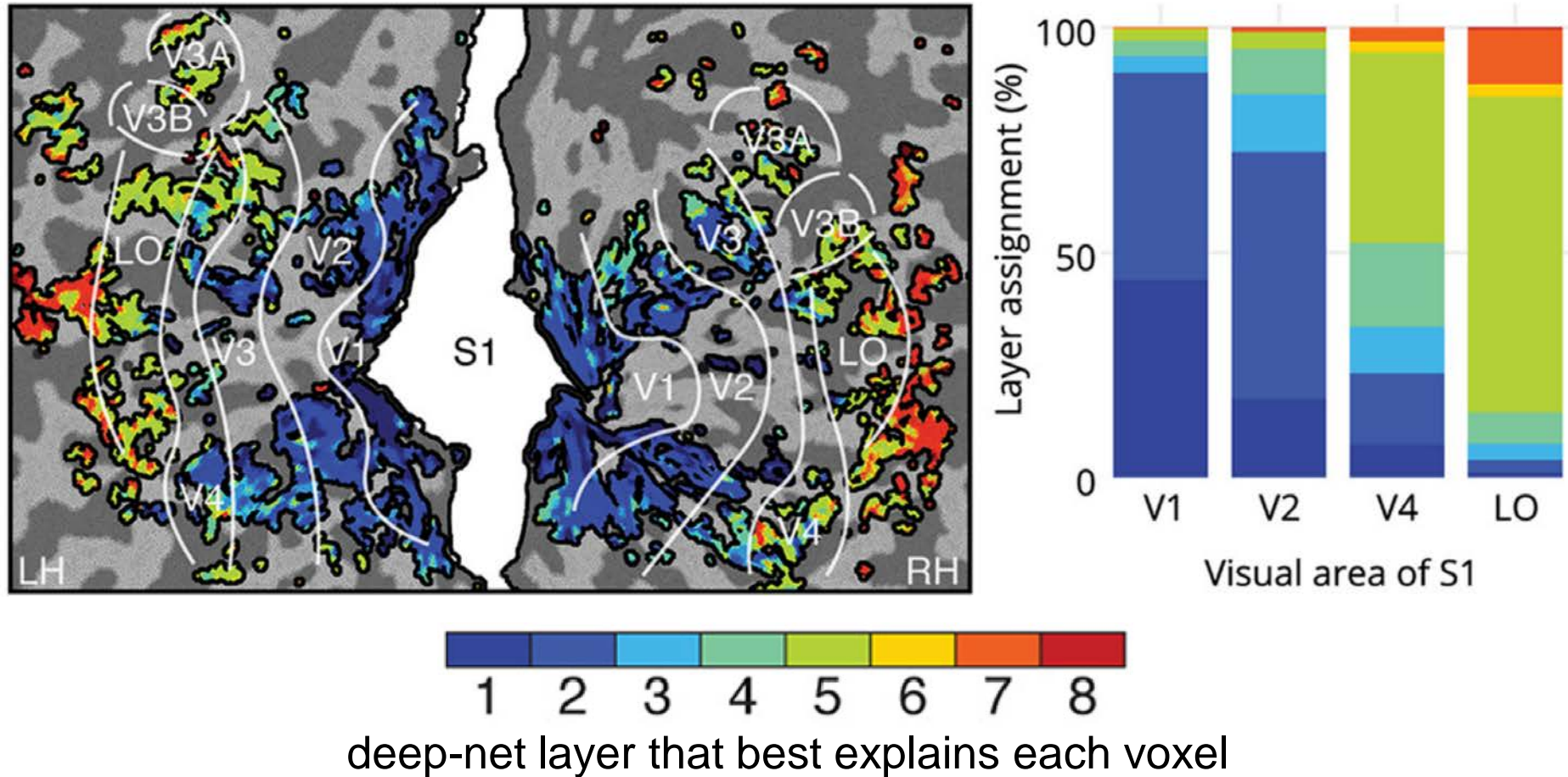
Deep convolutional feedforward networks predict IT representational geometry



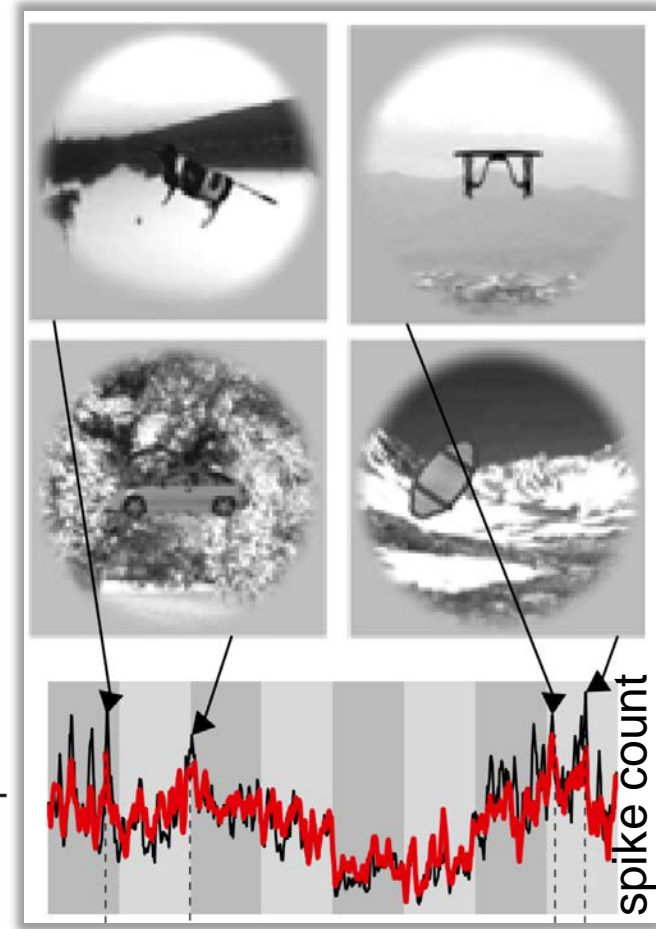
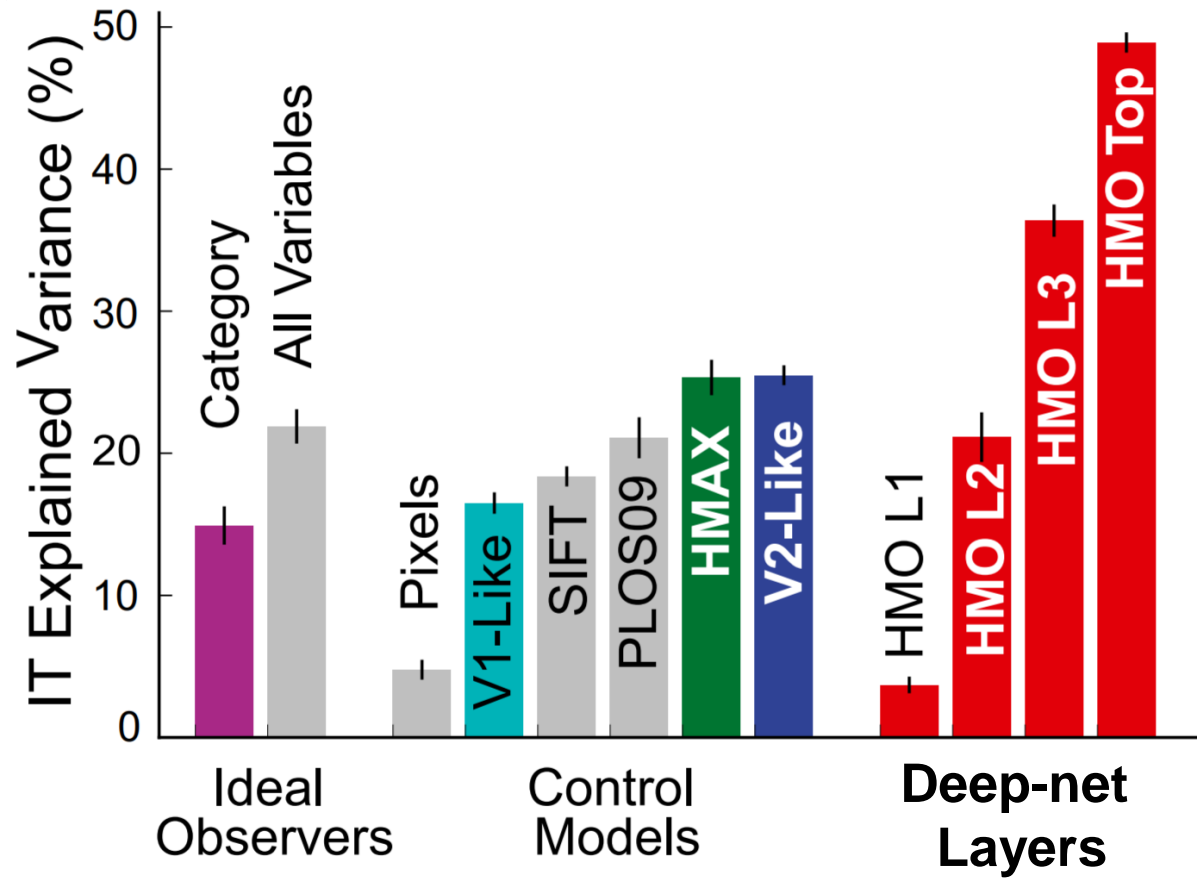
Deep-net layers correspond to stages of the ventral visual stream



Deep-net layers correspond to stages of the ventral visual stream

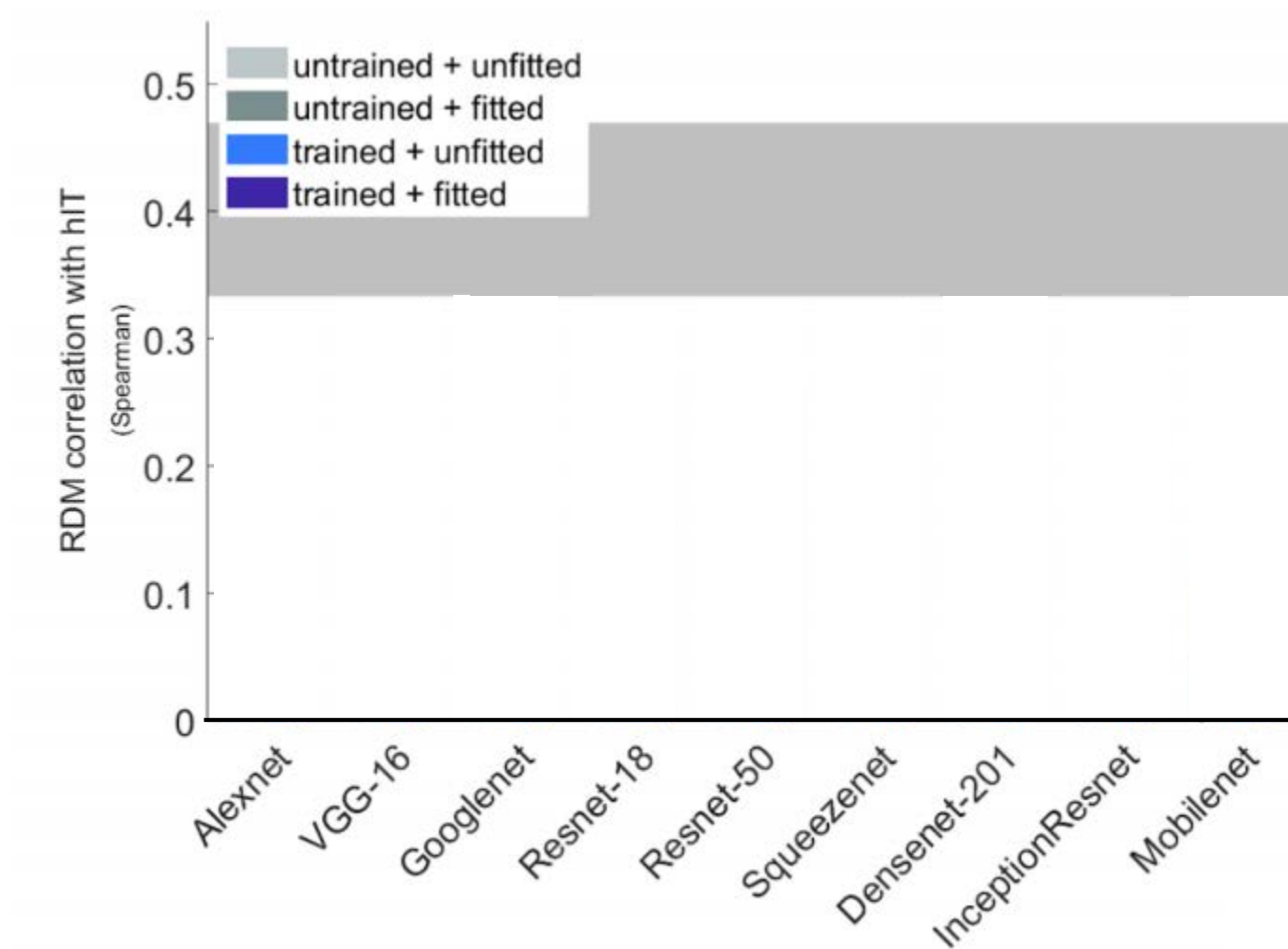


High explained variance for IT neuronal recordings

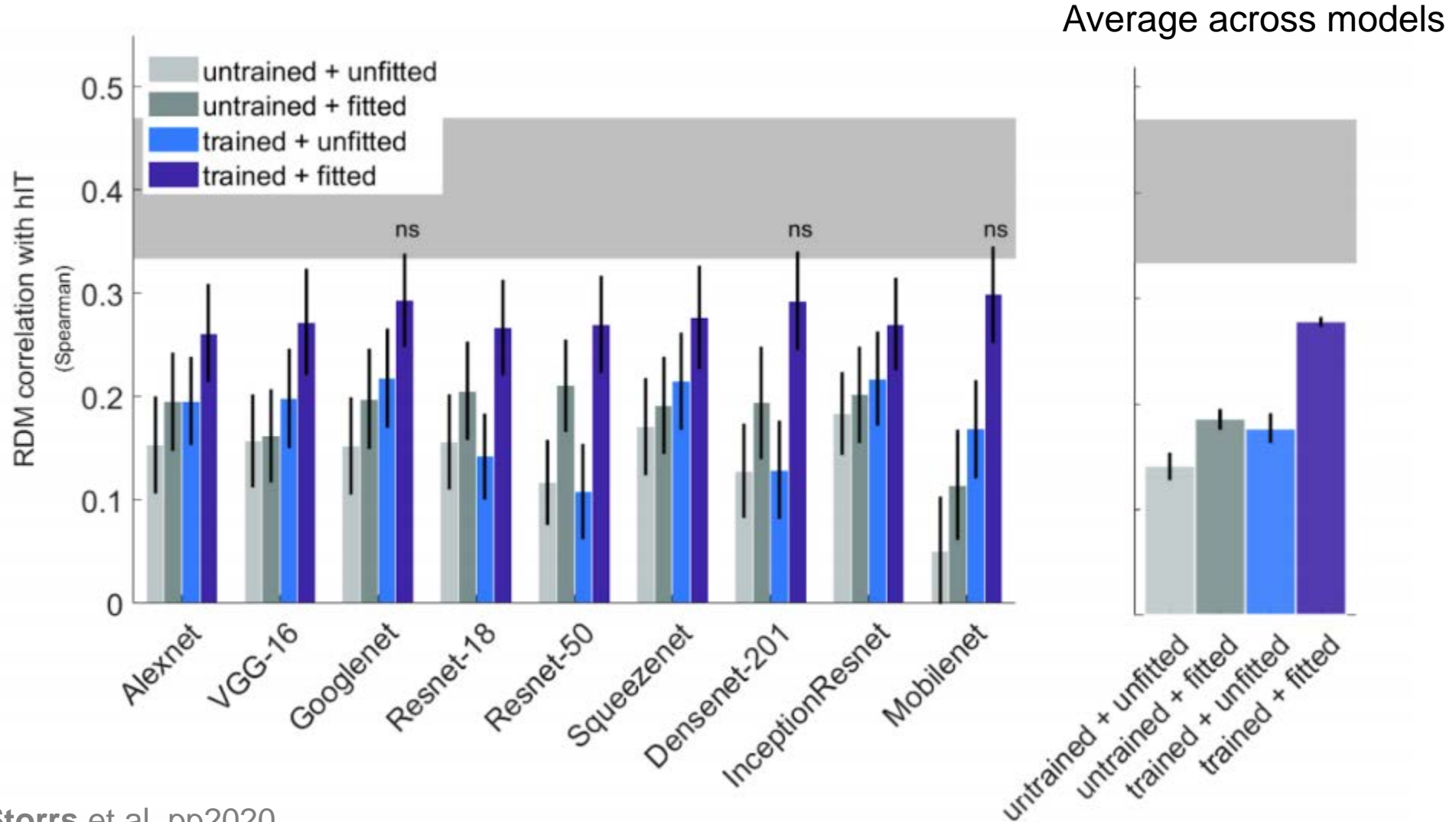


— single neuronal recording site
— deep-net predicted response for novel image

Diverse deep feedforward neural networks predict IT, after *task-training* and *IT-fitting*



Diverse deep feedforward neural networks predict IT, after *task-training* and *IT-fitting*



The converging feedforward story...

- Deep convolutional feedforward neural networks explain how the initial sweep through the primate visual hierarchy enables recognition at a glance.
- They predict representations of novel images better than any alternative current models.
- Both the *architecture* of the model and the *task training* contribute substantially to these successes.

However, we need to build models whose architecture more closely resembles the visual hierarchy.

A major feature of biological neural networks is recurrent signal flow.

Overview

1. Recurrent neural network models
2. Controversial stimuli

Overview

1. Recurrent neural network models

2. Controversial stimuli

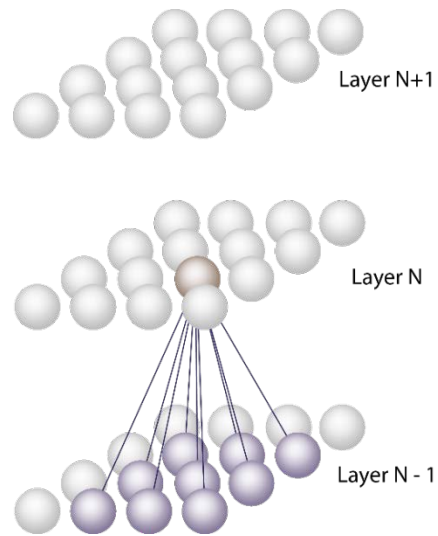
Do *recurrent* convolutional neural networks provide better models of vision?

Courtney Spoerer



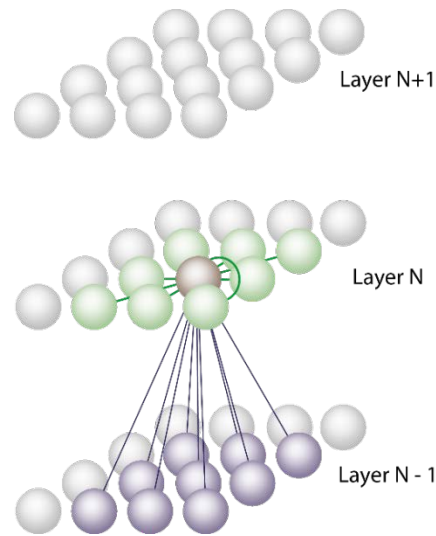
Recurrent convolutional neural networks

B (FF-CNN)



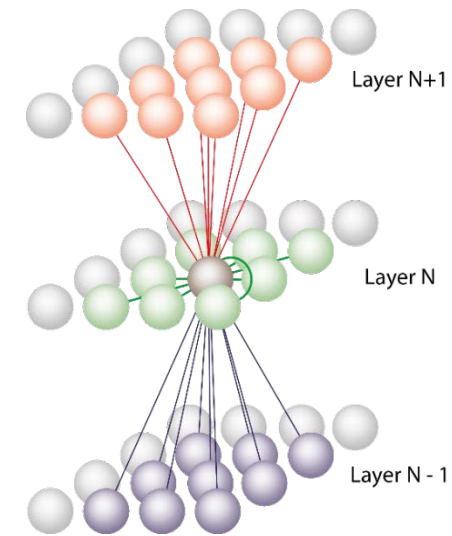
$$\sigma(h_B)$$

BL (RCNN)



$$\sigma(h_B + h_L)$$

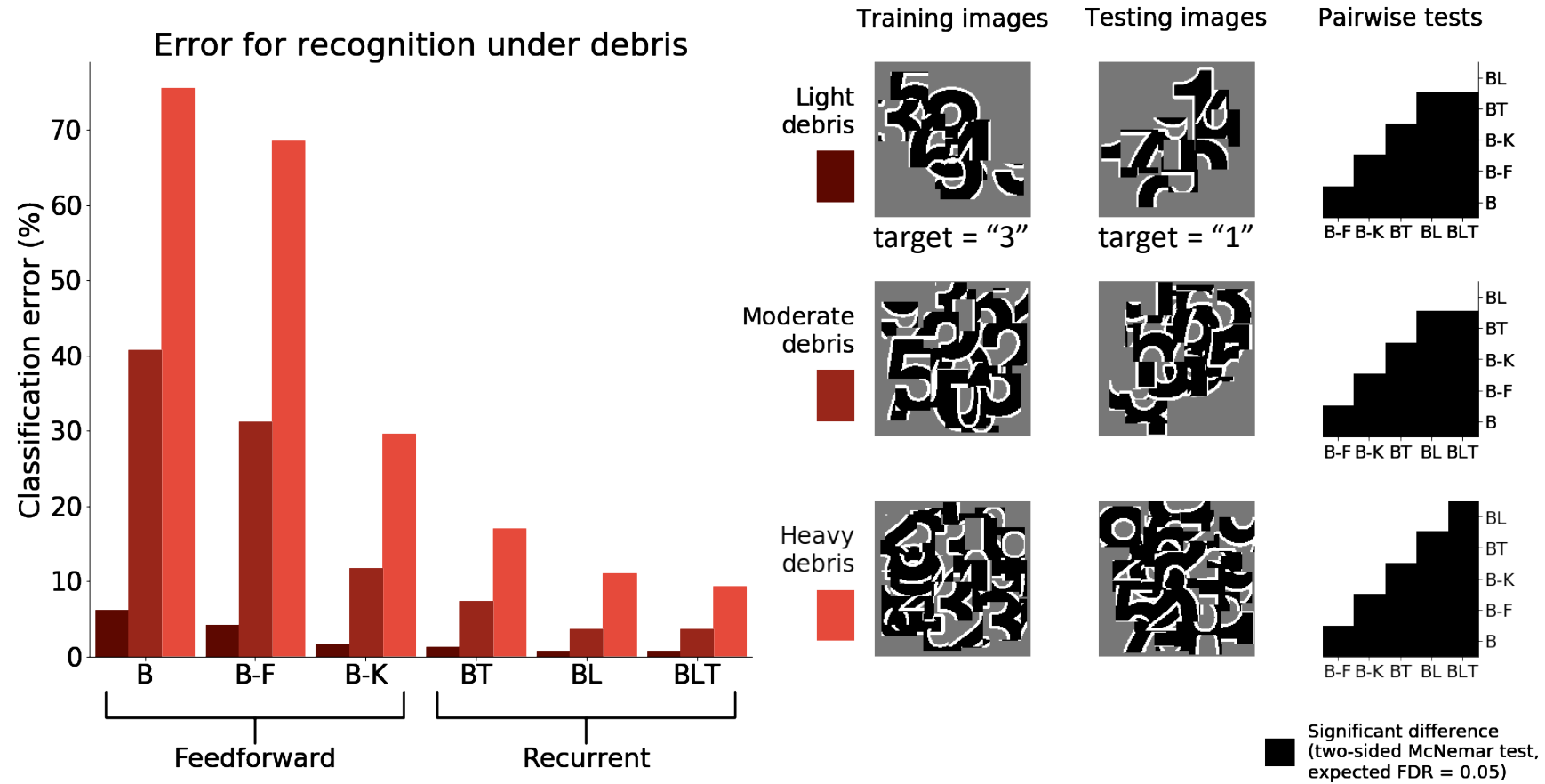
BLT (RCNN)



$$\sigma(h_B + h_L + h_T)$$

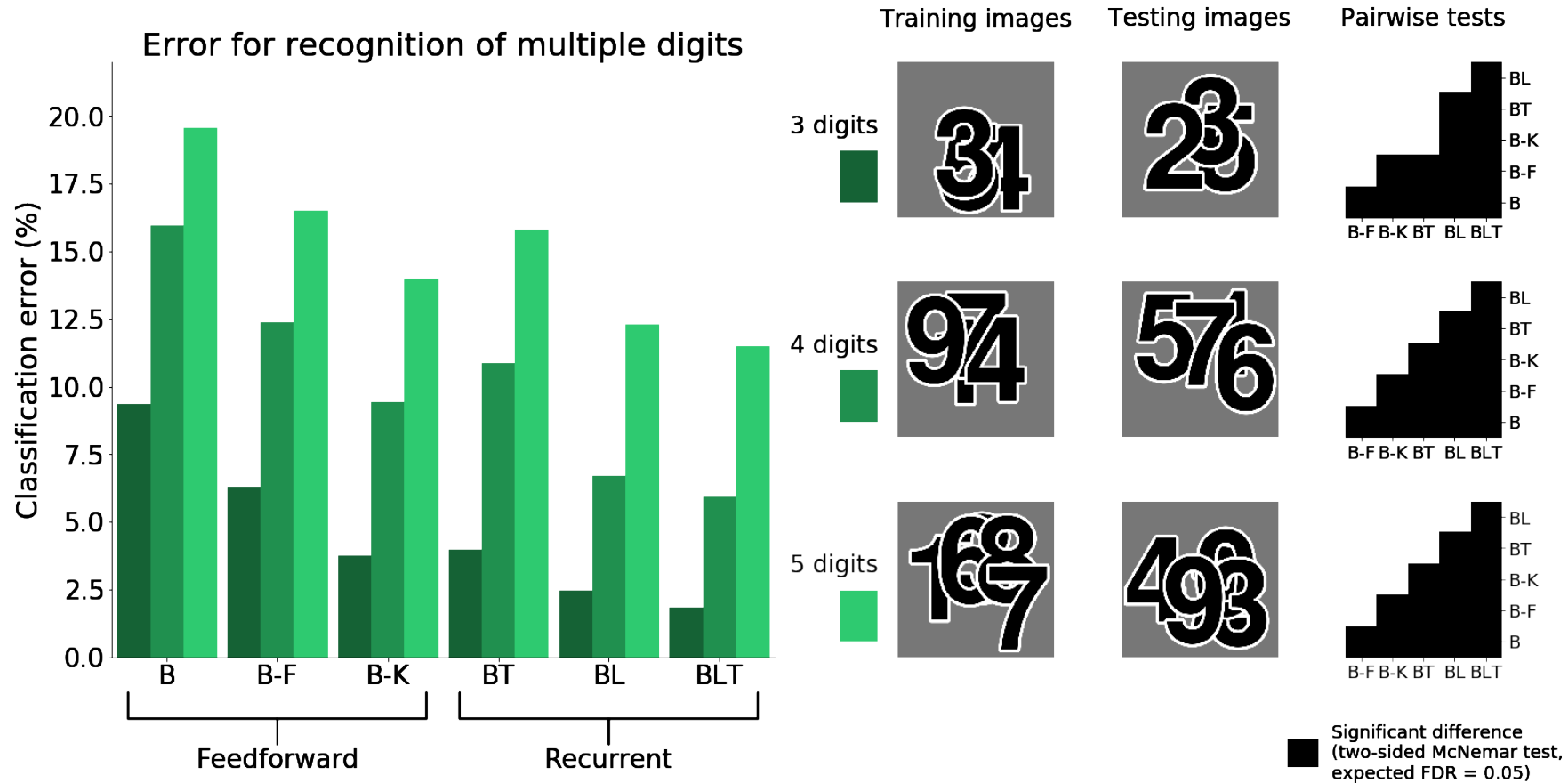
σ - non-linearity
 h_x - convolution

Digit debris: recognition under occlusion



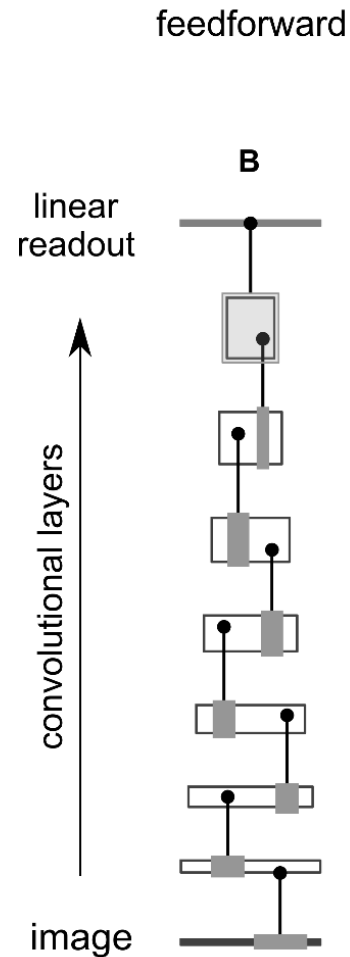
Digit clutter.

Multiple digit recognition

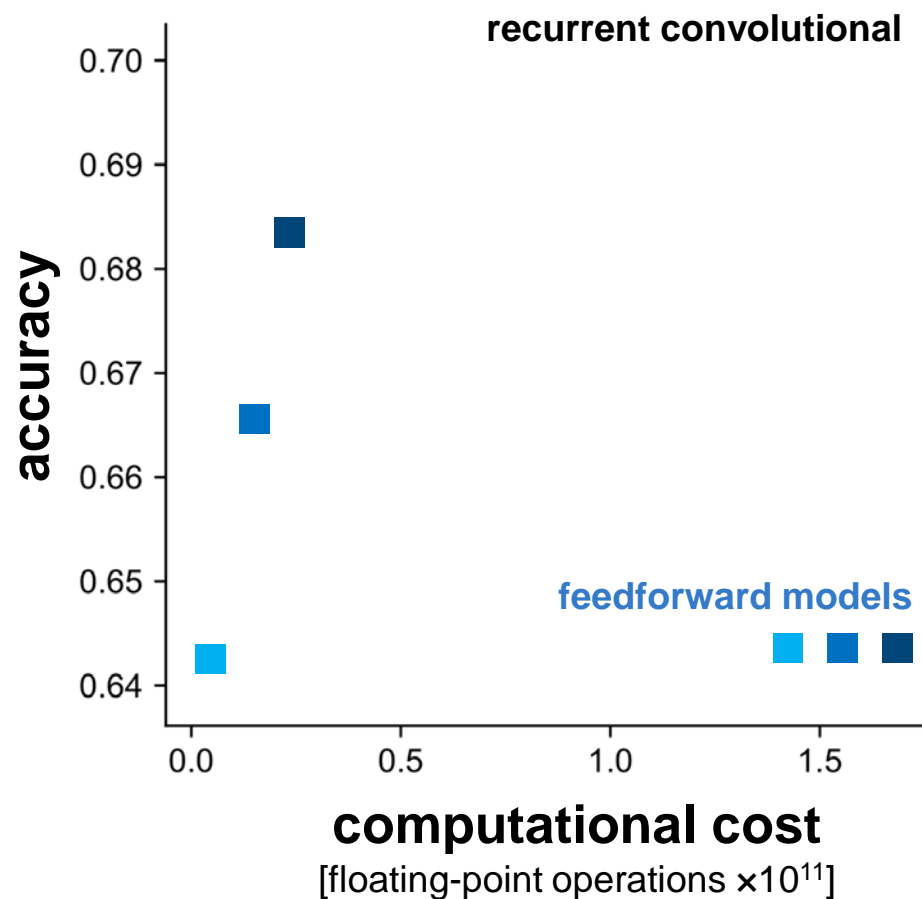


Can recurrent convolutional networks
be scaled up to process
natural images?

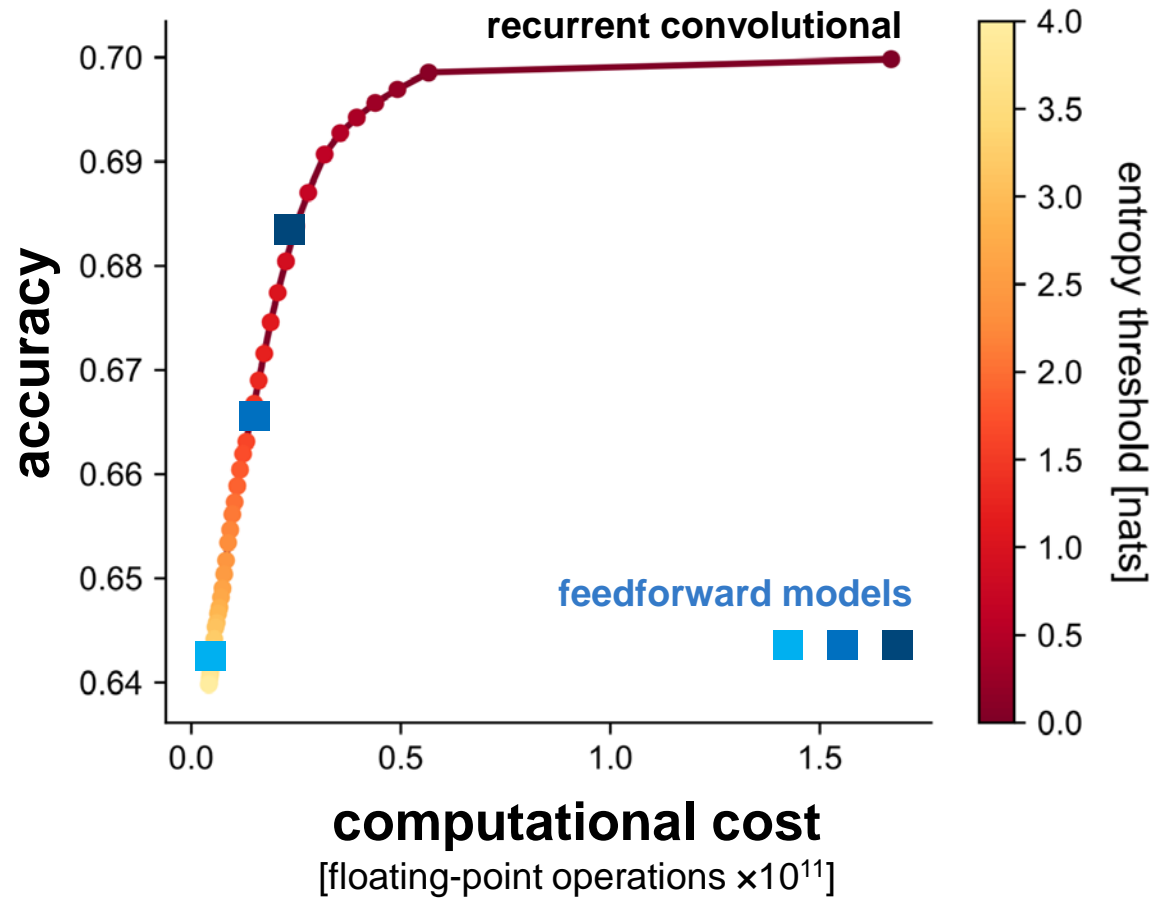
Recurrent convolutional networks trained to recognize natural images



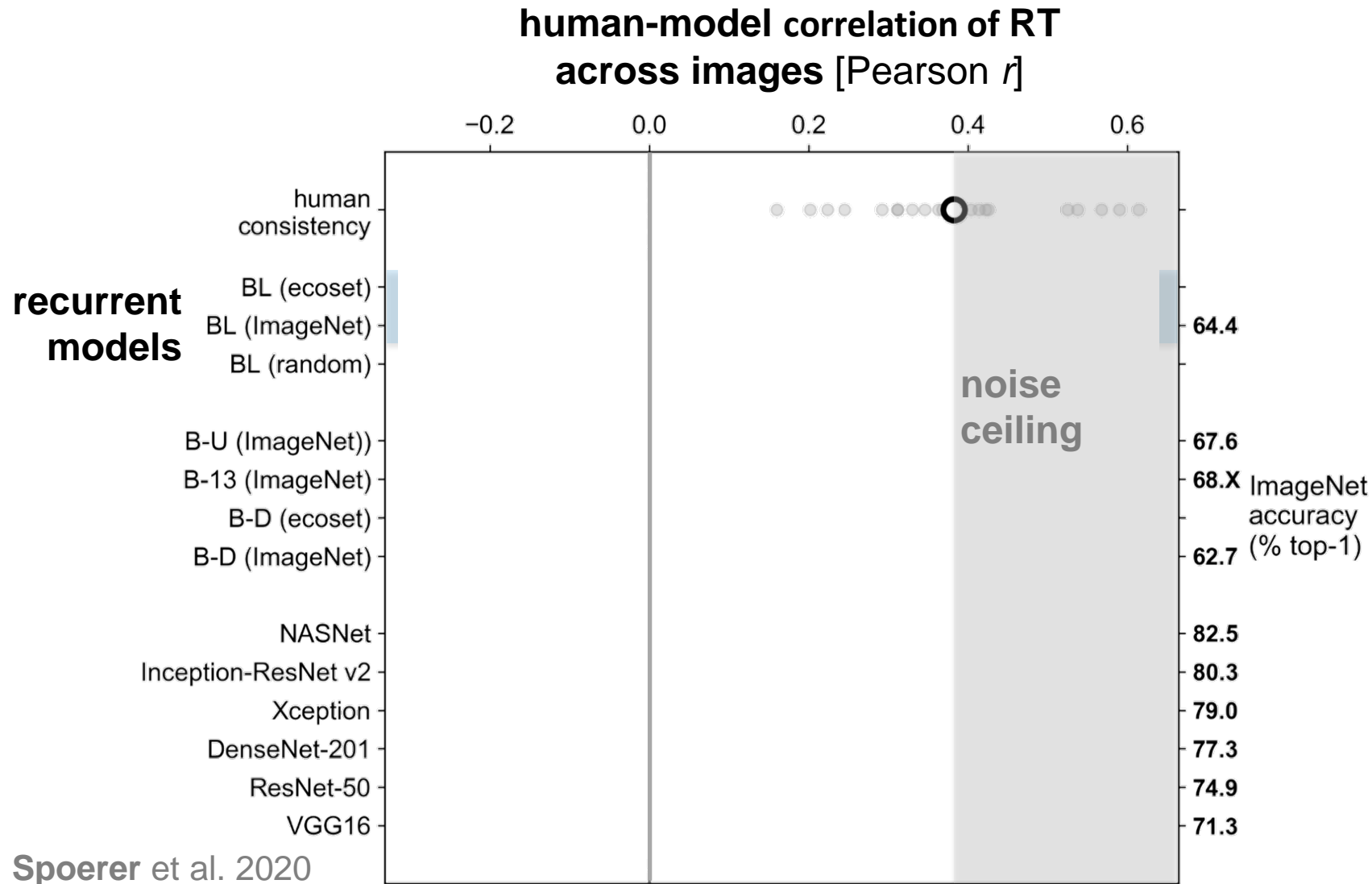
Recurrent models can trade off speed of computation for accuracy



Recurrent models can trade off speed of computation for accuracy



RCNNs predict human reaction times



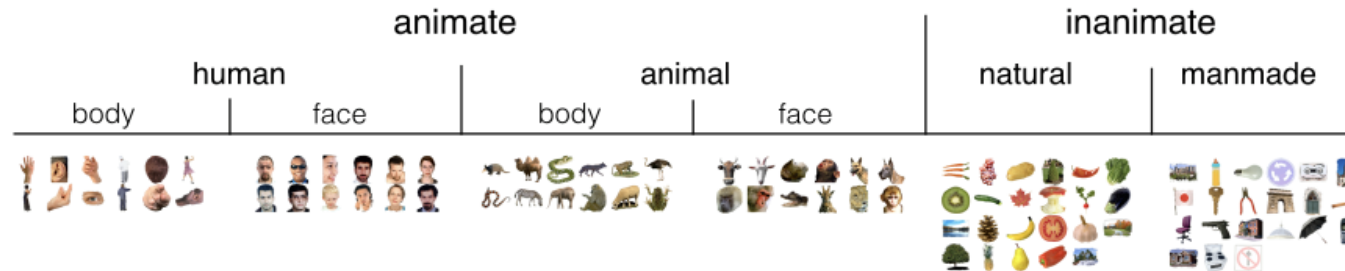
A portrait of a man with dark, curly hair and a light beard, smiling slightly. He is wearing a dark, collared shirt. The background is a solid, dark grey.

Tim Kietzmann

**Can recurrent neural network models
capture the representational dynamics
in the human ventral stream?**

Representational dynamics

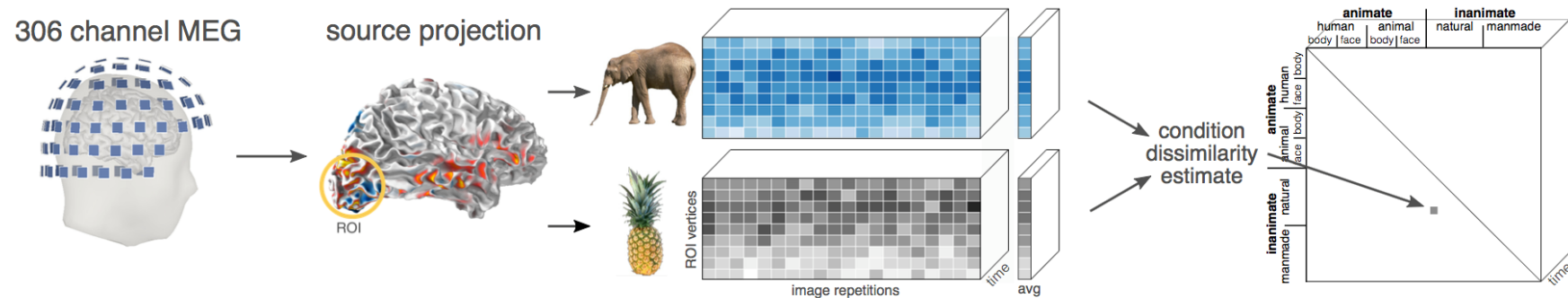
stimuli



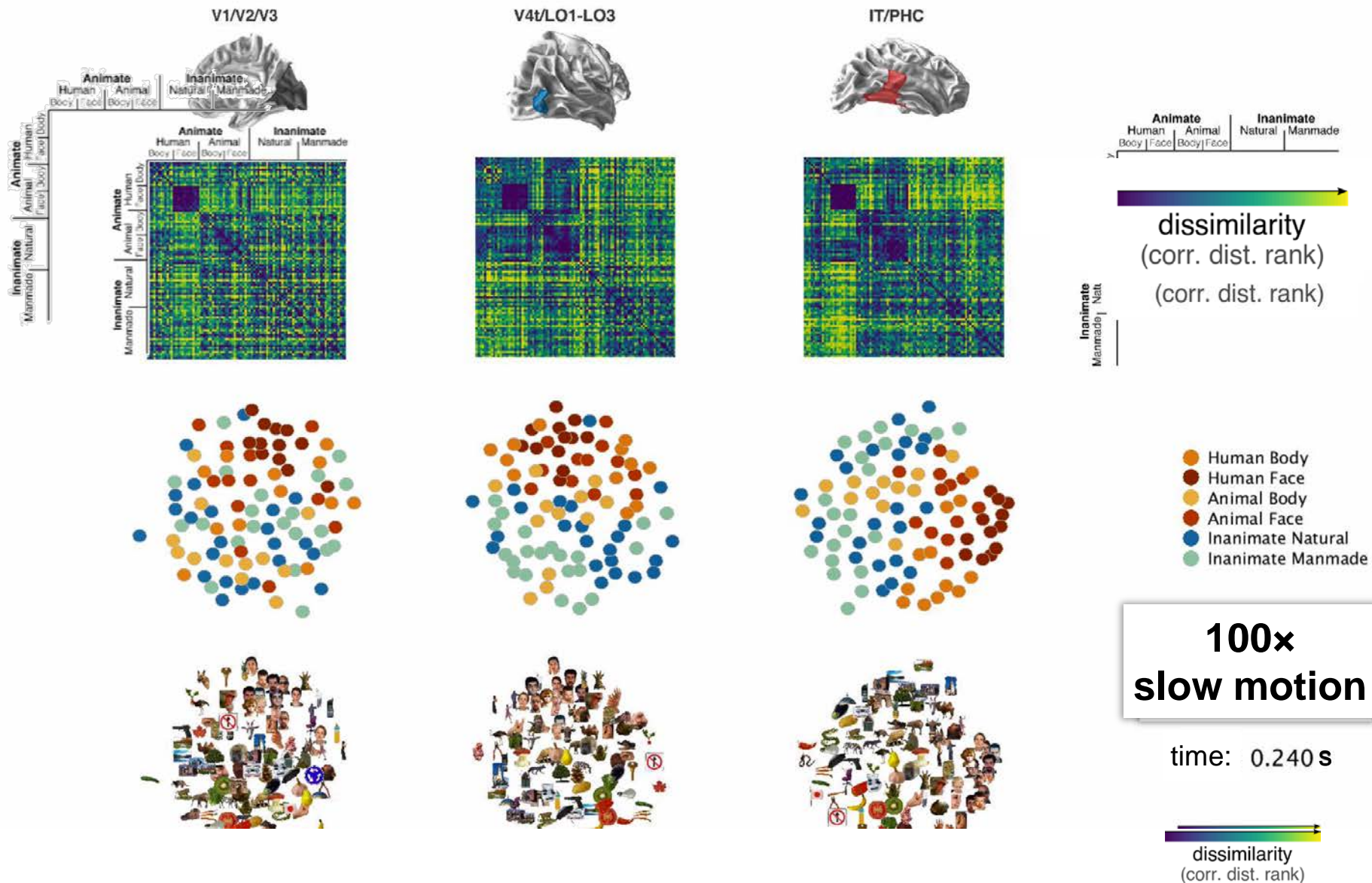
brain regions



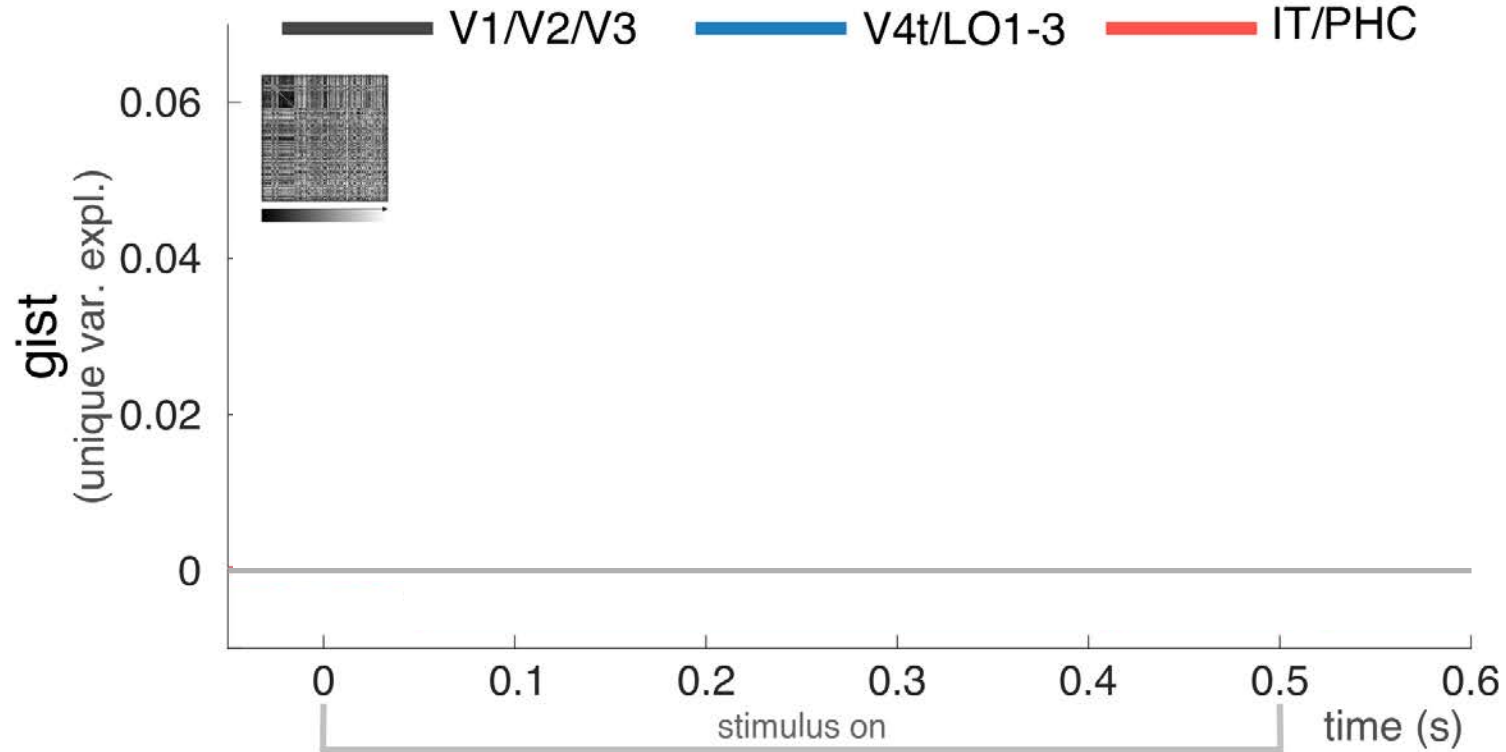
human magnetoencephalography



Movie time

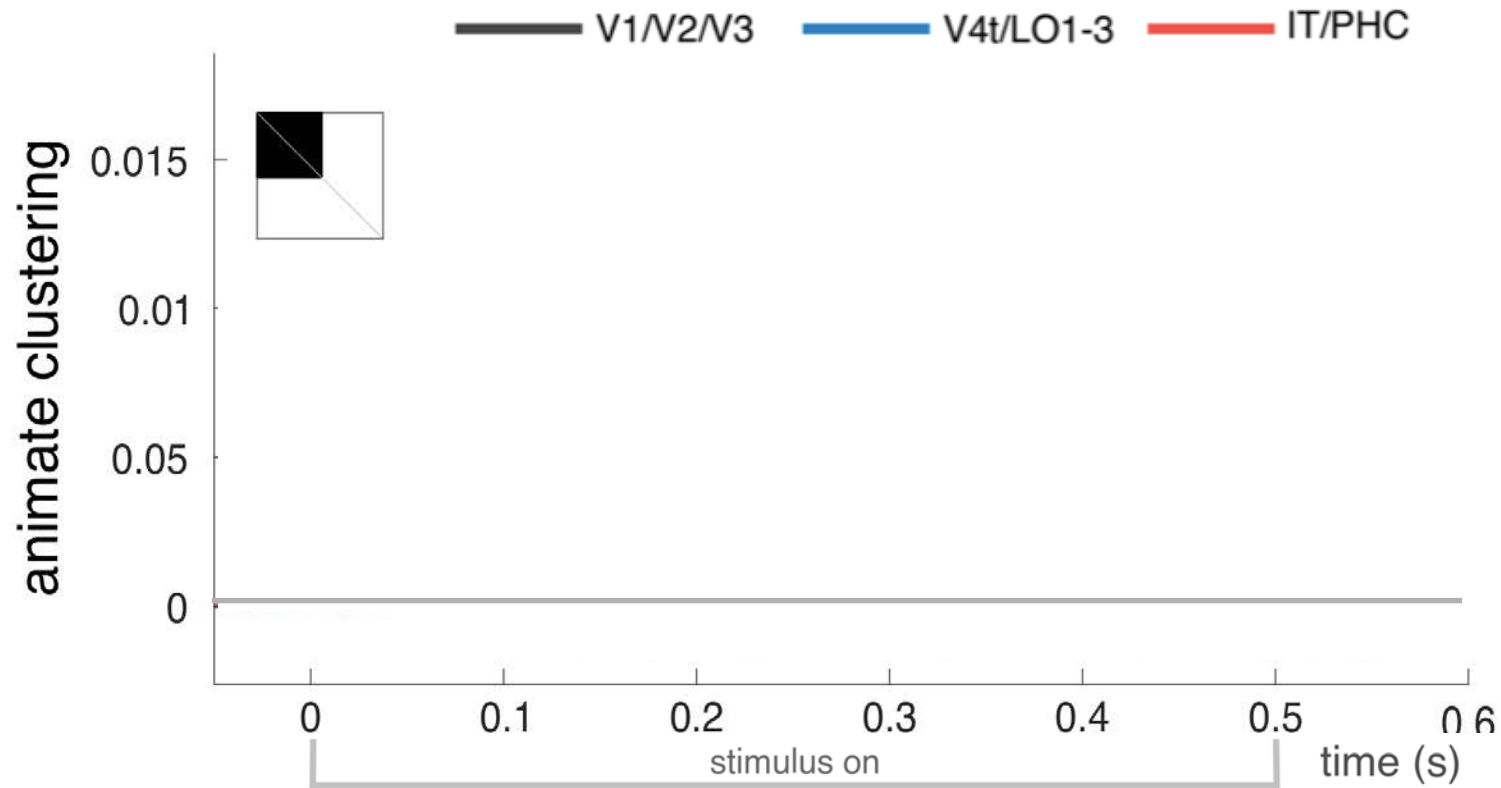


Low-level features: gist model



Gist-like geometries first emerge in early visual areas, where they remain stronger throughout.

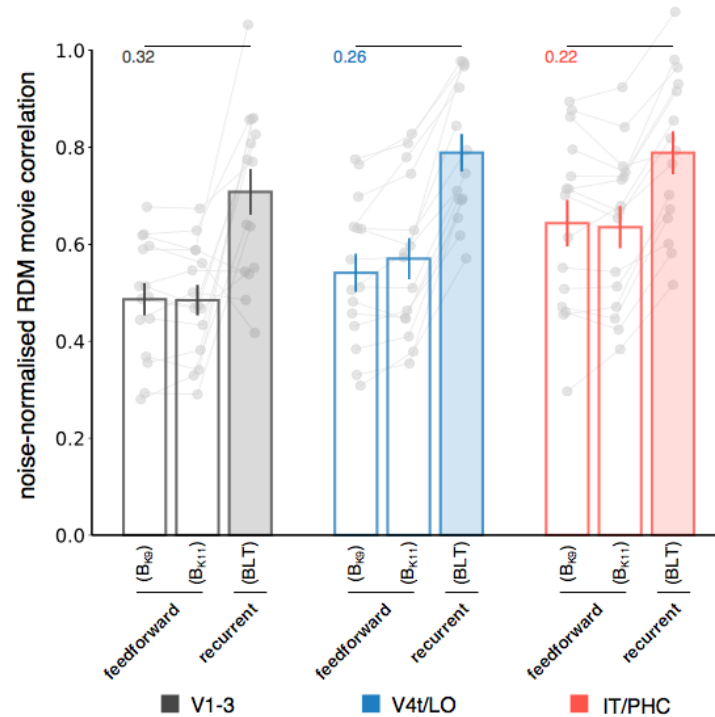
Categorical clustering: animacy



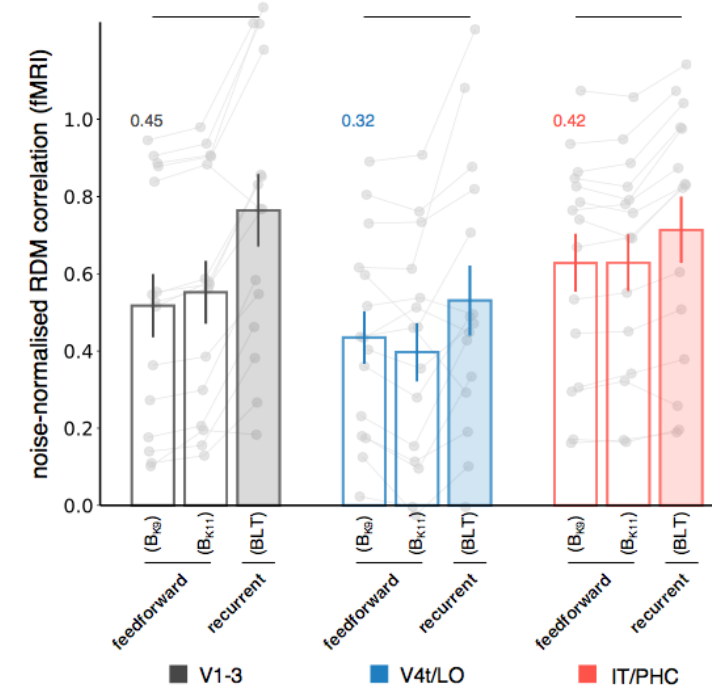
Animacy emerges first in IT/PHC, and only later in V4t/LO1-3.

Recurrent models better explain representations and their dynamics

magnetoencephalography



functional magnetic resonance imaging



□ feedforward

■ recurrent

The emerging recurrent story...

- Recurrent neural networks provide a **more neurobiologically realistic and computationally powerful** modeling framework.
- Recurrent processing can enable a network to
 - **recycle its computational resources**,
 - perform **more robust inferences**, and
 - **flexibly trade off speed and accuracy**.
- Recurrent models also **better explain the representational dynamics** of the human ventral stream.

Controversial stimuli:
pitting neural networks
against each other
as models of
human recognition



Tal Golan

Controversial stimuli

Controversial stimuli: motivation

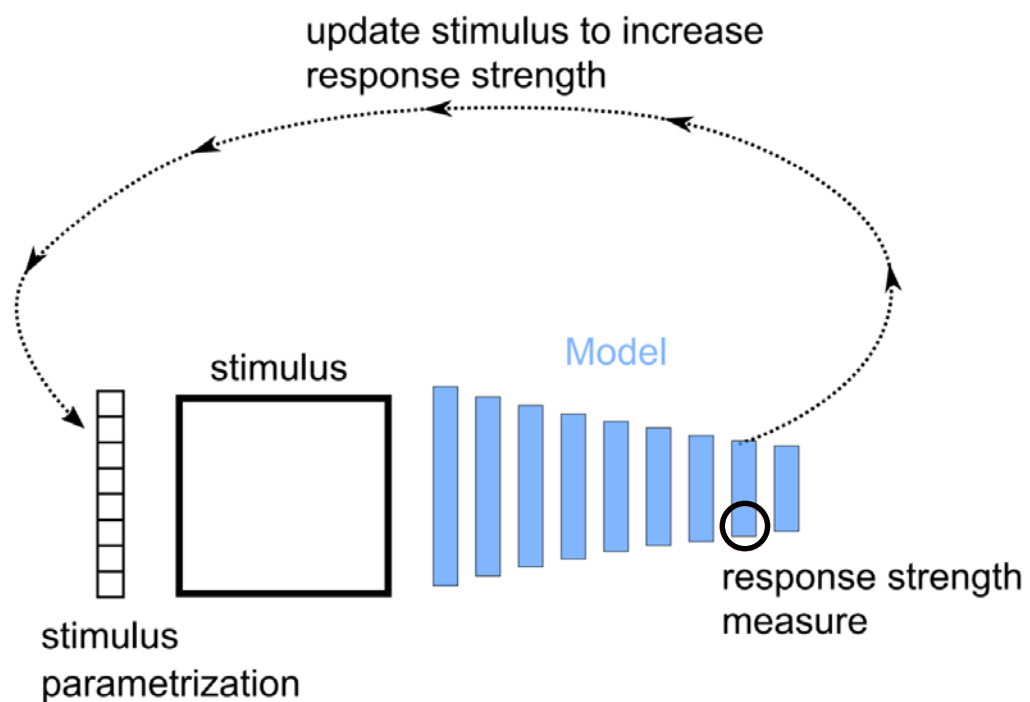
- Theoretical progress depends on experiments for which competing theories make distinct predictions.
- We can implement competing theories in testable NN models.
- However, NN models have many parameters, and theoretically distinct models often make similar predictions for natural stimuli.

Insight 1: To elicit models' distinct *inductive biases* we can test models on a population of stimuli not used in training (*out of distribution*).

- natural stimuli drawn from a different stimulus population
- synthetic stimuli (optimized to elicit bolder predictions, e.g. superstimuli, adversarial stimuli, and metamers)

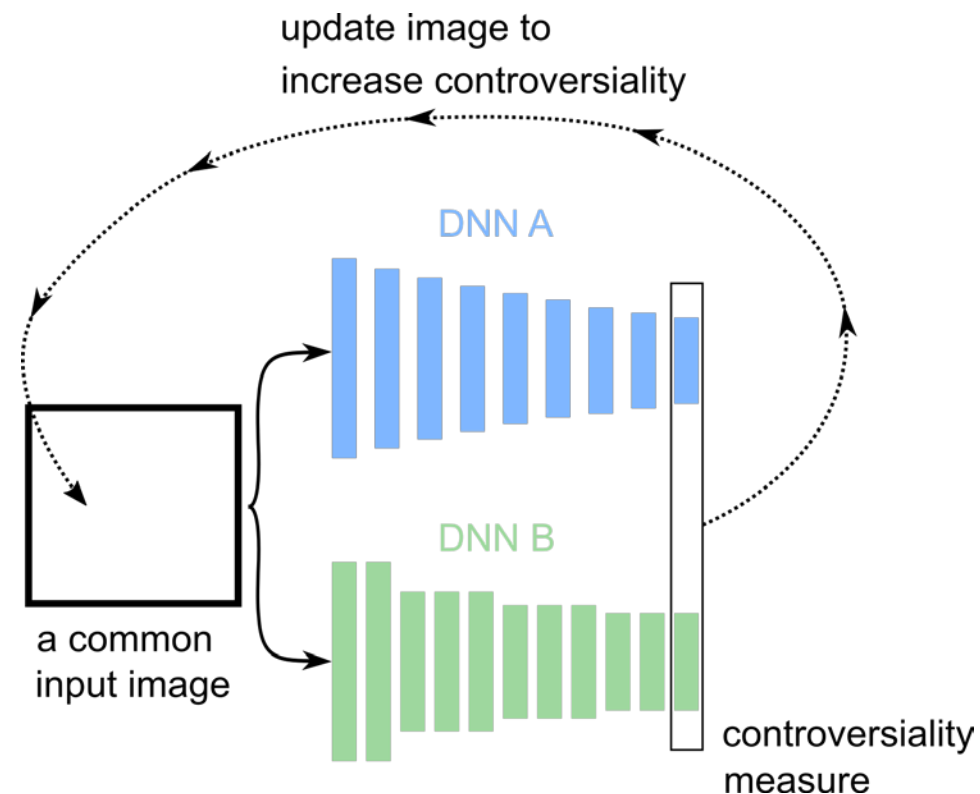
Insight 2: Since our goal is to adjudicate among models, we can create synthetic stimuli optimized to elicit distinct predictions from different models: stimuli that are ***controversial*** among the models.

Superstimulus



Abbasi-Asl et al. 2018, Malakhova 2018,
Ponce et al. 2019, Bashivan et al. 2019,
Walker et al. 2019

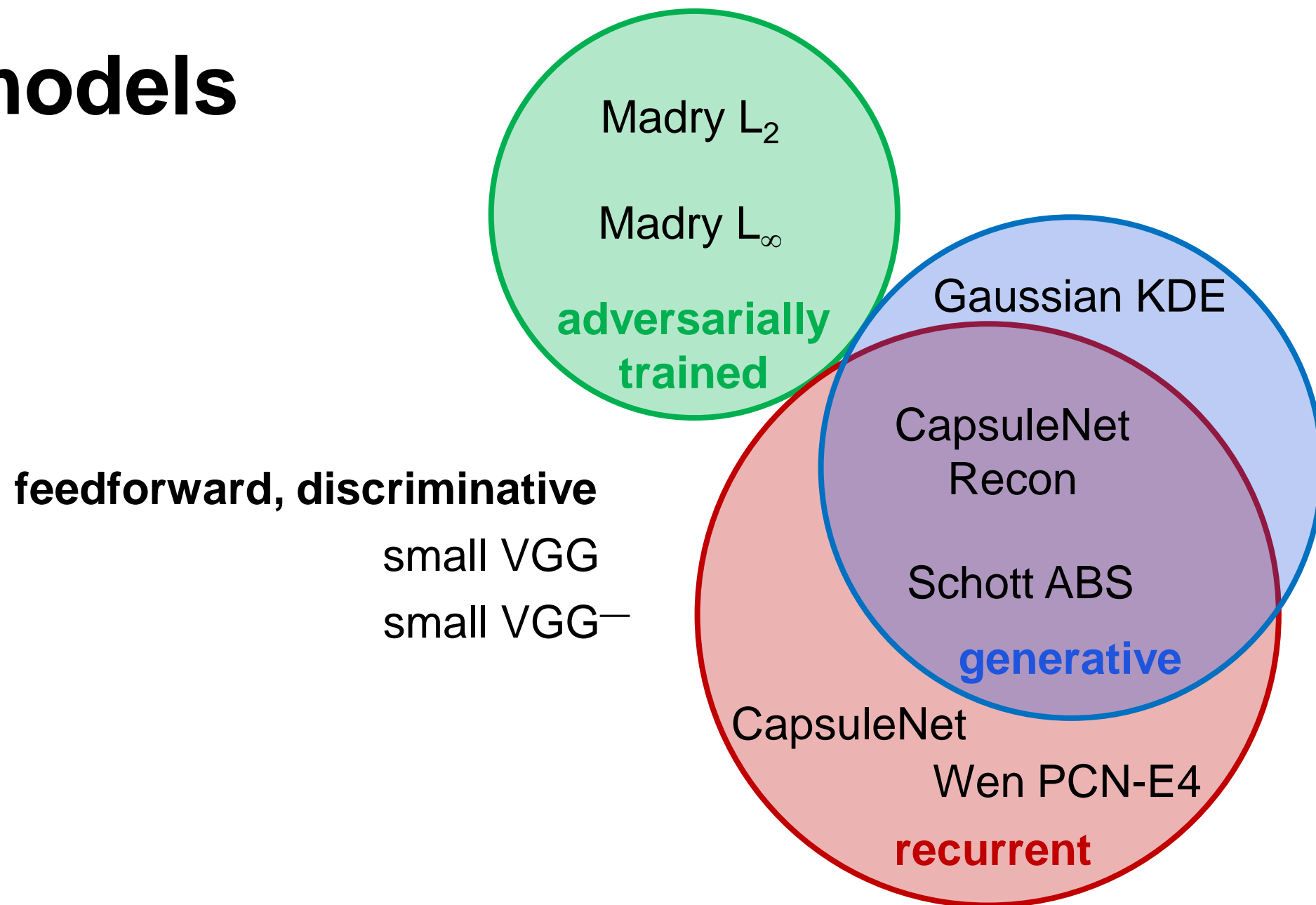
Controversial stimulus



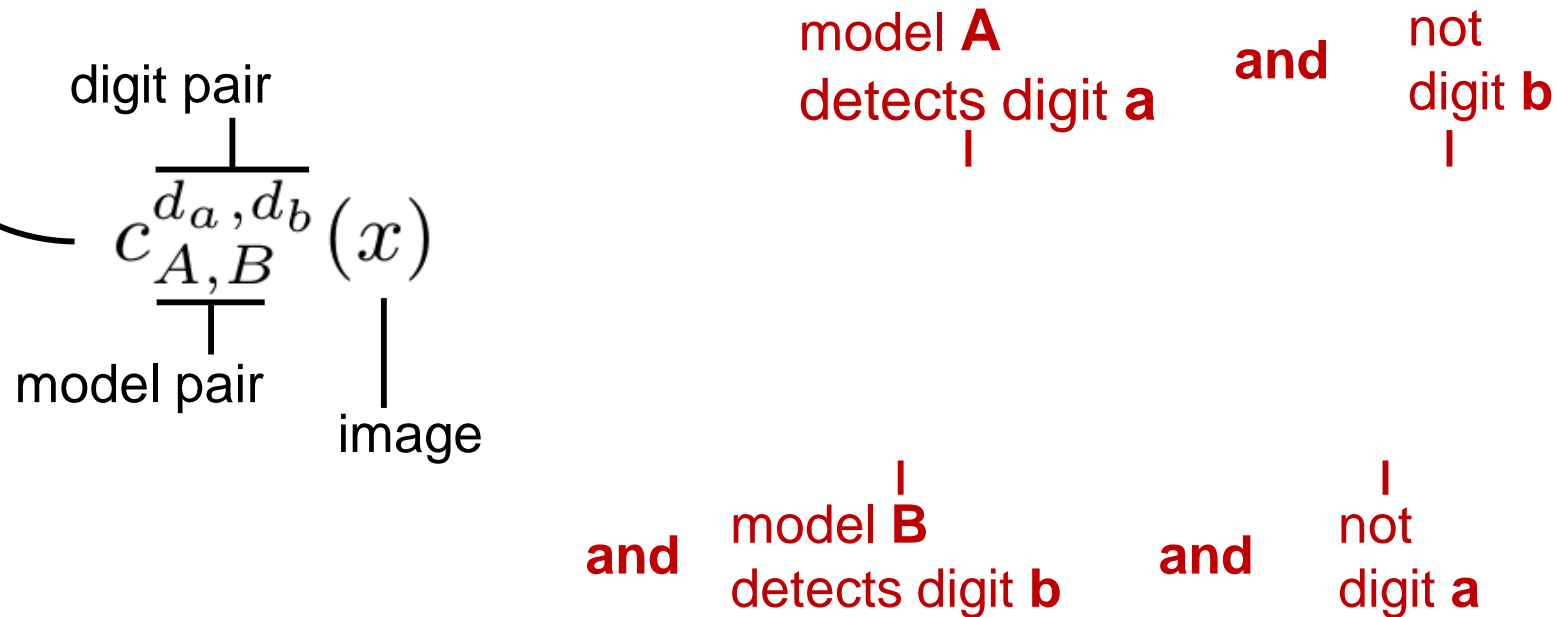
Golan et al. 2020

[illegible]

Tested models



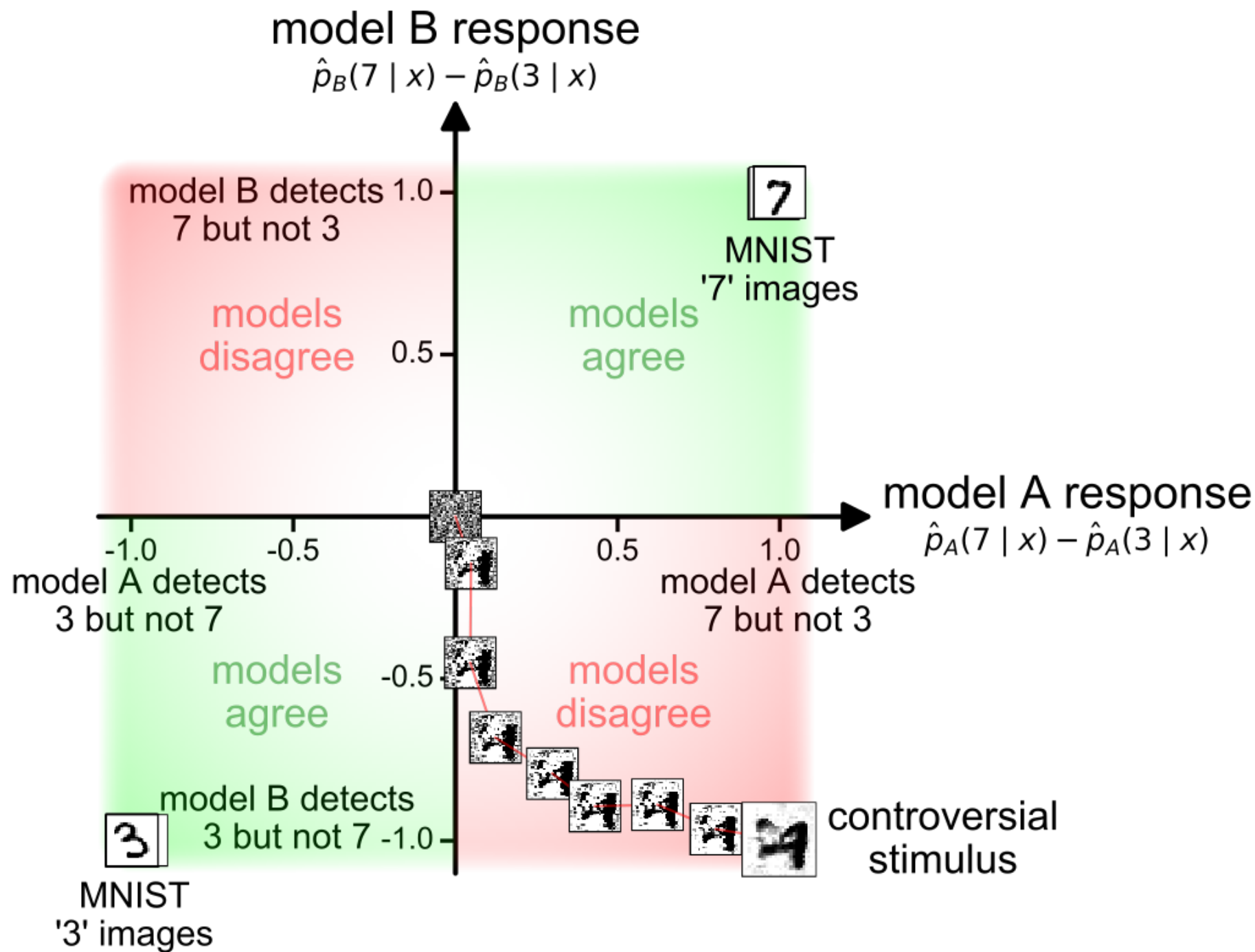
Controversiality index



Controversiality index

The diagram illustrates the components of the controversiality index formula. A bracket labeled "digit pair" spans d_a, d_b . A bracket labeled "model pair" spans A, B . A vertical line labeled "image" points to x . The formula is:

$$c_{A,B}^{d_a, d_b}(x) = \min \left\{ \overbrace{\hat{p}_A(d_a | x)}^{\text{model A detects digit a}}, \overbrace{1 - \hat{p}_A(d_b | x)}^{\text{and not digit b}}, \right. \\ \left. \overbrace{\hat{p}_B(d_b | x)}^{\text{and model B detects digit b}}, \overbrace{1 - \hat{p}_B(d_a | x)}^{\text{and not digit a}} \right\},$$



Controversial stimuli

(optimized by gradient descent)

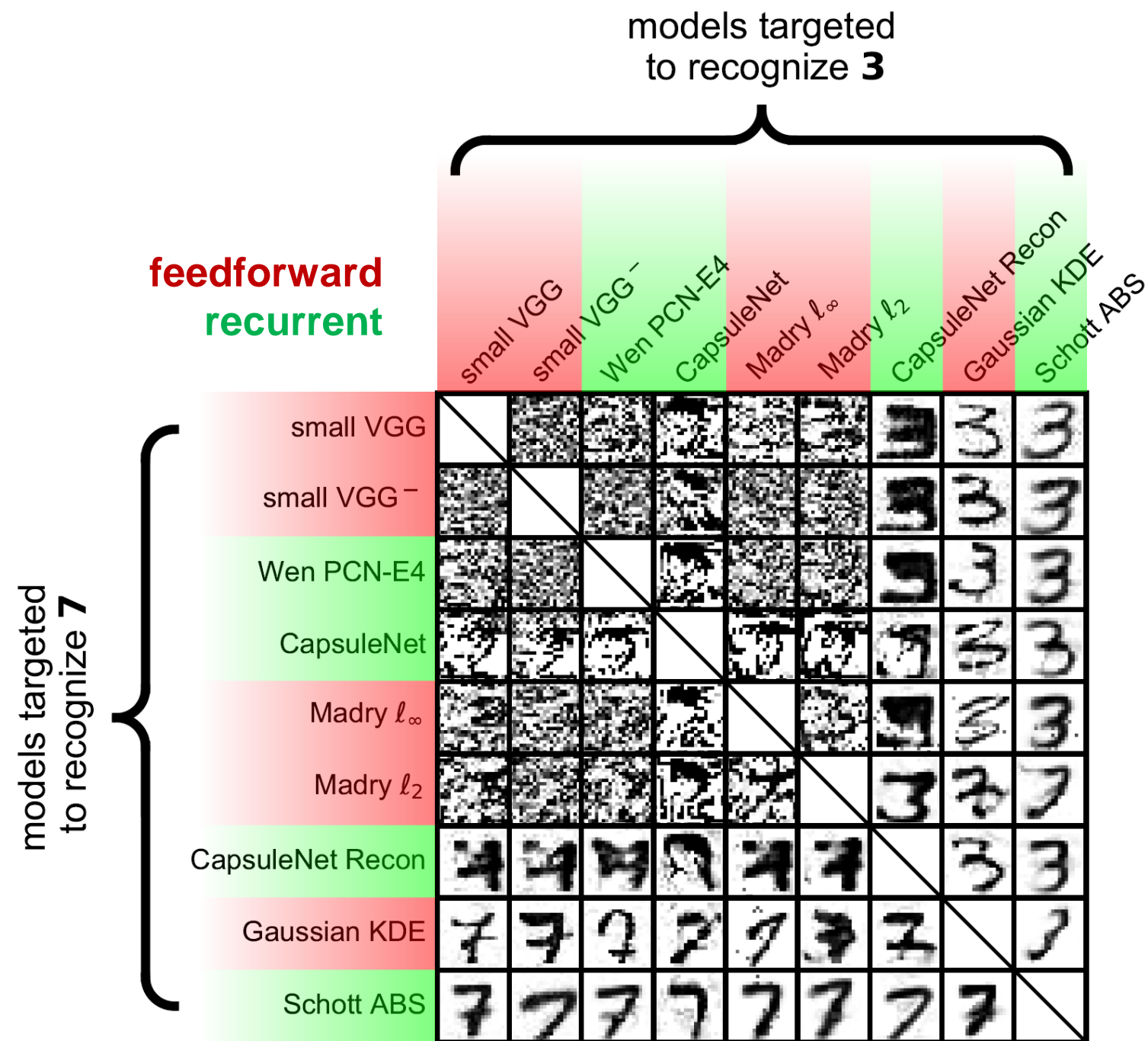
adversarial example
a stimulus that is controversial
between a model and ground truth

models targeted
to recognize **7**

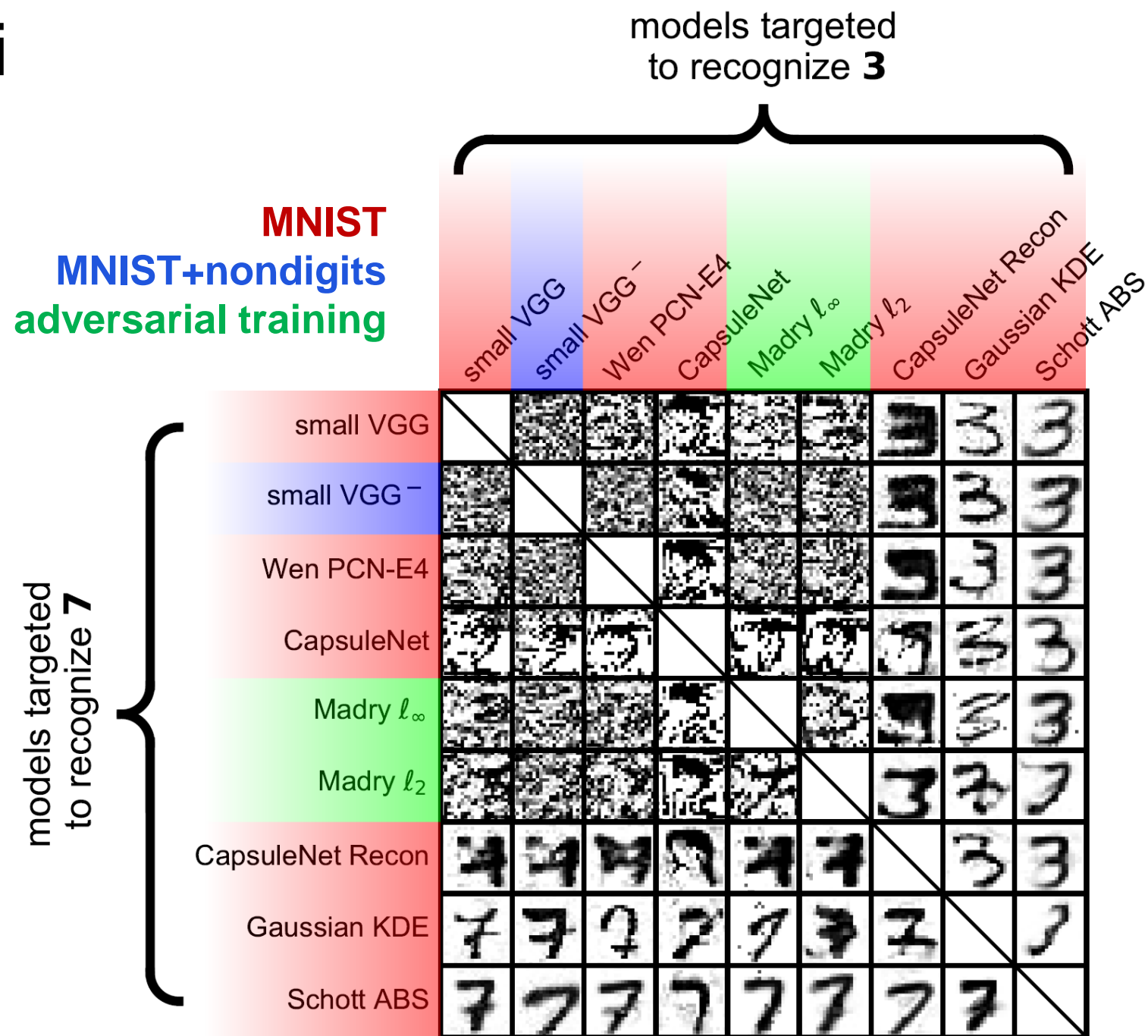
models targeted
to recognize **3**

	small VGG	small VGG ⁻	Wen PCN-E4	CapsuleNet	Madry ℓ_∞	Madry ℓ_2	CapsuleNet Recon	Gaussian KDE	Schott ABS
small VGG									
small VGG ⁻									
Wen PCN-E4									
CapsuleNet									
Madry ℓ_∞									
Madry ℓ_2									
CapsuleNet Recon									
Gaussian KDE									
Schott ABS									

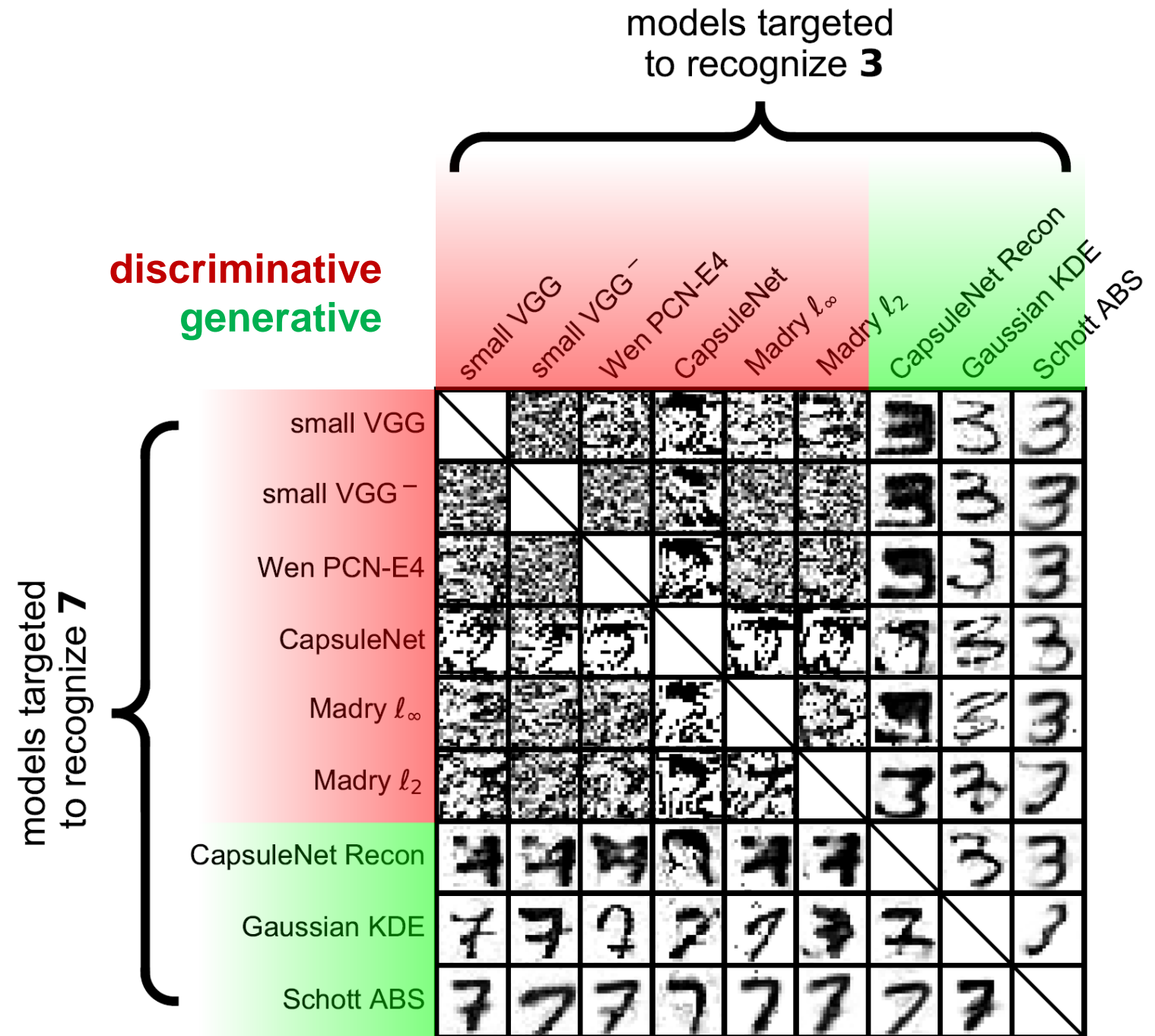
Controversial stimuli

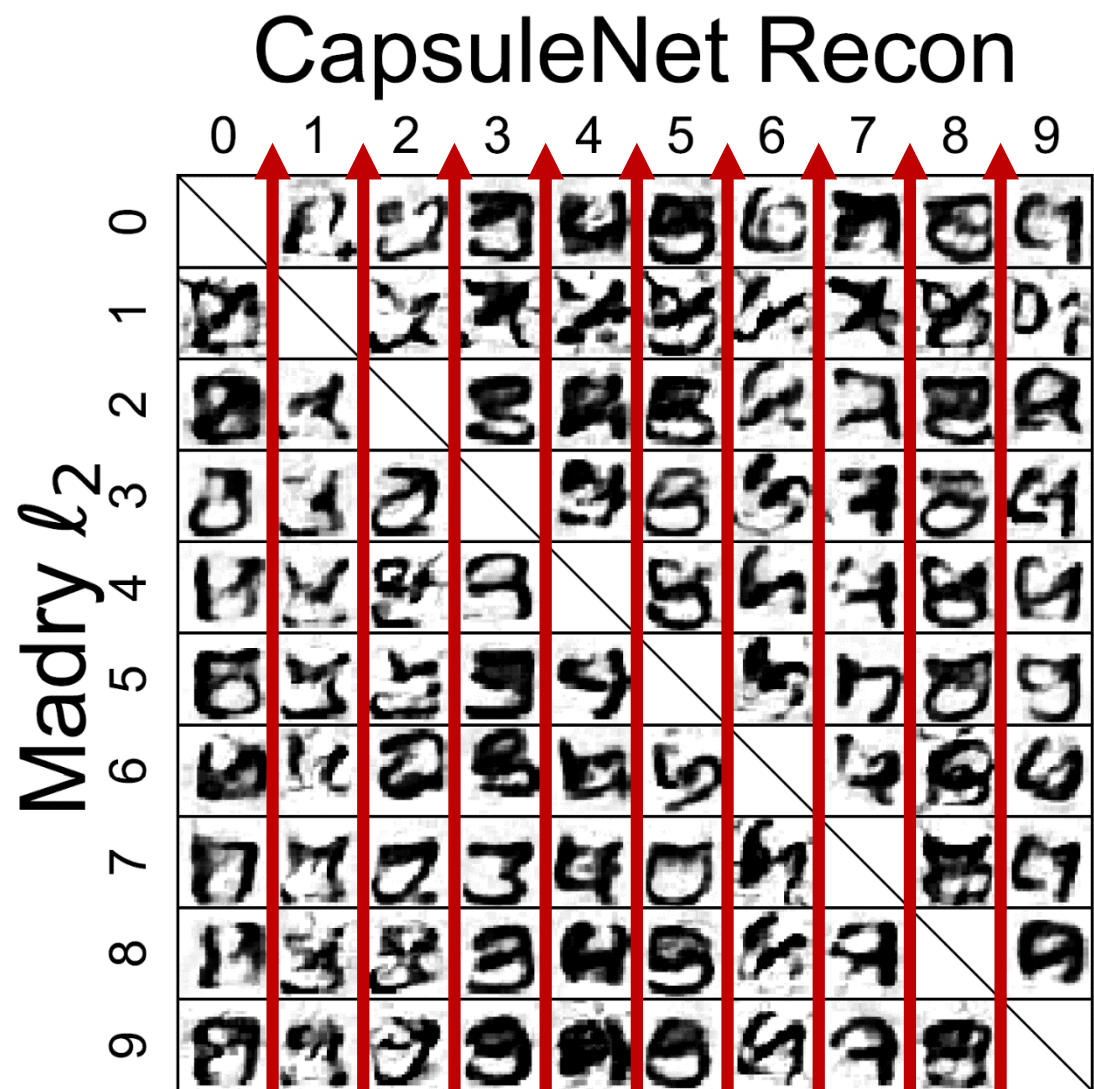
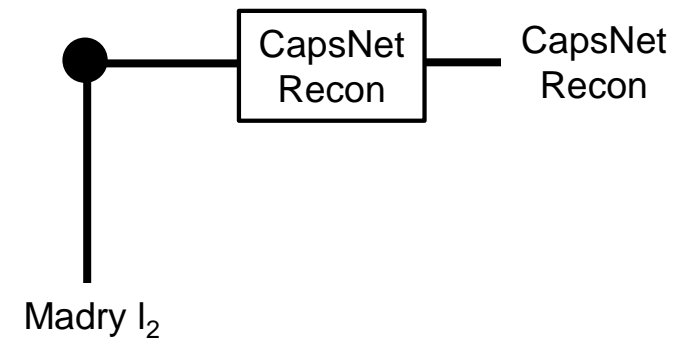


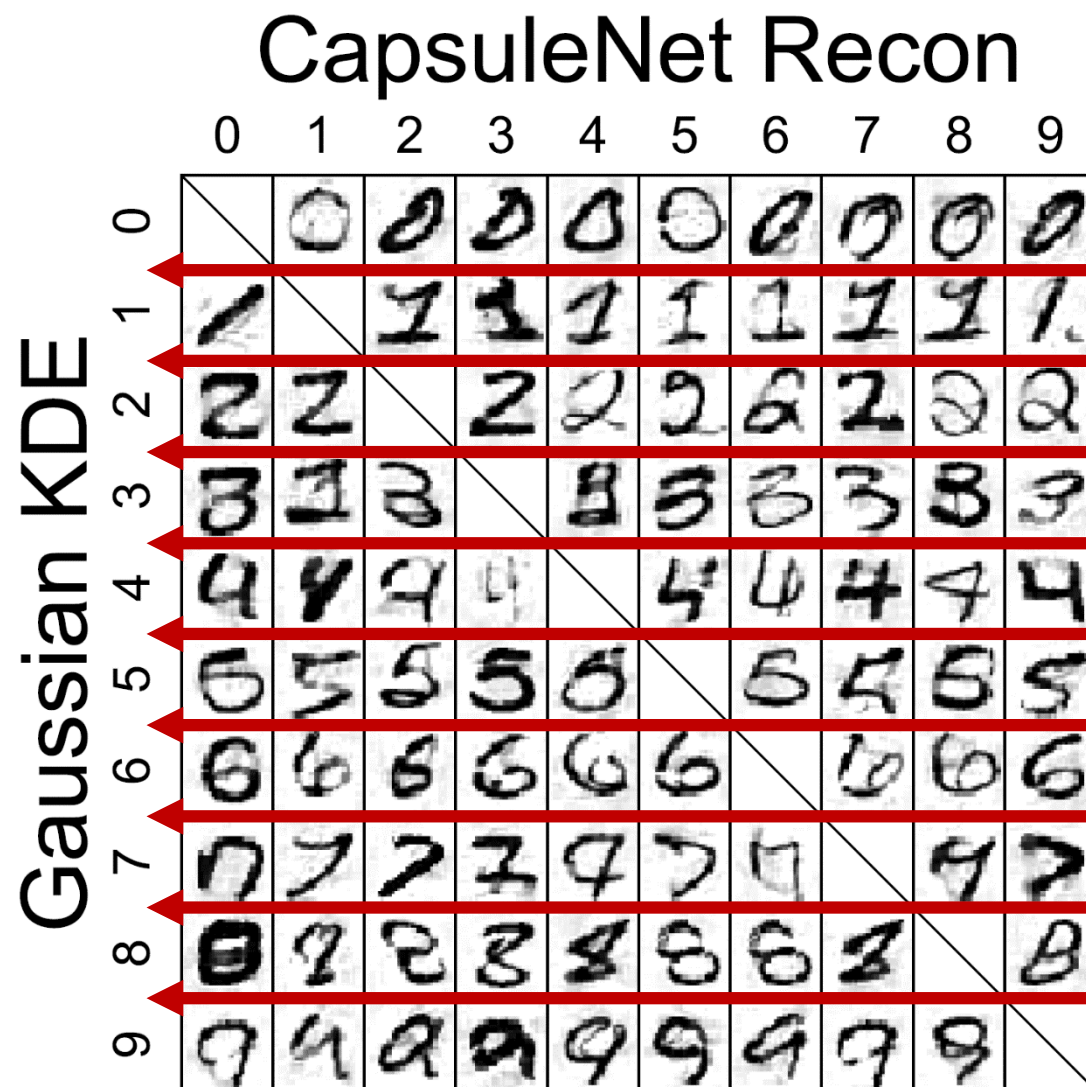
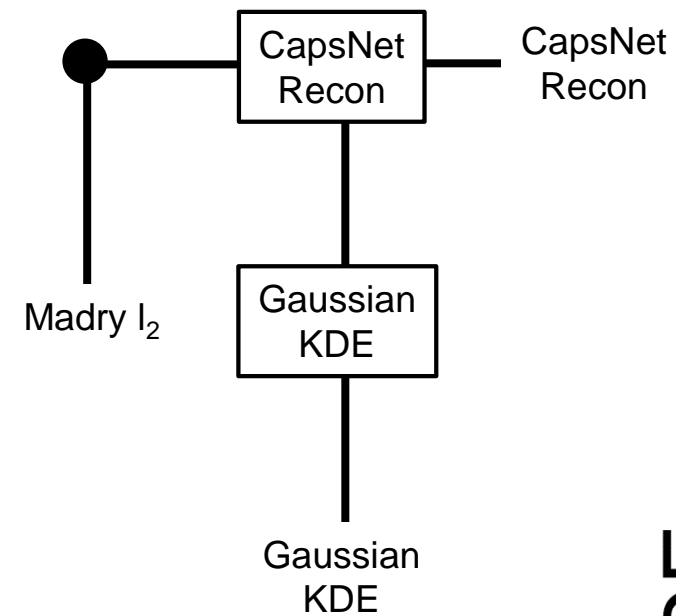
Controversial stimuli

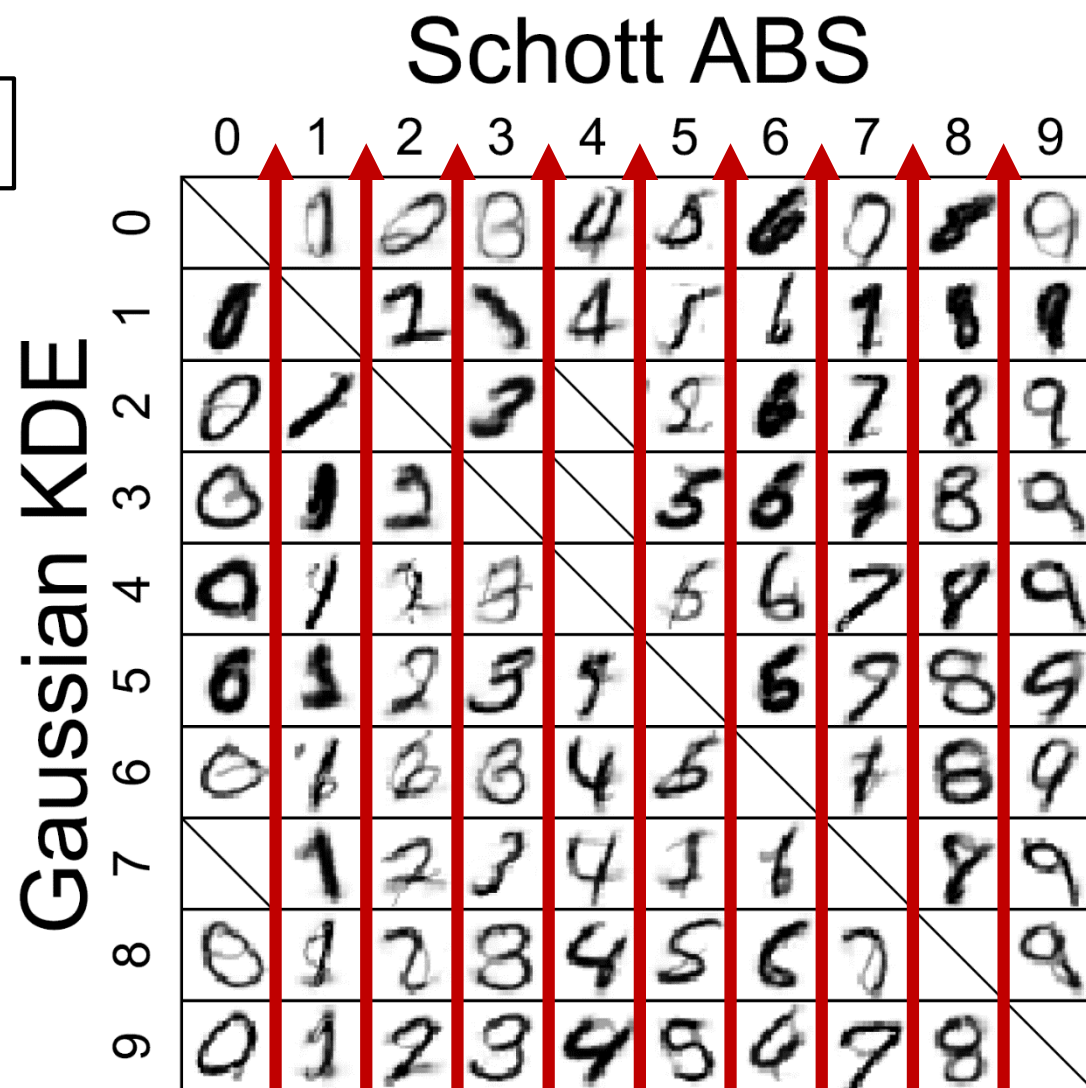
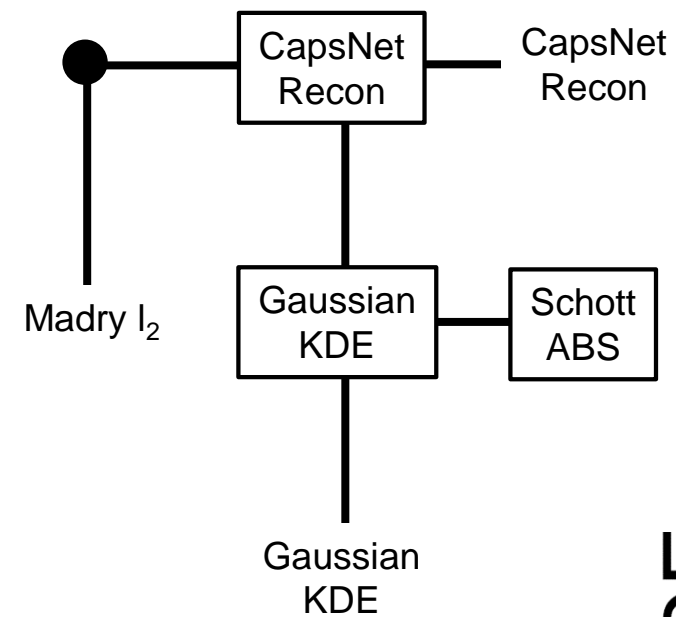


Controversial stimuli

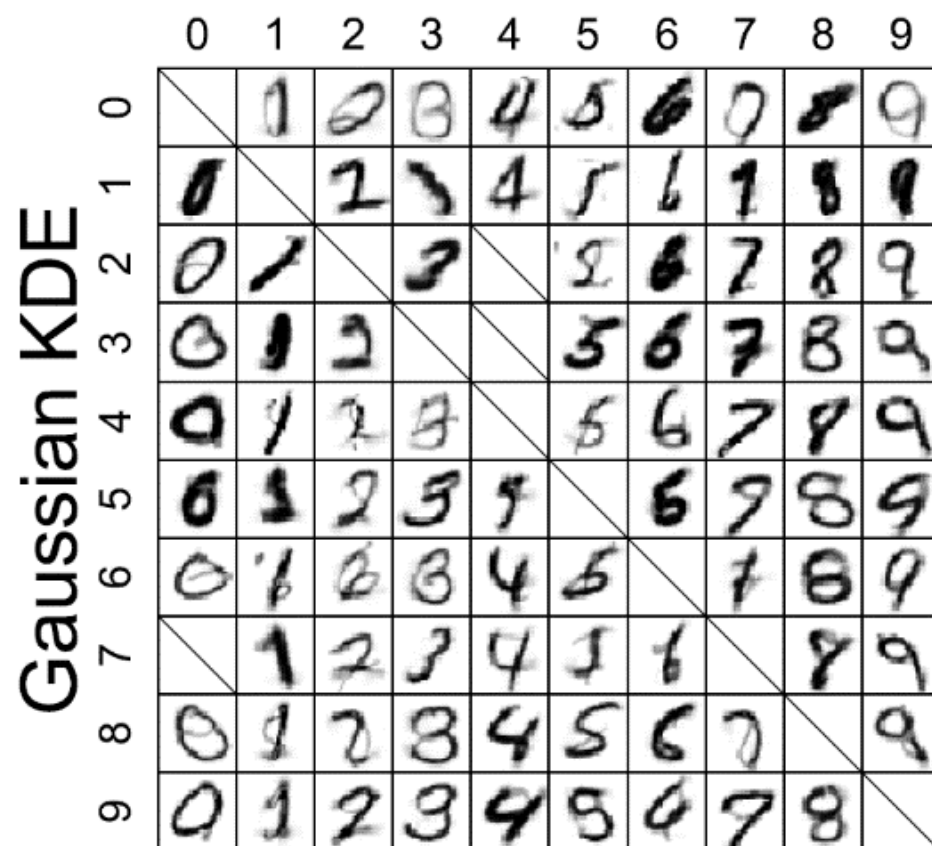






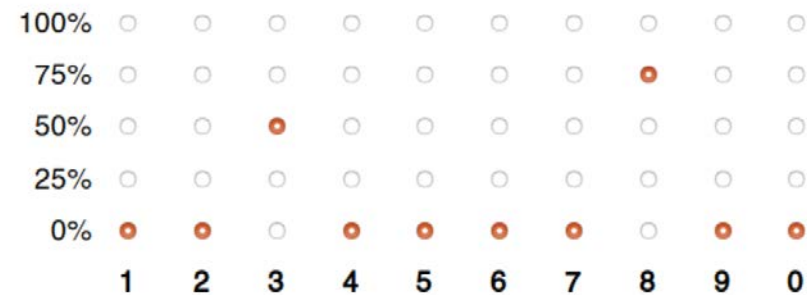


Schott ABS



Behavioral experiment

What number does this look like?



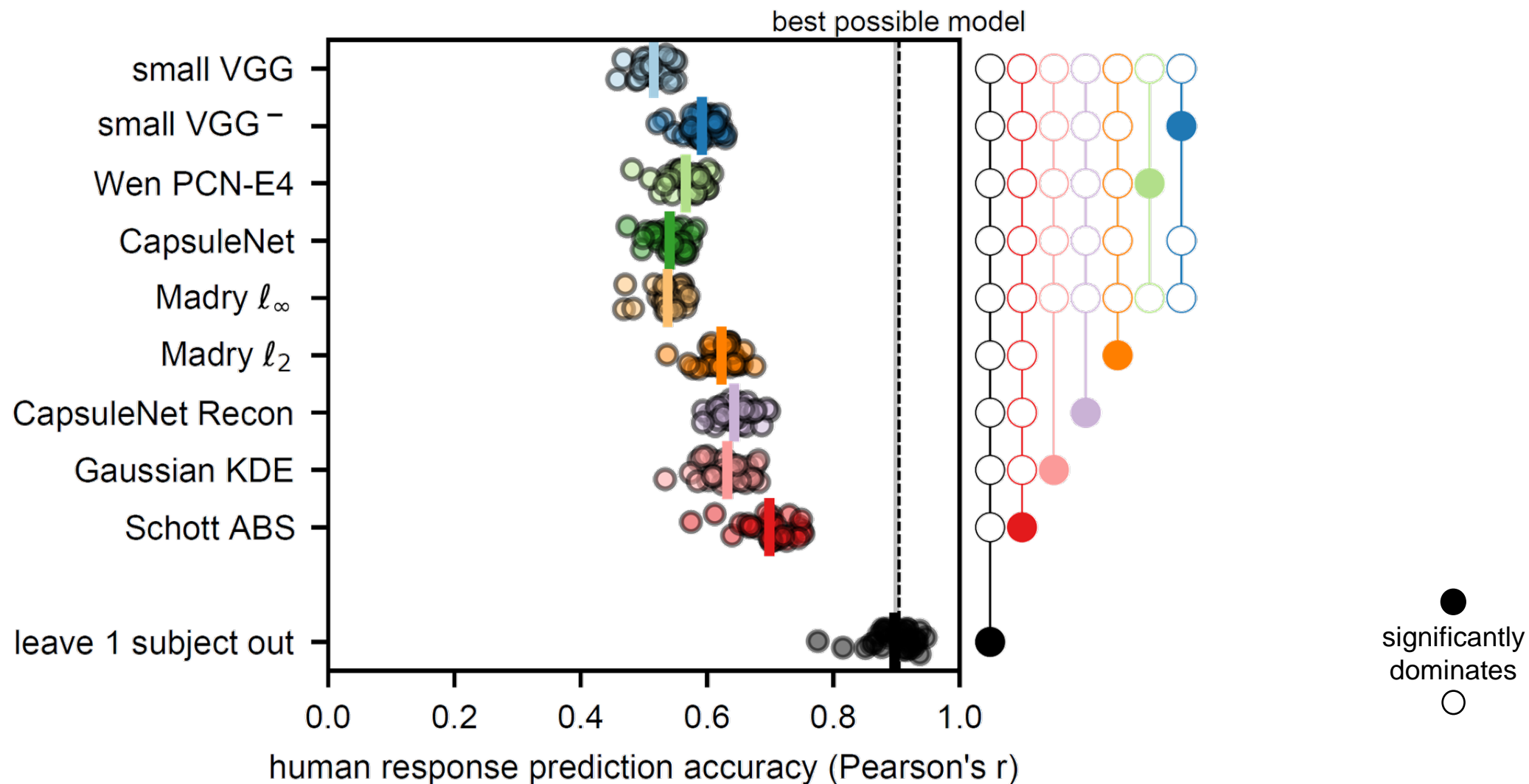
← Previous

Next →

Behavioral experiment

- 30 subjects (tested via *Prolific*)
- stimuli included 20 controversial stimuli per model pair
(36×20) + 100 MNIST images = 820 images per subject
- stimuli presented in a randomized order
- 820 stimuli x 10 scales x 30 subjects
(246,000 data points)

Controversial stimuli



Controversial stimuli

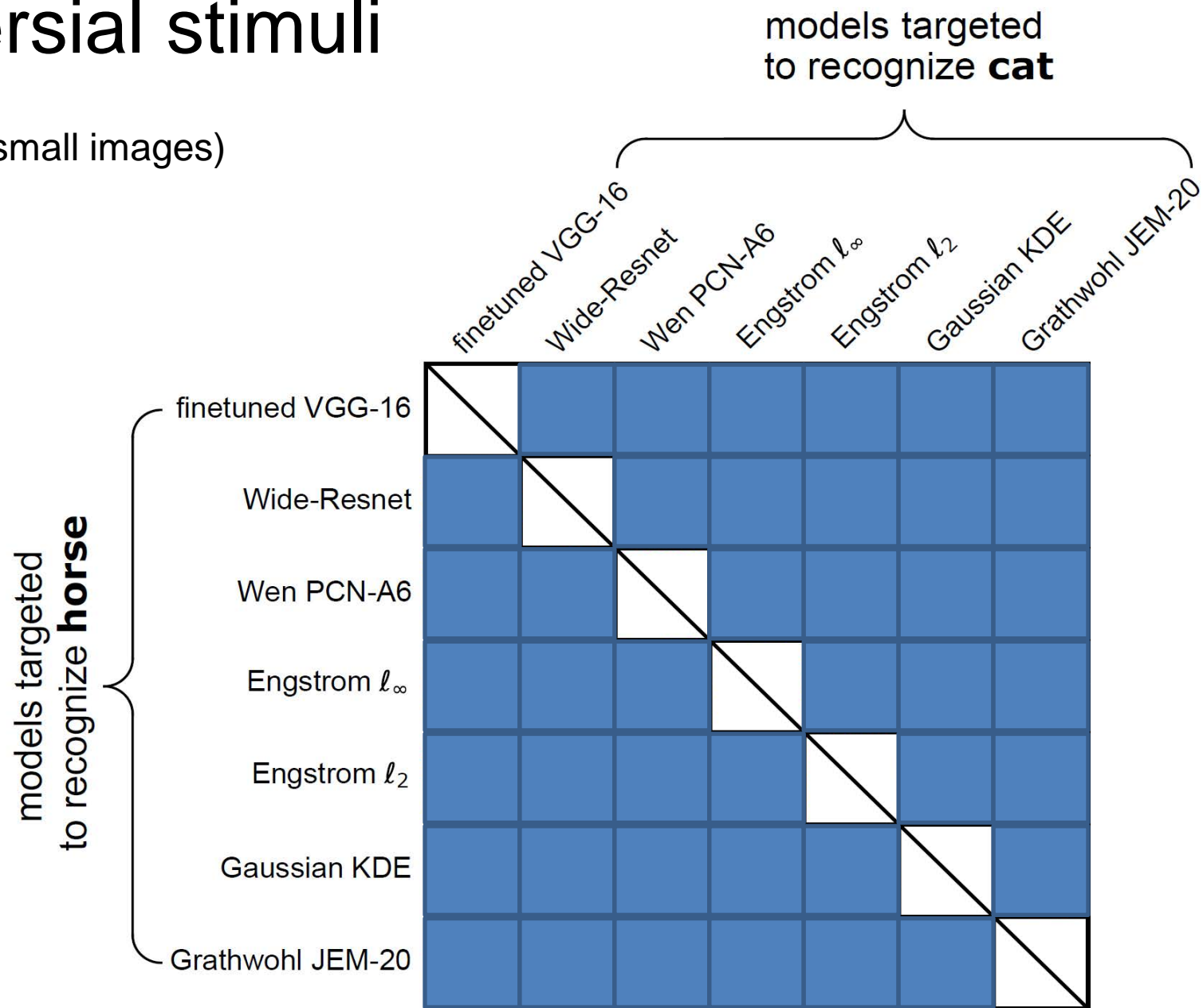
natural images

(CIFAR-10 set of small images)

Controversial stimuli

natural images

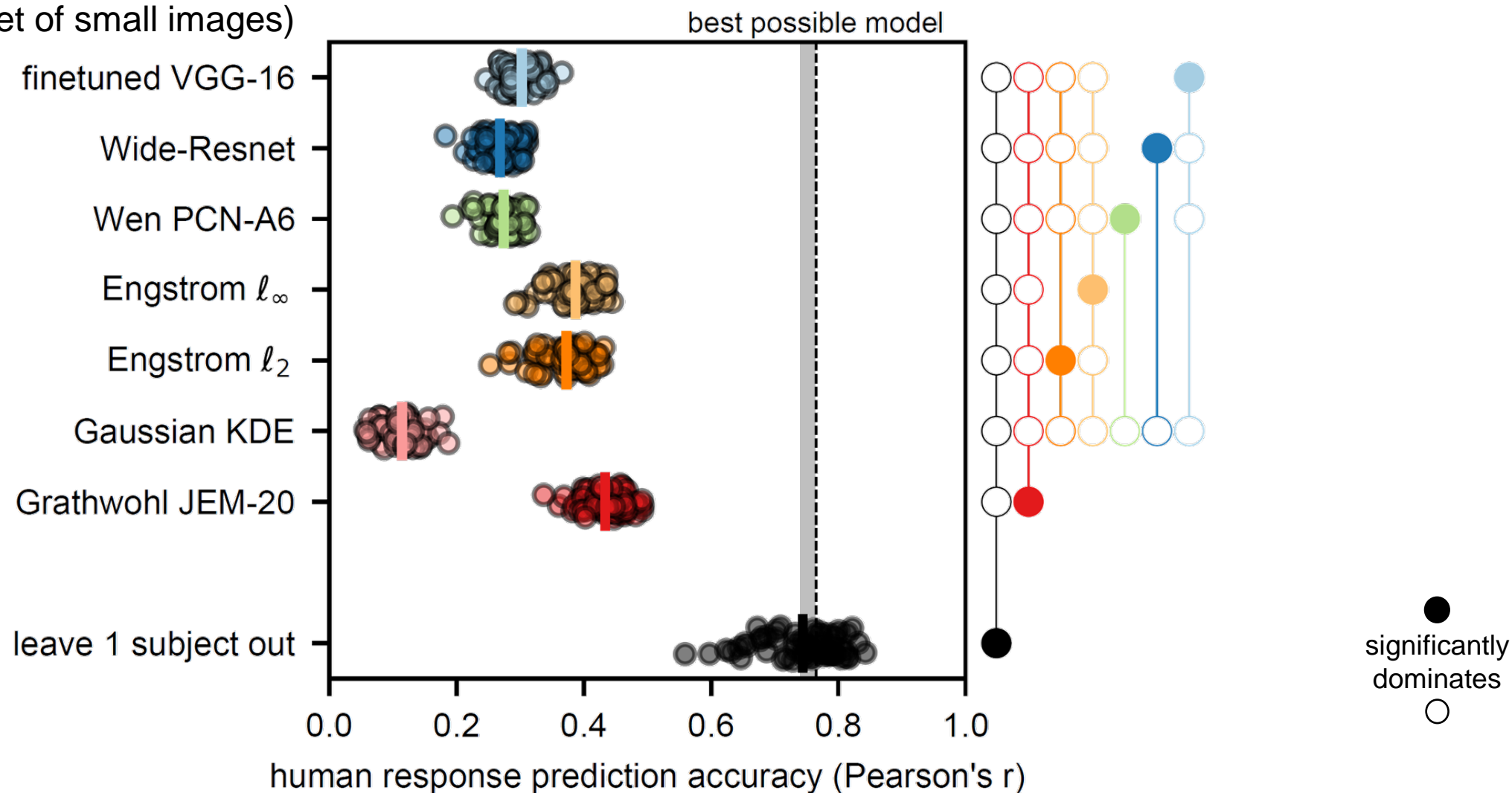
(CIFAR-10 set of small images)



Controversial stimuli

natural images

(CIFAR-10 set of small images)

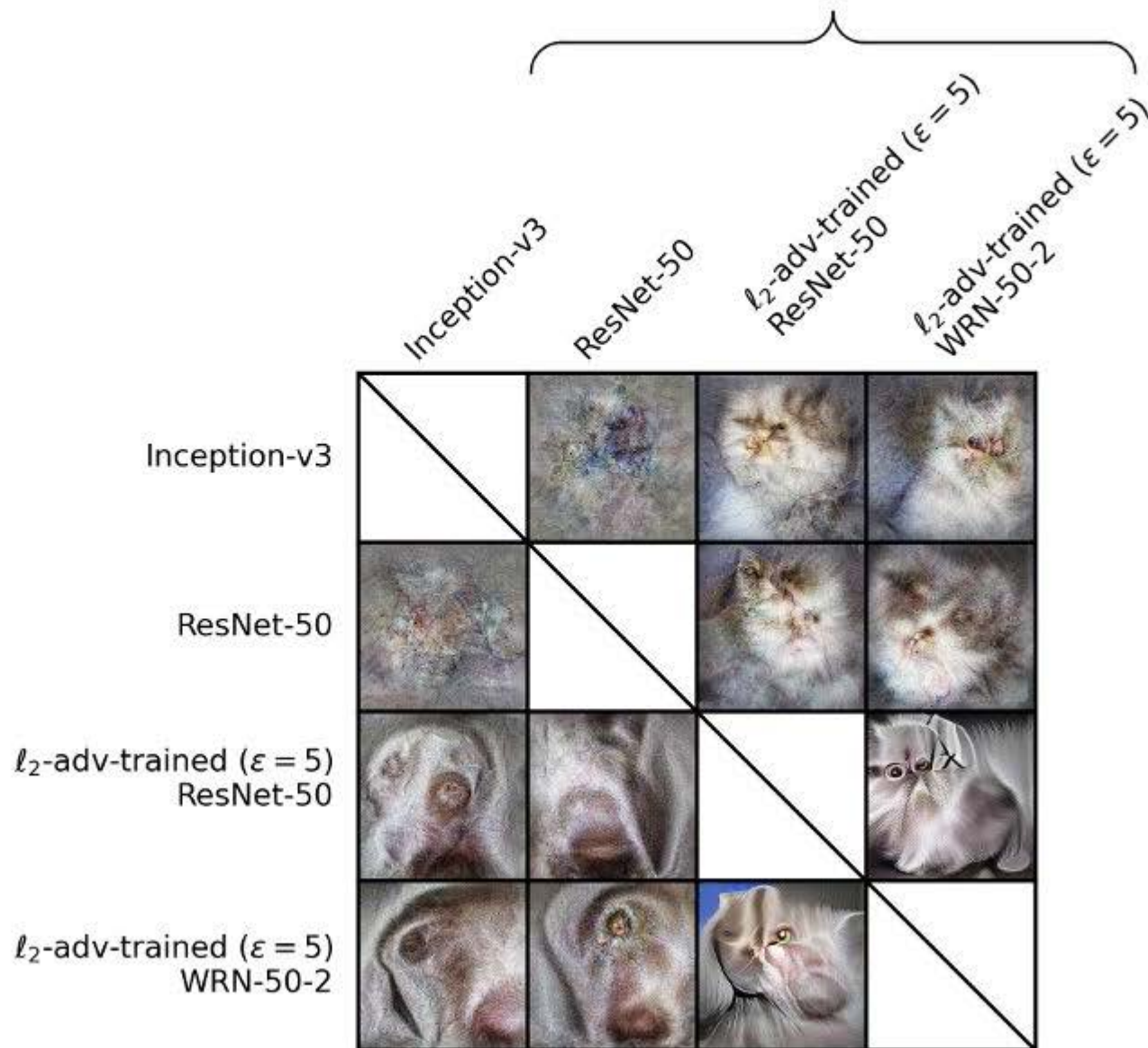


Controversial stimuli

natural images
(ImageNet)

models targeted
to recognize **Persian cat**

models targeted
to recognize **Weimaraner**



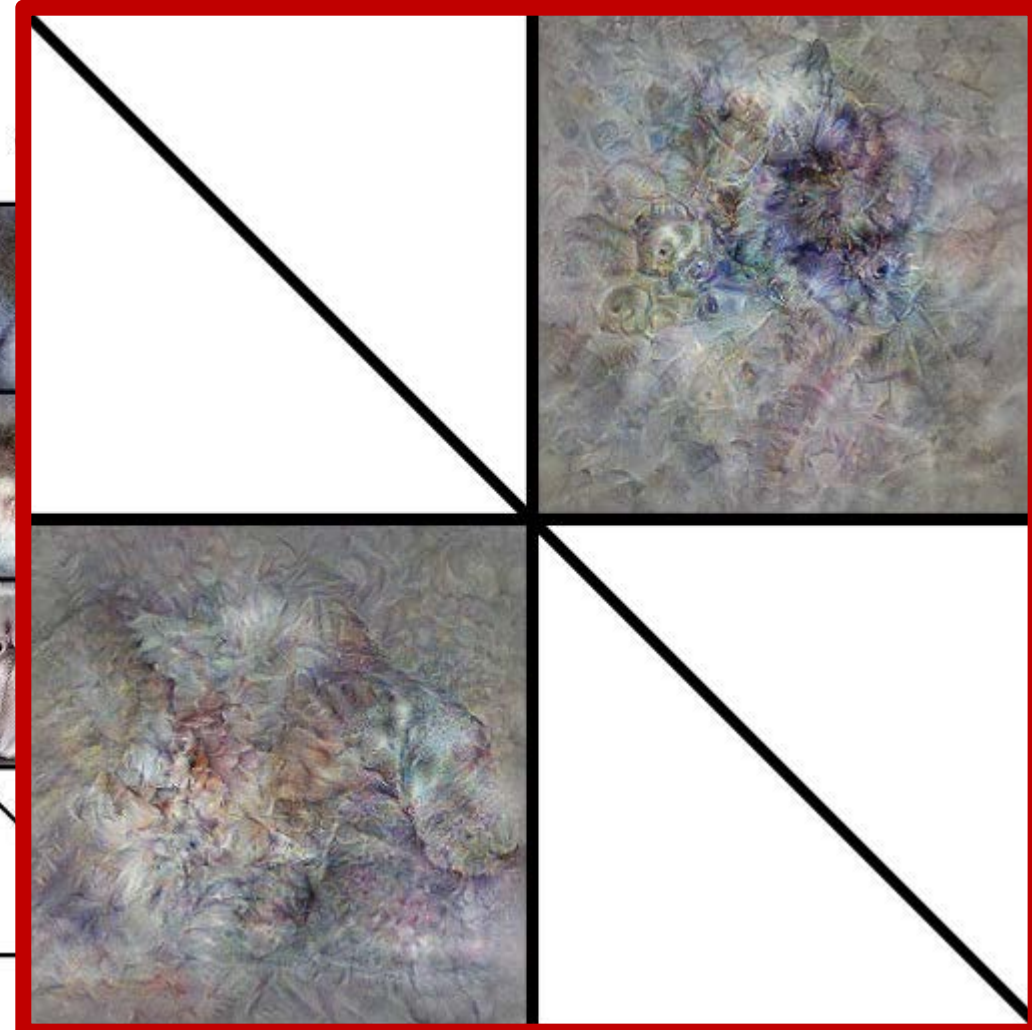
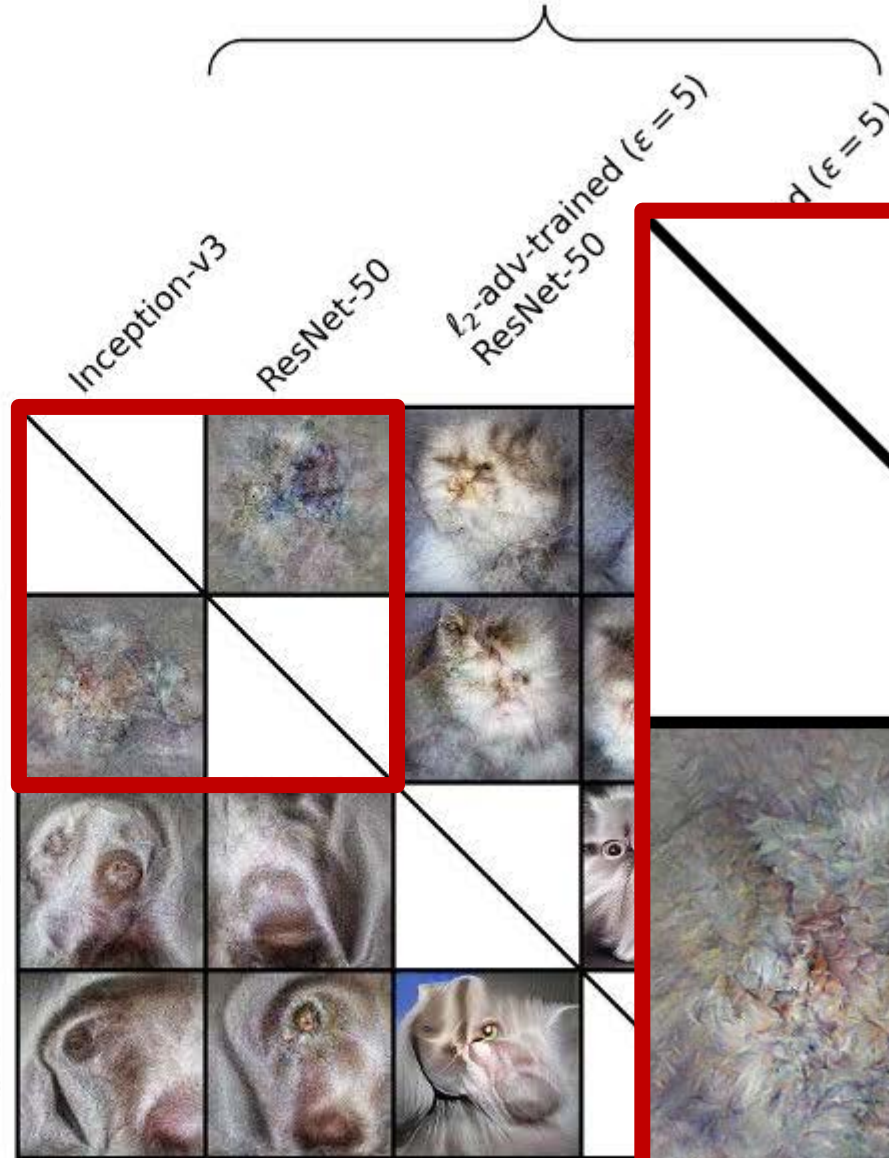
Controversial stimuli

natural images
(ImageNet)

models targeted
to recognize **Persian cat**

models targeted
to recognize **Weimaraner**

Inception-v3
ResNet-50
 l_2 -adv-trained ($\epsilon = 5$)
ResNet-50
 l_2 -adv-trained ($\epsilon = 5$)
WRN-50-2

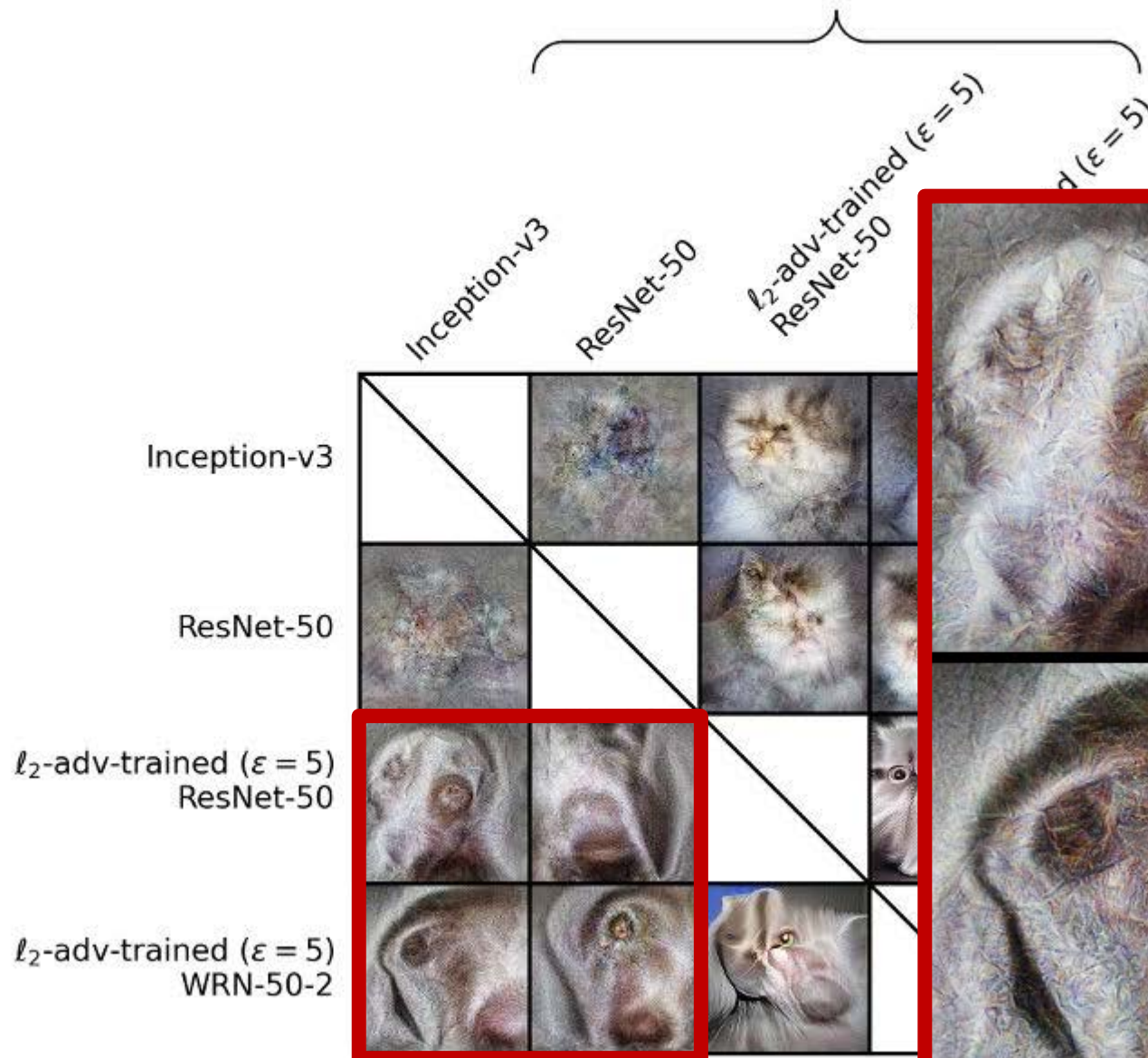


Controversial stimuli

natural images
(ImageNet)

models targeted
to recognize **Persian cat**

models targeted
to recognize **Weimaraner**

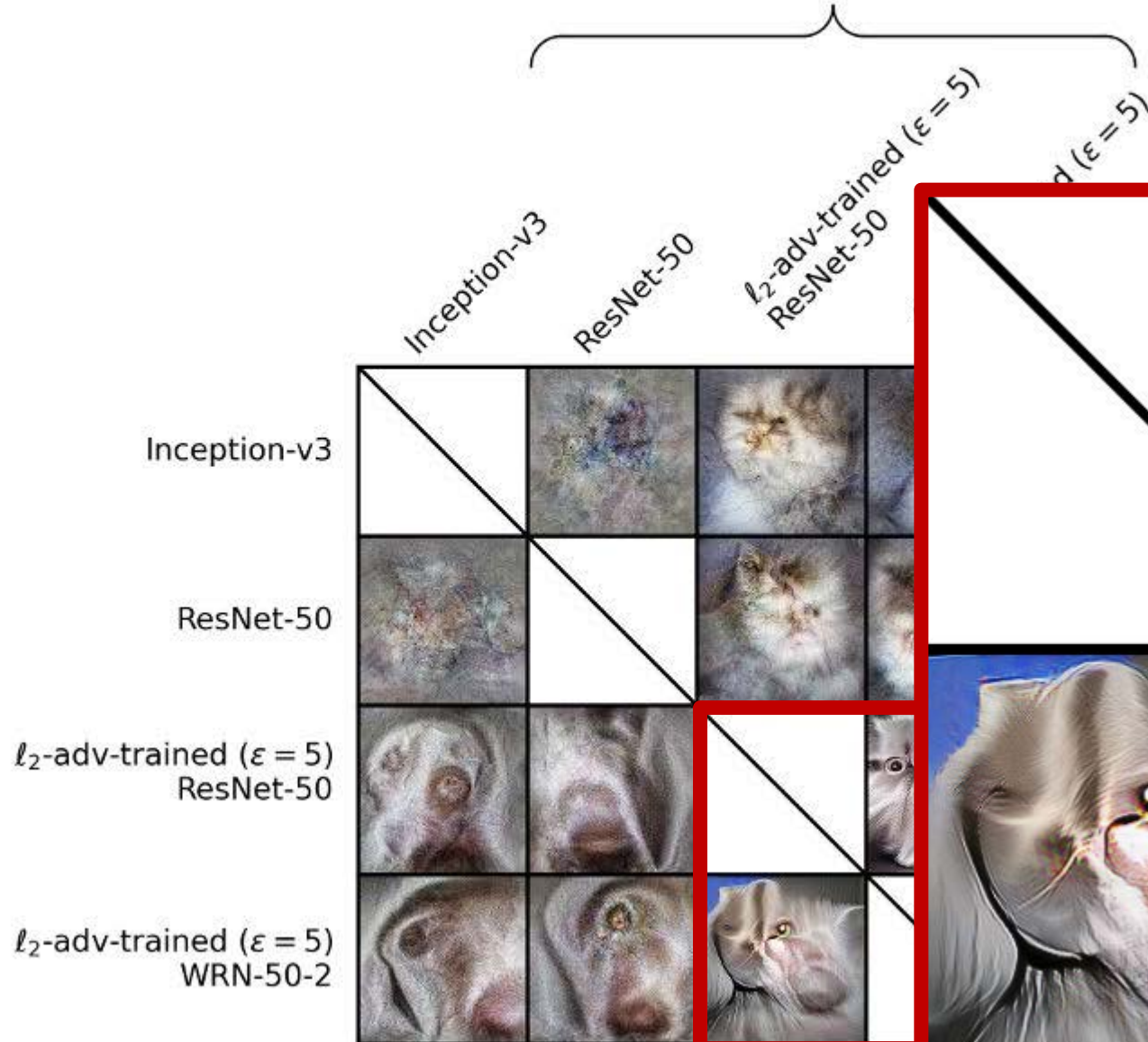


Controversial stimuli

natural images
(ImageNet)

models targeted
to recognize **Persian cat**

models targeted
to recognize **Weimaraner**



Overall conclusions

1. We can adjudicate among task-performing deep net models by inferentially comparing their representations to brain representations.
Nili et al. 2014, Kriegeskorte & Diedrichsen 2019
2. Recurrent convolutional vision models better predict human ventral stream representational dynamics and reaction times
Kietzmann et al. 2019, Spoerer et al. 2020
3. Controversial stimuli enable us to elicit differences in the inductive biases of deep net model.
Golan et al. 2020
4. Human vision may rely on a computational mechanism that combines elements of discriminative and generative inference.
Golan et al. 2020