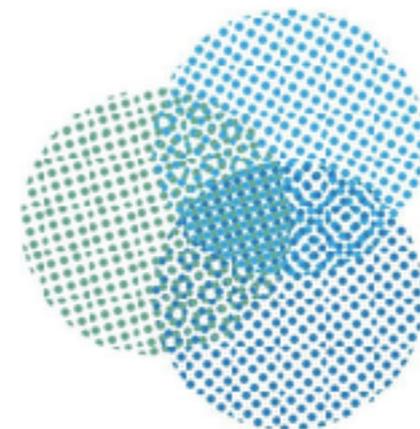


Making a Science from the Computer Vision Zoo

Xavier Boix

xboix@mit.edu



CENTER FOR
Brains
Minds+
Machines

Computer Vision Zoo

Brains

SLAM

Super Pixels

ResNet

HMAX

Primates

Human

CRFs

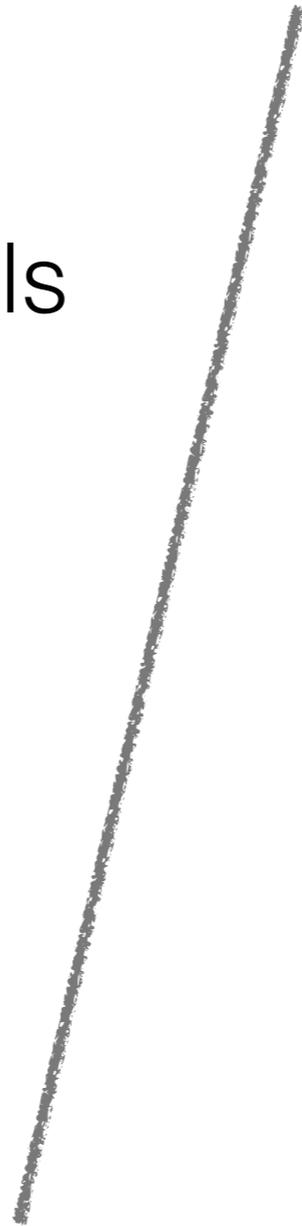
VGG

Zebra Fish

AlexNet

Fly

Rat



The Computer Vision Zoo

Engineering of algorithms

New algorithms = Better Performance

(“Darwinian evolution”)

Evolution of the Computer Vision Zoo

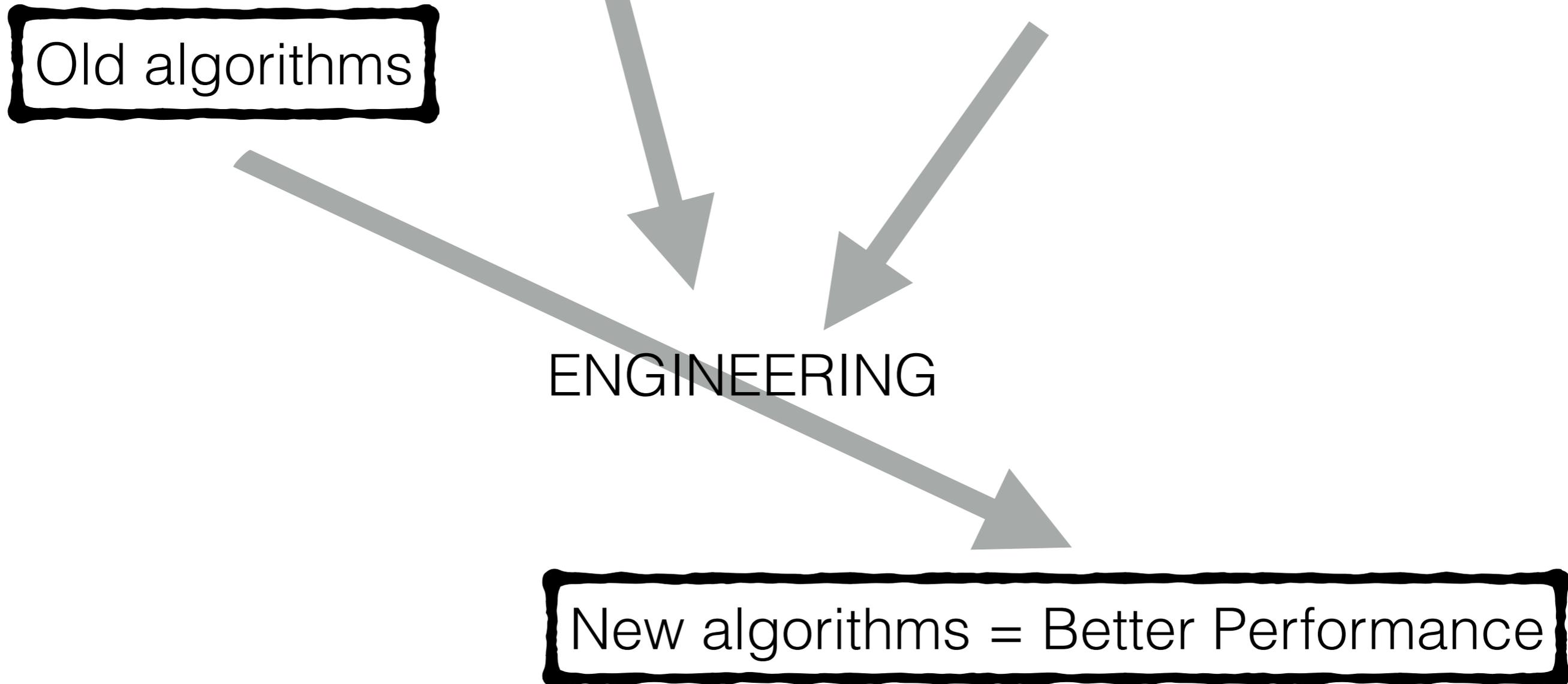
Machine Learning

Brain &
Cognitive Sciences

Old algorithms

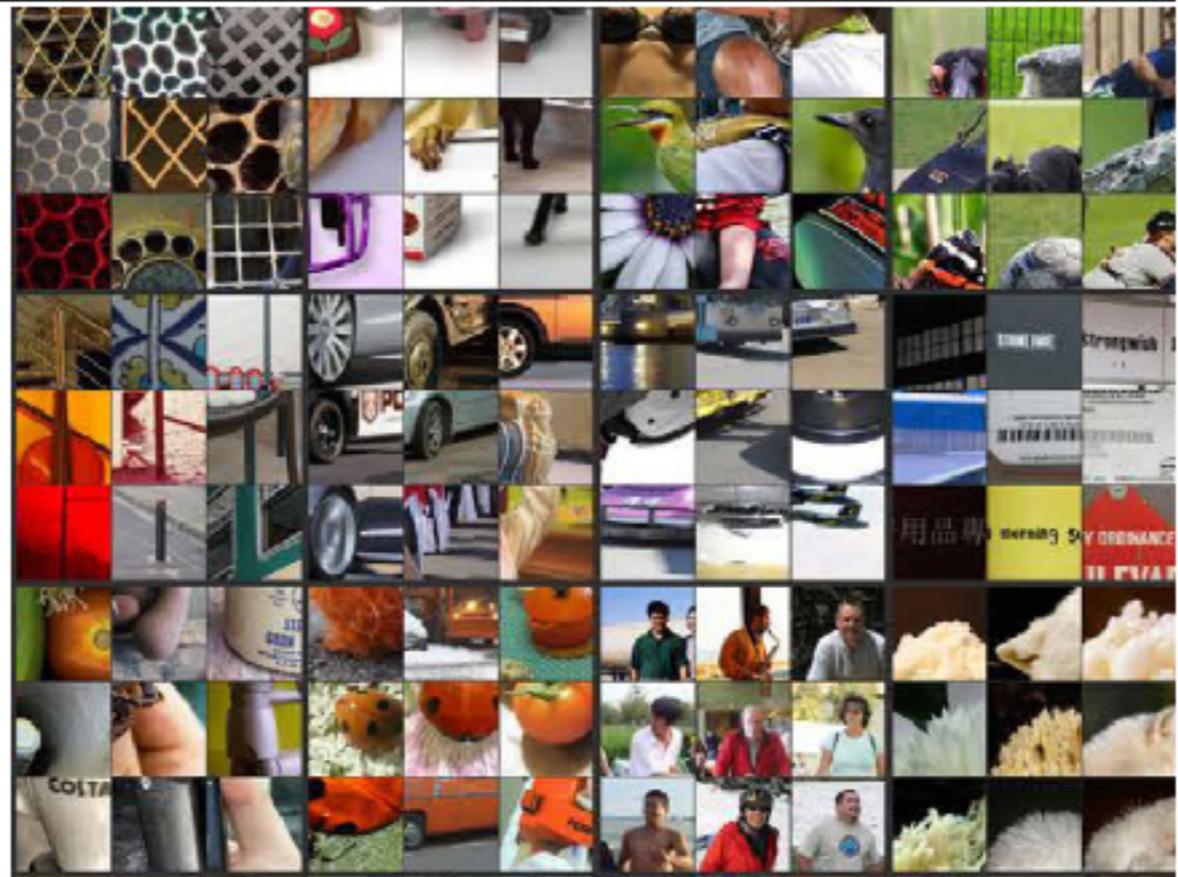
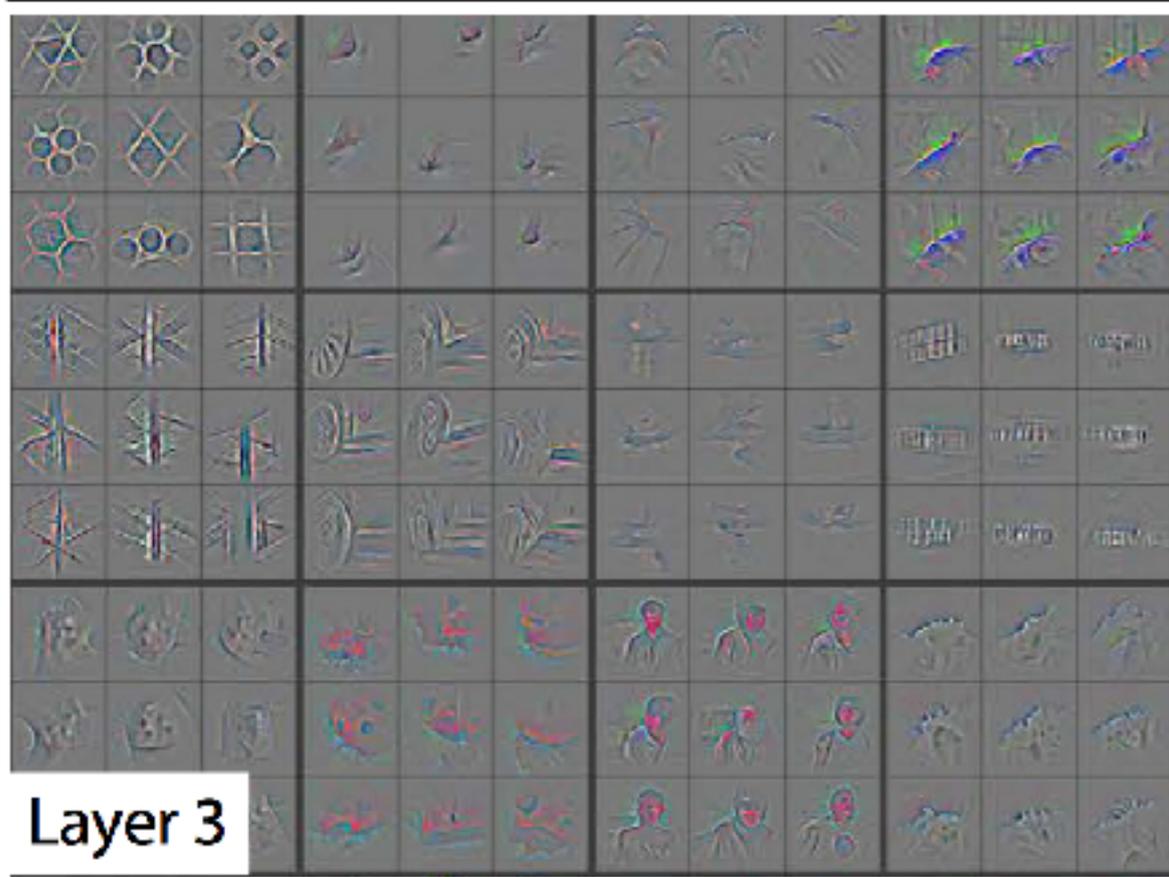
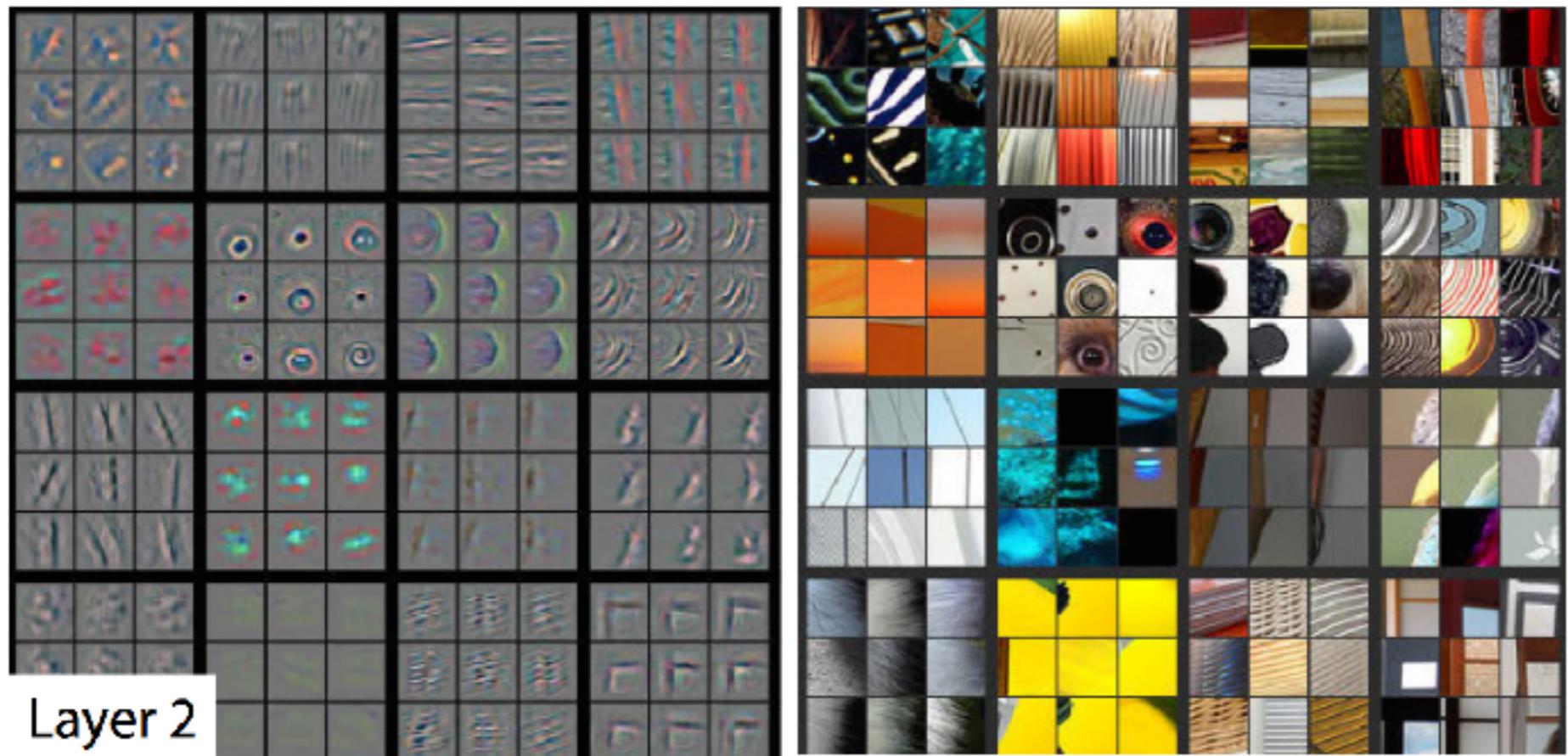
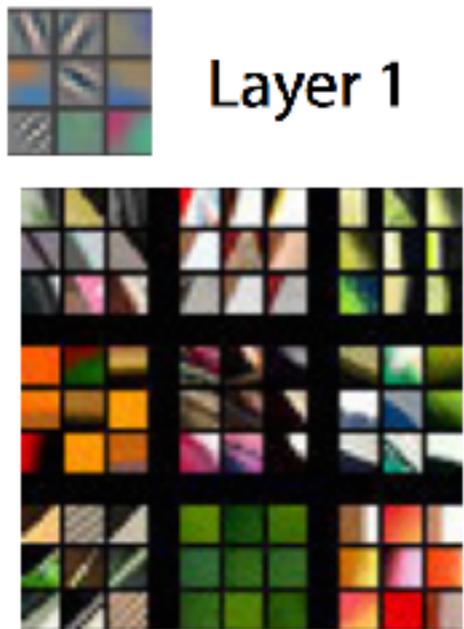
ENGINEERING

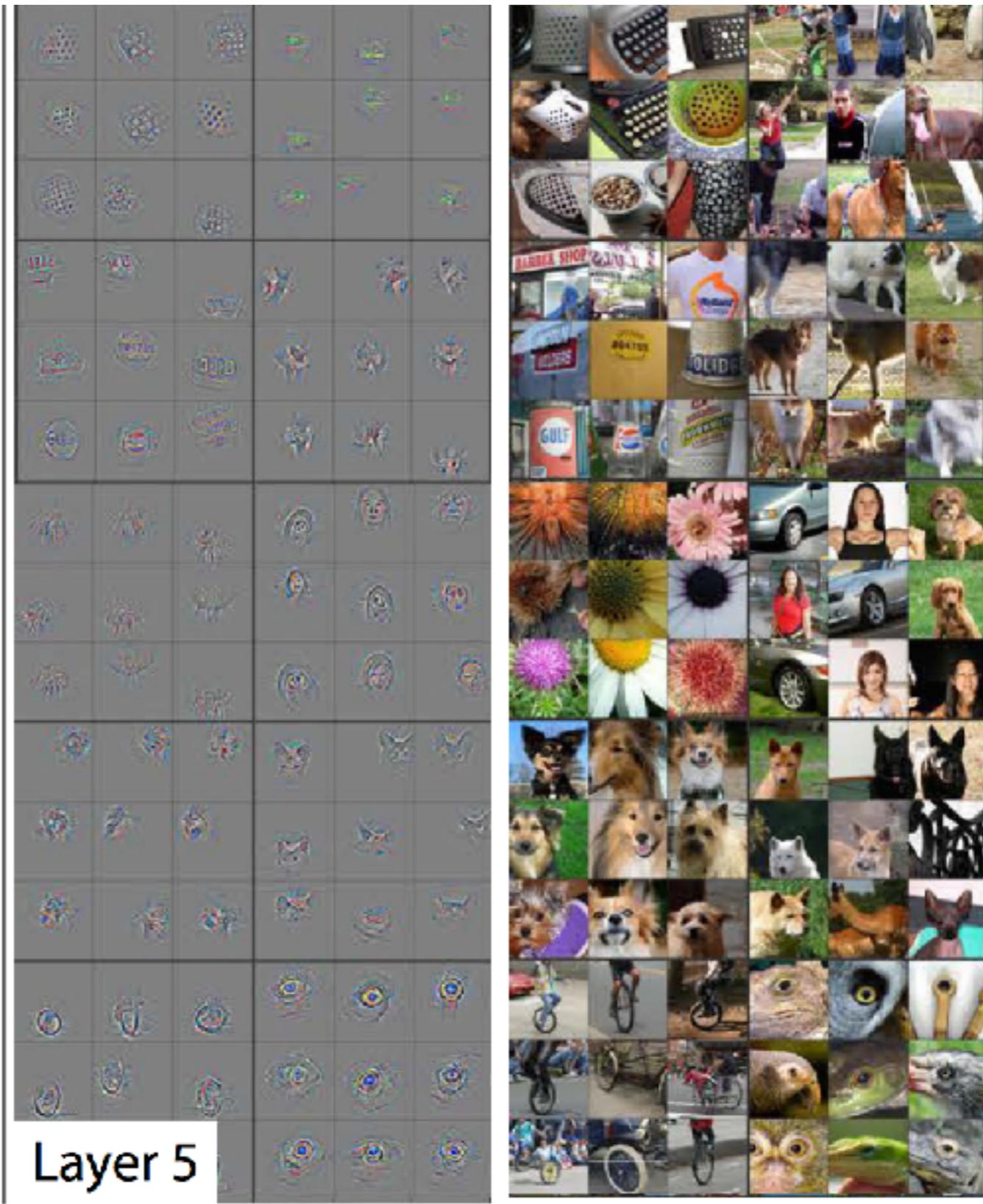
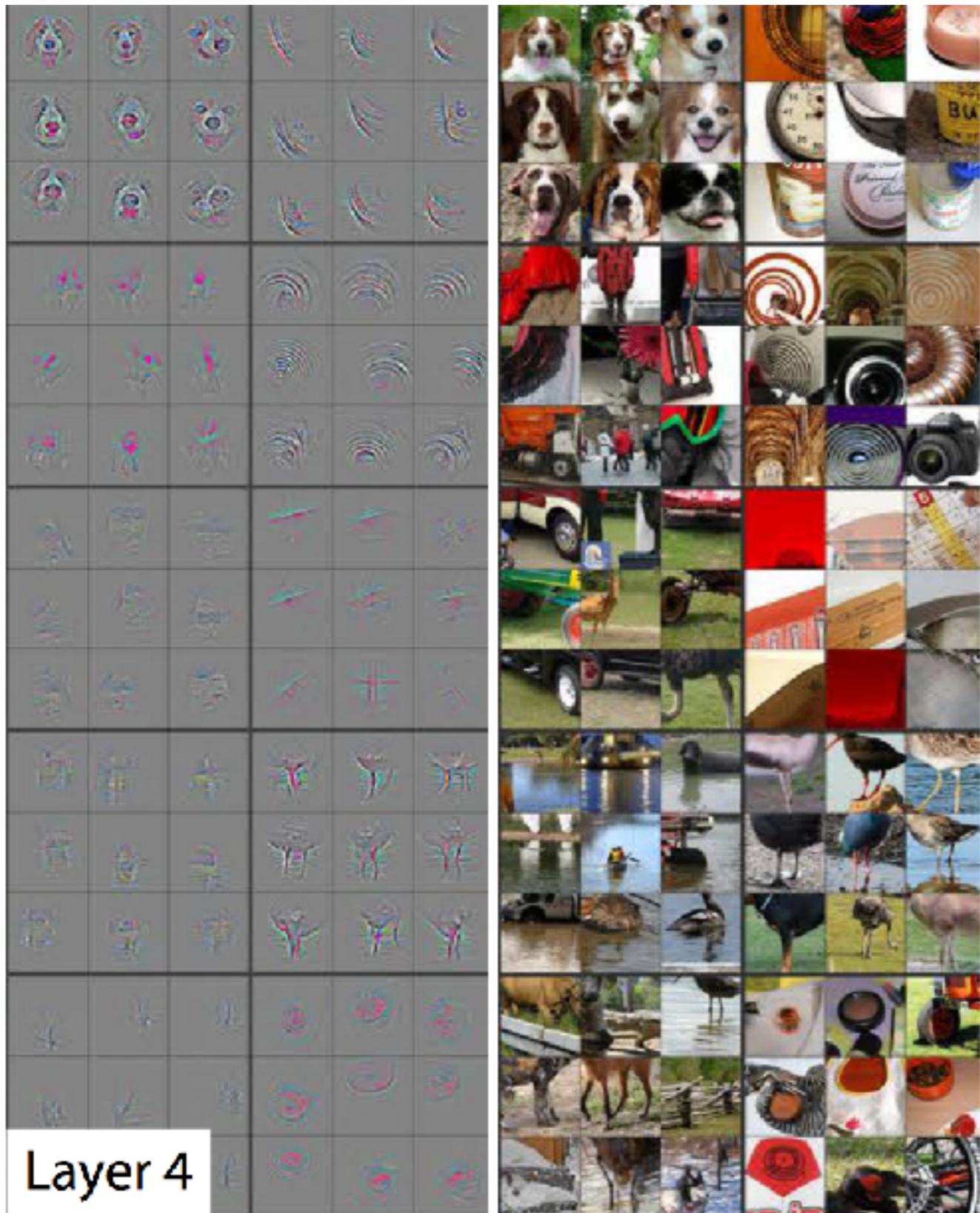
New algorithms = Better Performance

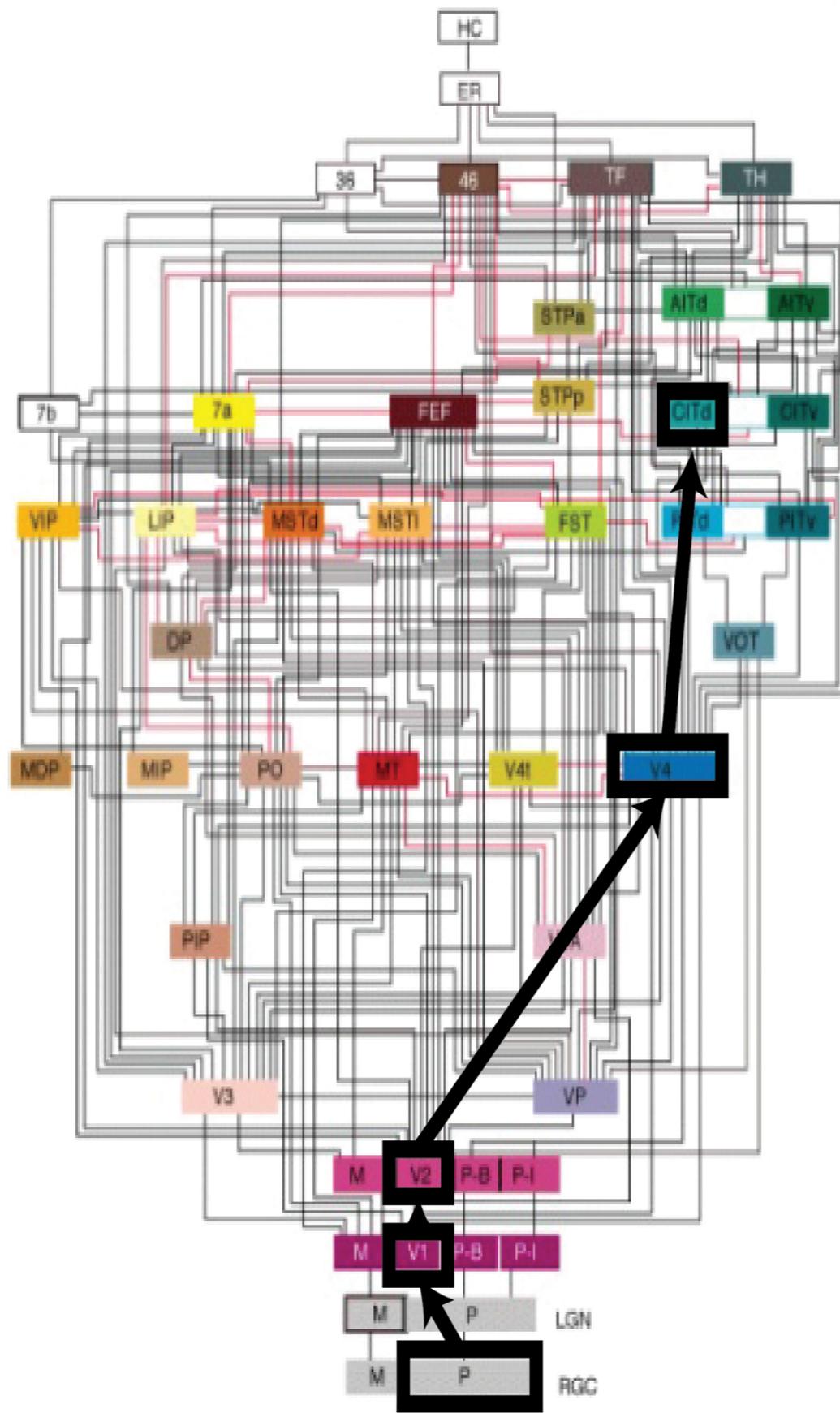


**Are we going to understand intelligence by building
AIs that pass the Turing test**

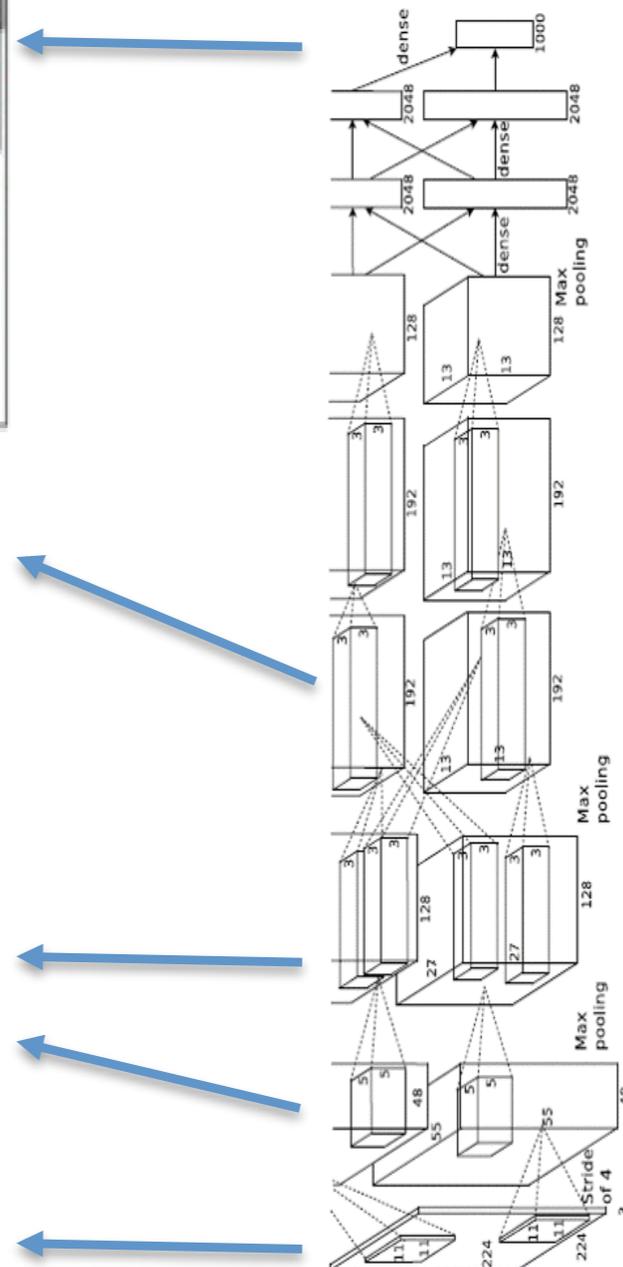
?







Felleman and Van Essen 1991



Krizhevsky et al 2012

Marr Levels of Understanding

1) Computational

Computer Vision
(Build better algorithms)

2) Algorithmic

3) Implementation

Science of Computer Vision

Make a science from the Zoo

computational principles among algorithms

- Develop testable hypothesis



- Test and challenge the hypothesis

At the Computational Level

Develop scientific hypothesis...

- Meta-algorithms, unification of algorithms

- Visualization analysis

- Mathematical theories

etc.

...to predict properties of the algorithms before testing them

Computer Vision Zoo

Brains

SLAM

Super Pixels

Human

Primates

ResNet

HMAX

Models

Zebra Fish

CRFs

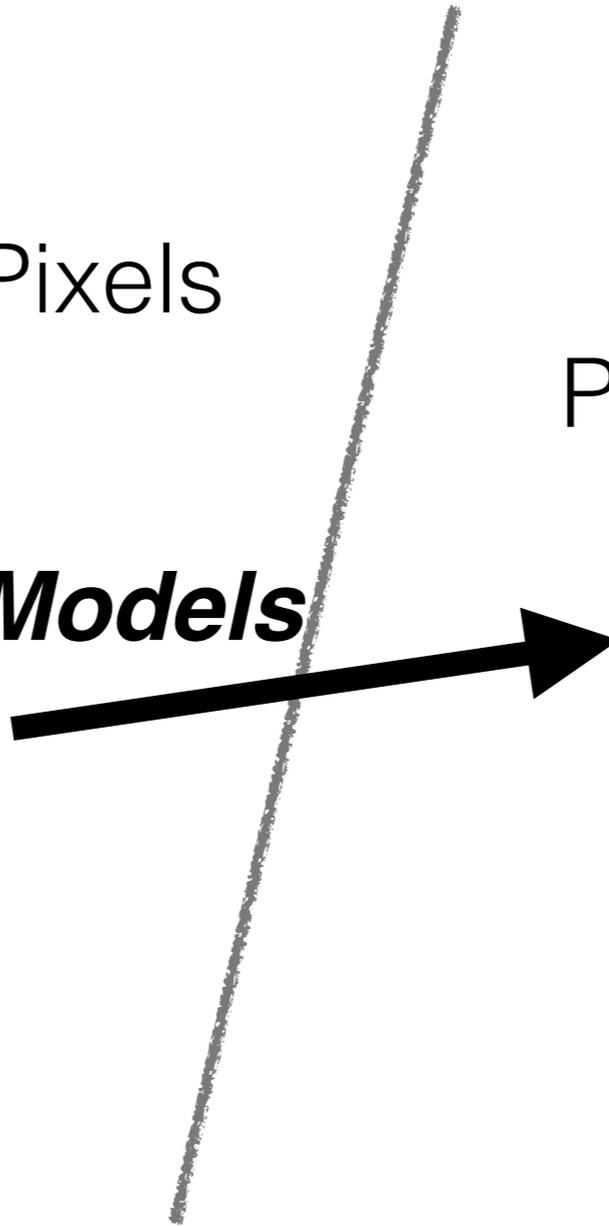
VGG

Fly

Rat

AlexNet

Computational hypothesis



1) Computational

Same principles should explain
artificial and biological intelligence

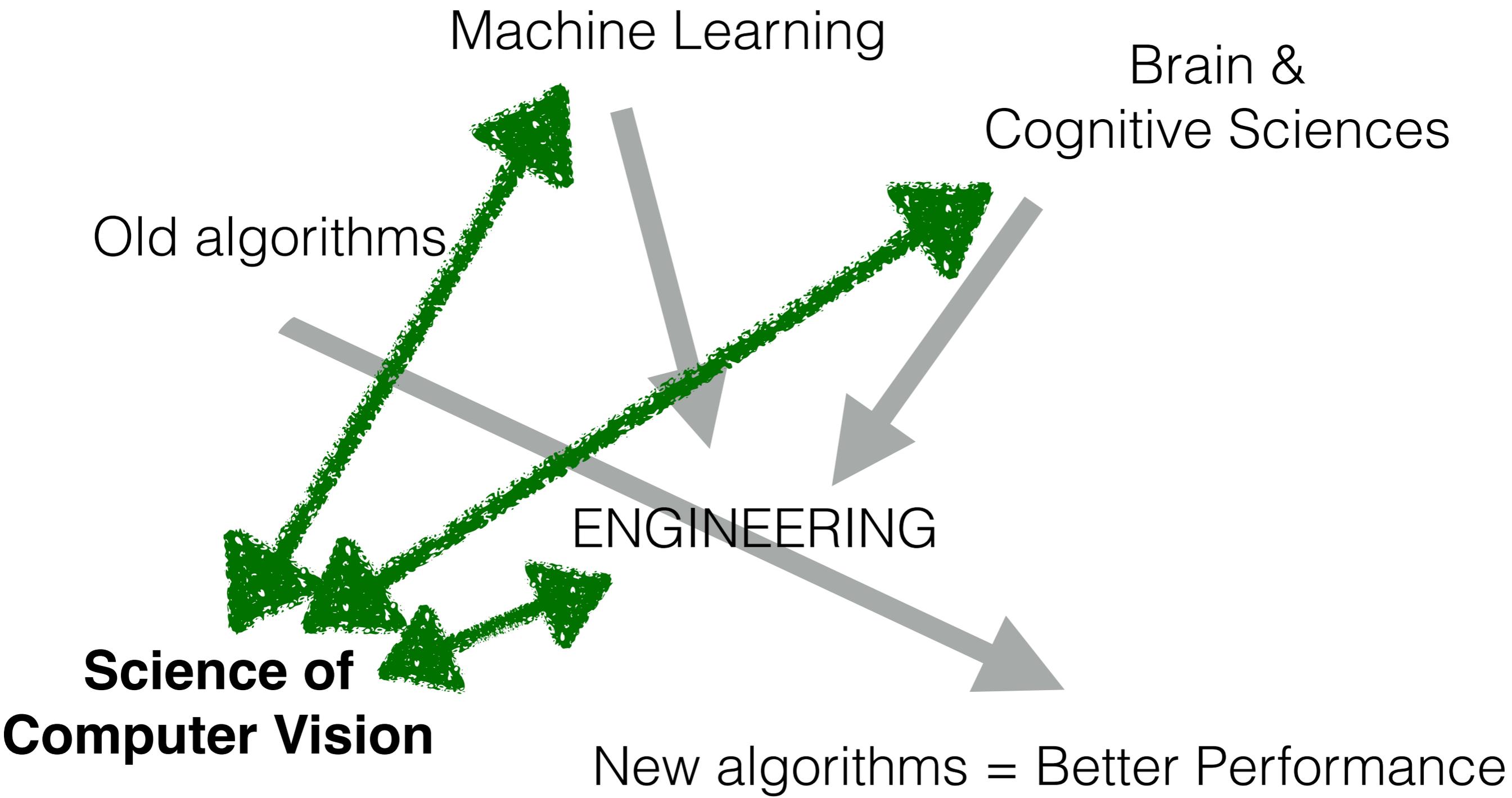
(eg. Why the brain is how it is?
Why an algorithm works better than others?)

2) Algorithmic

Computer Vision
(Build better algorithms)

Modelling the Brain
(How does the brain work?)

3) Implementation



3 Snapshots

1) Failures of DNNs

2) Beyond Object Recognition

3) DNNs to explain the Brain

Minimal Images in Deep Neural Networks

with

Sanjana Srivastava and Guy Ben-Yosef



Adversarial examples



Original

Adversarial

Adversarial
(exaggerated to show
perturbation)



(Ullman et al., 2016)



(Ullman et al., 2016)



(Ullman et al., 2016)



(Ullman et al., 2016)

Minimal images in humans

- Smallest region of an image that is still recognizable to humans (Ullman et al., 2016)
- DNNs unable to recognise human minimal images (Ullman et al., 2016; Ben-Yosef et al., 2018)

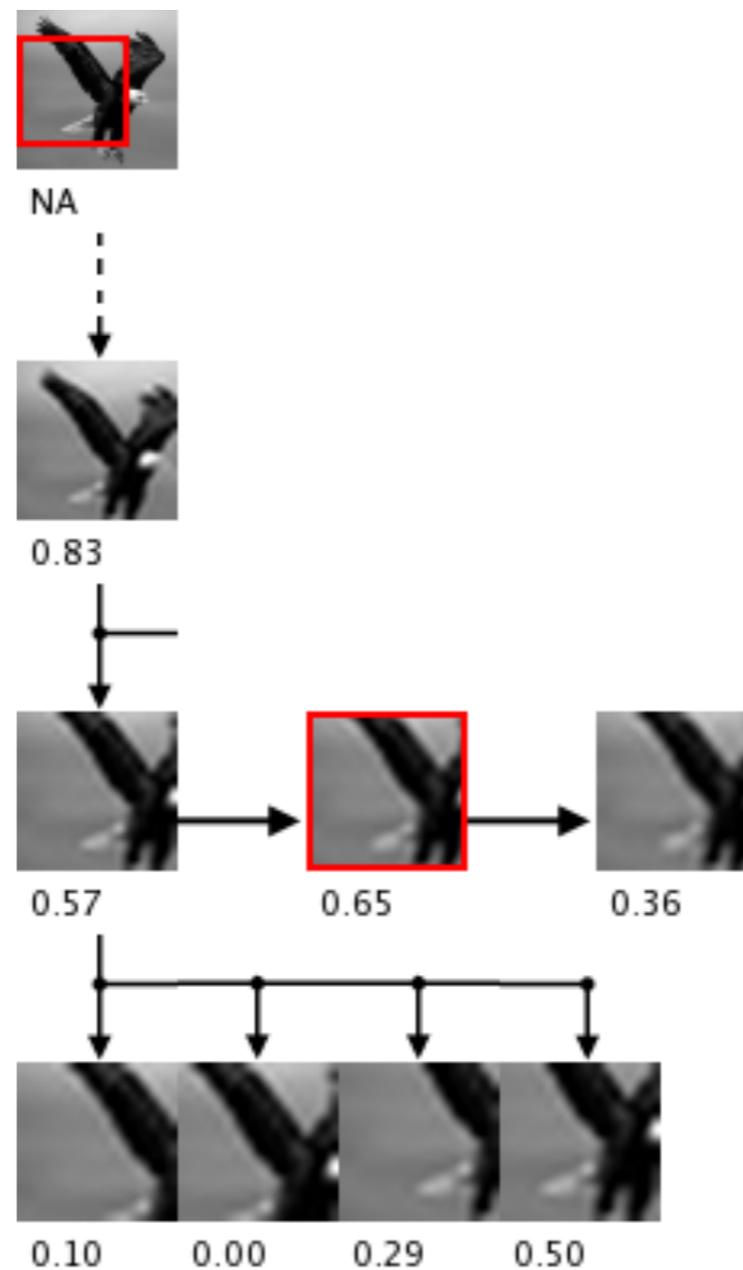
Unrecognizable to humans



Recognizable to humans

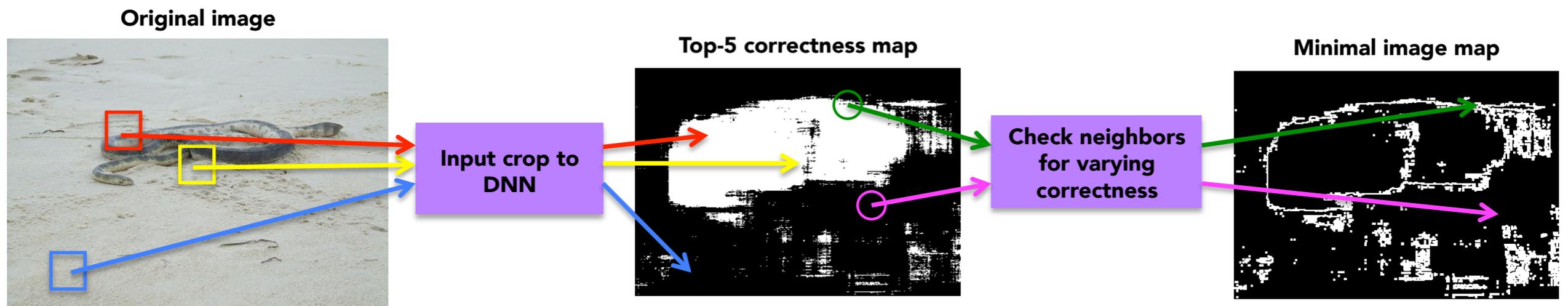


Extracting Human Minimal Images

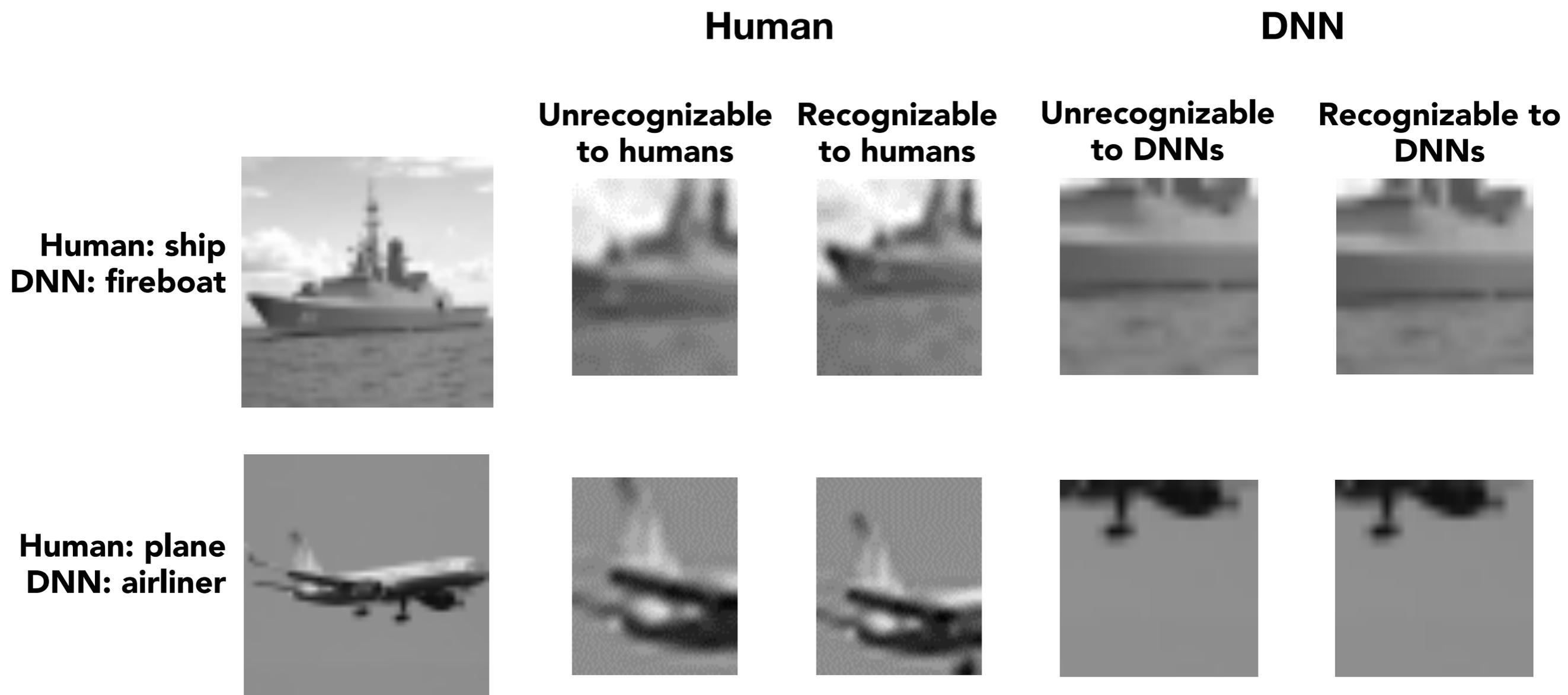


(Ullman et al., 2016)

Extracting DNNs Minimal Images



Do DNNs have their own set of minimal images?



Fragile Recognition Image

FRI: *Image region that is correctly categorised but a small shift of the region location or shrink of the region scale produces misclassification.*

Minimal Images \subset FRIs

Different ways to evaluate fragile object recognition

Original image



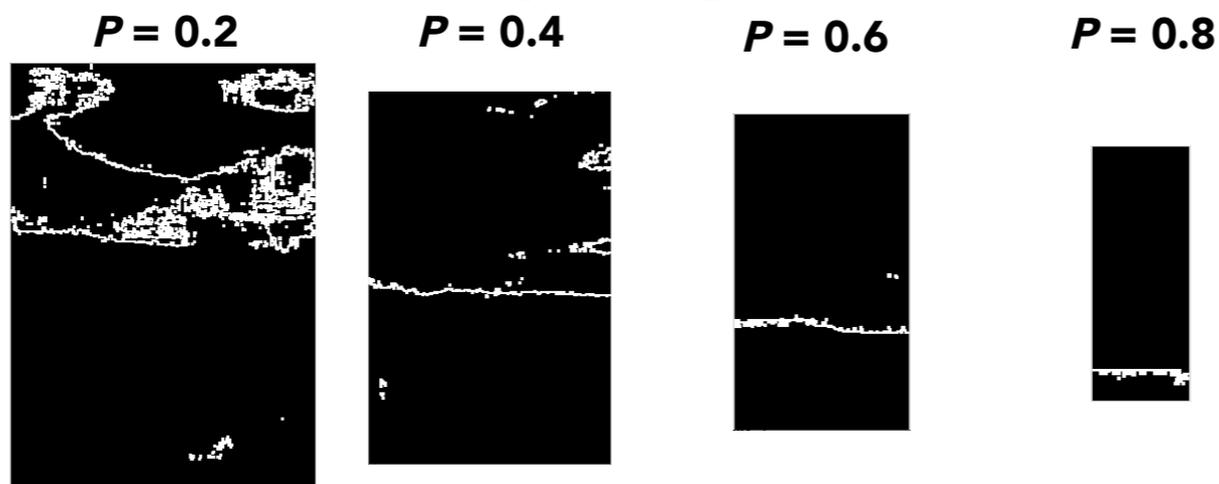
Label: coffee



Loose vs. strict

Scale vs. shift

$P = \text{crop side length} / \text{image lesser dimension}$



FRI occurs at boundaries of useful regions

Loose FRI: there exists a slight change that will cause failure

Strict FRI: any slight change will cause a failure

Experimental Setup

VGG, ResNet, Inception: time to process 1 image in 8 K80 GPU is ~5 minutes

500 images randomly chosen from the validation set

DNNs fail on smaller changes



ImageNet class: Electric guitar
Human option: Instrument



ResNet score: ~0.8
Human success rate: 65%



ResNet score: ~0.4
Human success rate: 55%



ImageNet class: Broccoli
Human option: Vegetable



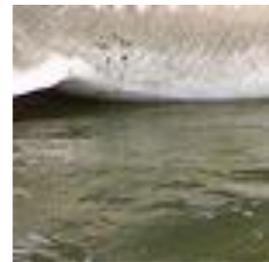
ResNet score: ~0.9
Human success rate: 60%



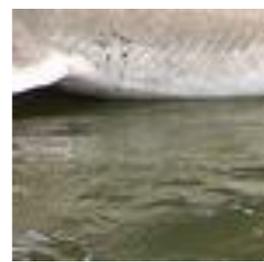
ResNet score: ~0.5
Human success rate: 80%



ImageNet class: Sturgeon
Human option: Fish



ResNet score: ~0.7
Human success rate: 5%



ResNet score: ~0.3
Human success rate: 15%



ImageNet class: Artichoke
Human option: Vegetable



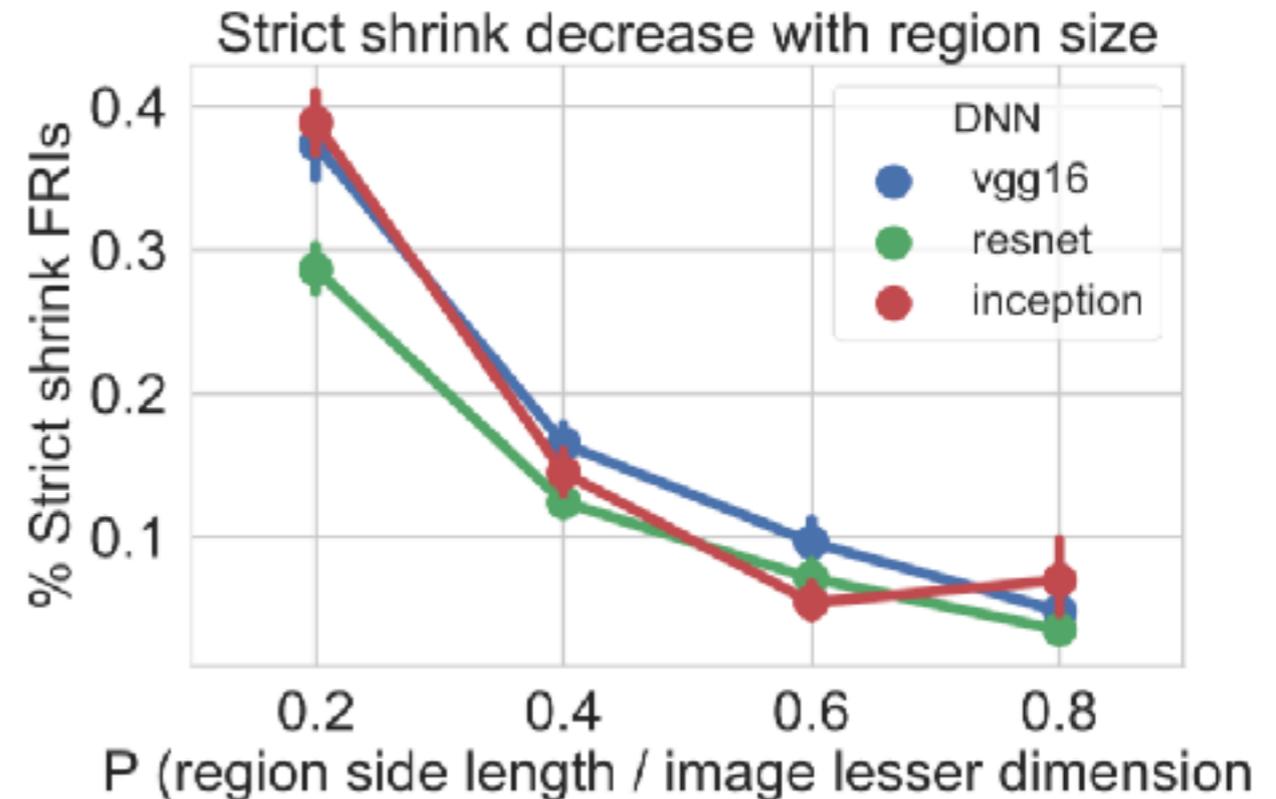
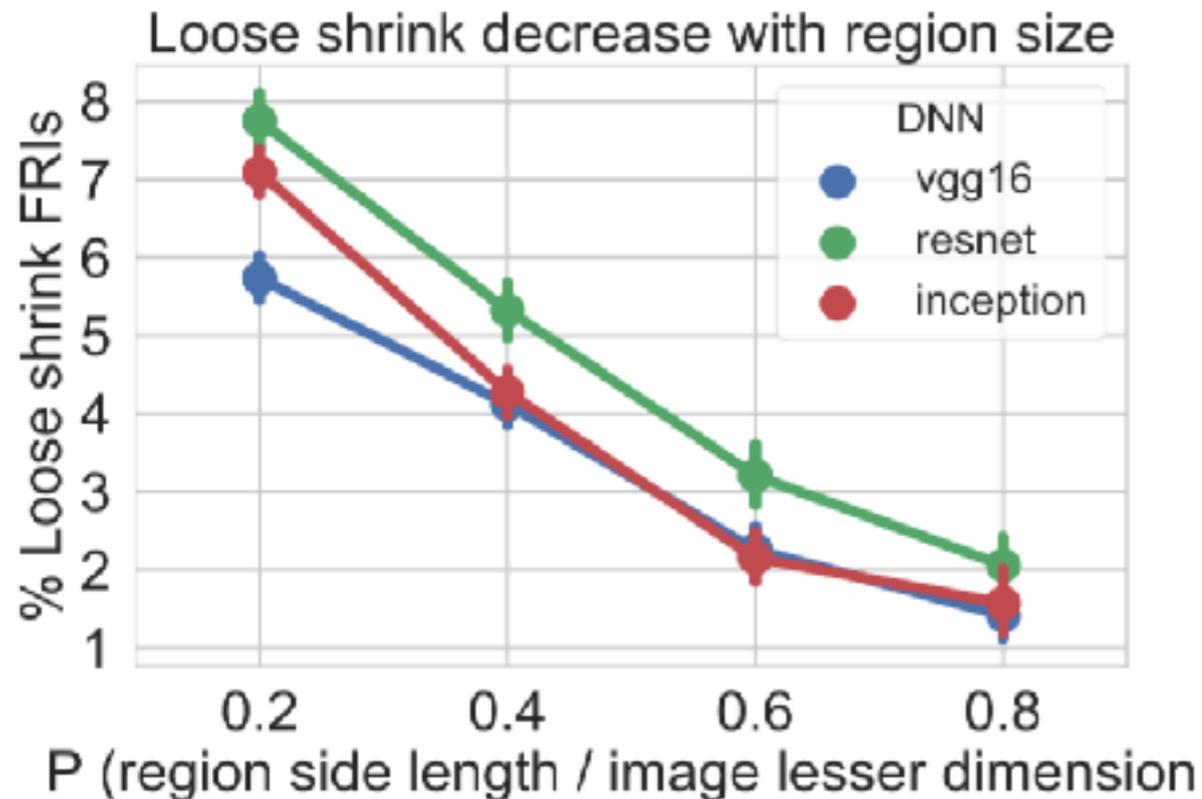
ResNet score: ~0.8
Human success rate: 65%



ResNet score: ~0.4
Human success rate: 75%

Difference from humans (1): DNNs can be affected by a single-pixel shift or shrink of the visible image region

DNNs have large FRIs



Difference with humans (2): DNNs are generally affected by visible image regions that are larger than human minimal images

Is this a new type of adversarial example without synthetic patterns?

In synthetic images:

Luo et al. 2016



Original

Adversarial

Adversarial
(exaggerated to show
perturbation)

Targeted synthetic changes that
cause DNN failure

Generally imperceptible to humans

In **natural images**:



Various translations of a single image

Translations can cause DNNs to fail (Engstrom et al. 2016) (Azulay, Weiss 2018)

Location Invariance and Minimal Images

Experiment:

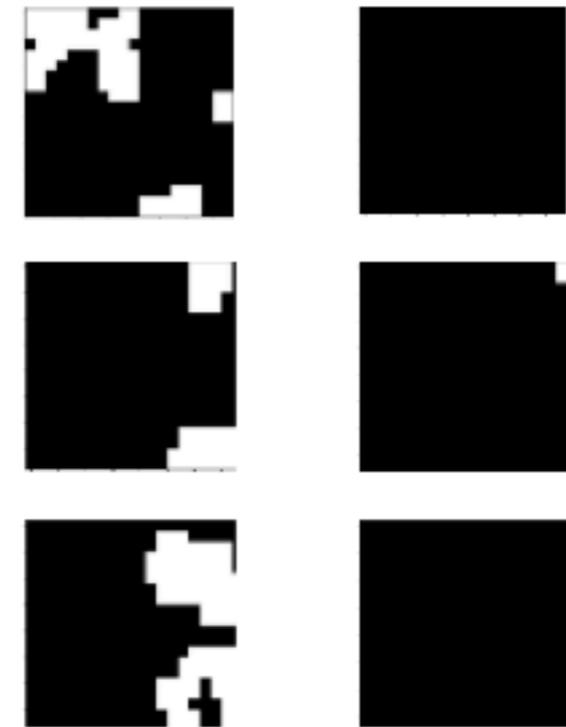
- CIFAR10
- Network: 2 convolutions + 2 fully connected (standard, ~75% accuracy)
- We vary the *pooling region size*:
 - Larger pooling sizes lead to more position invariance.
 - Pooling size = 32: entire image pooled (maximum invariance, similar accuracy ~72%)



Pooling improves location invariance

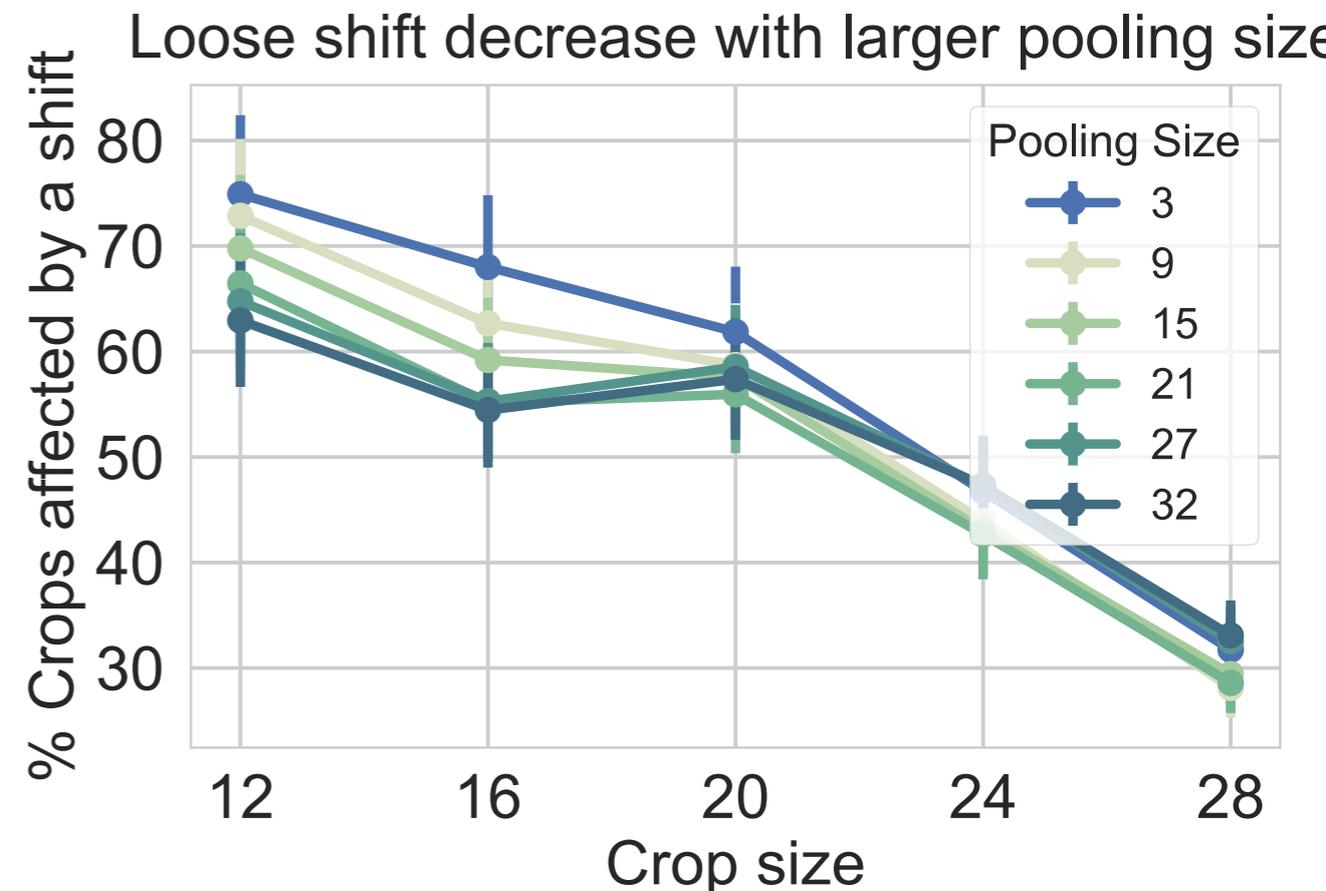
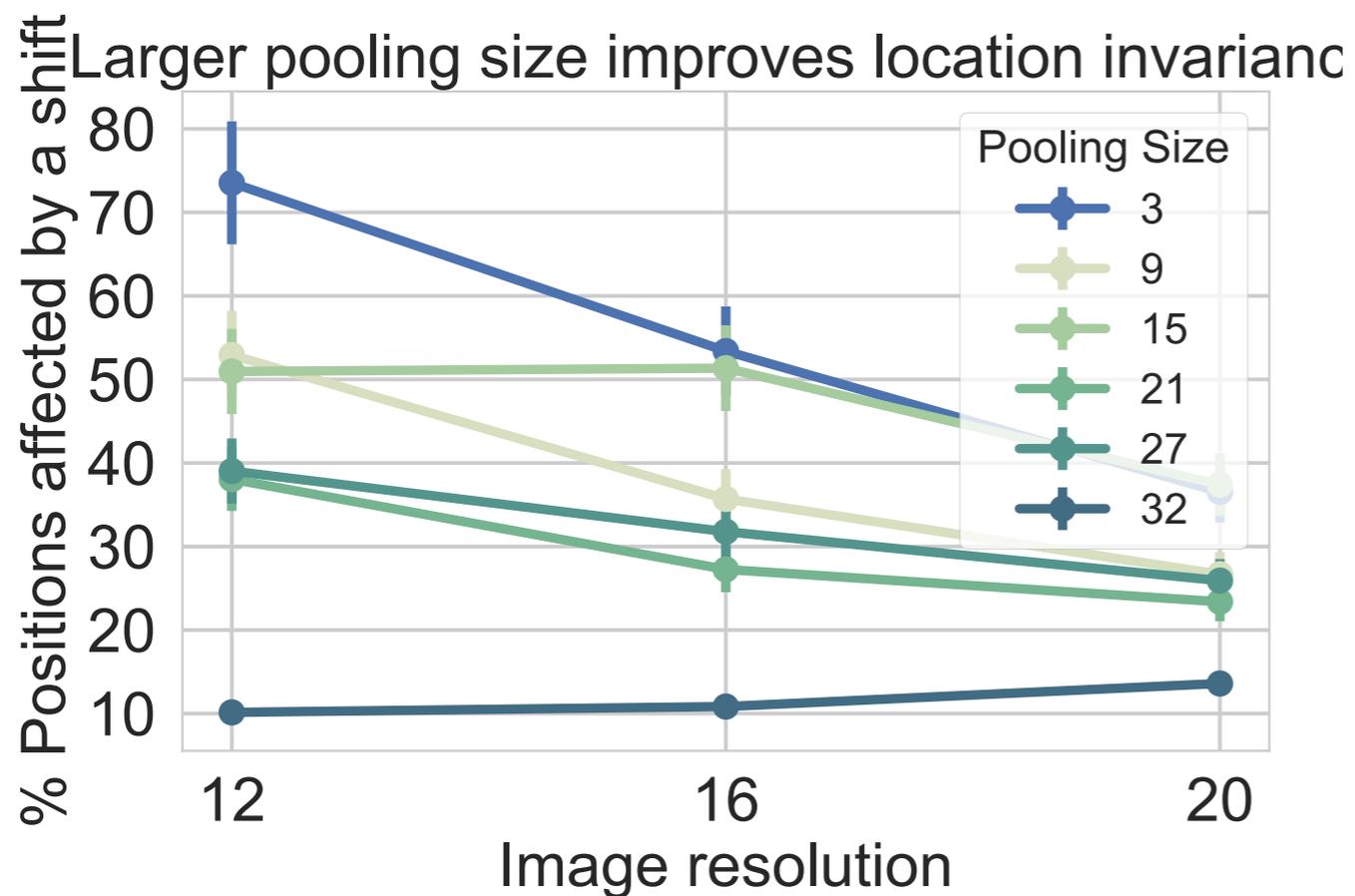
Pooling size = 3

Pooling size = 32



White pixels indicate positions for which a shift of 1 pixel will change network accuracy

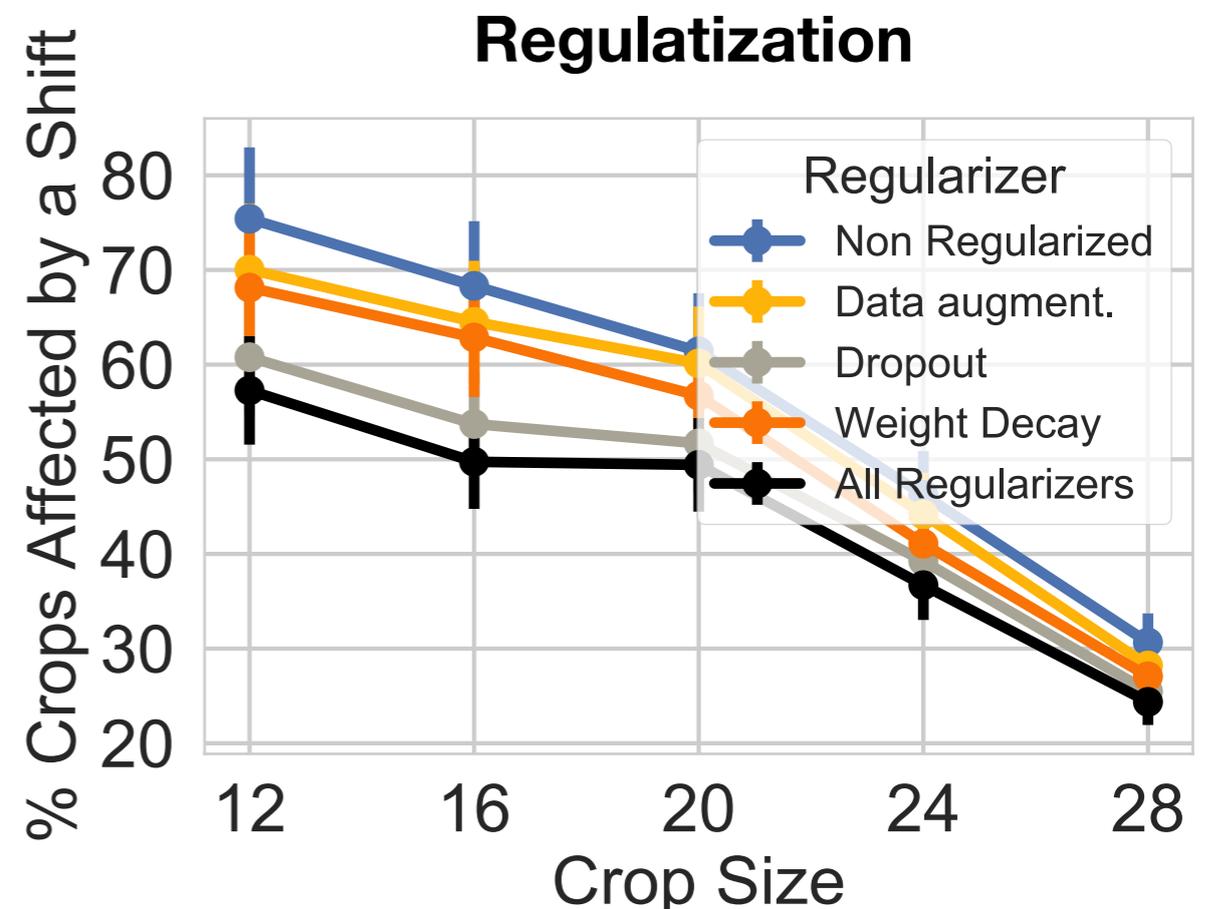
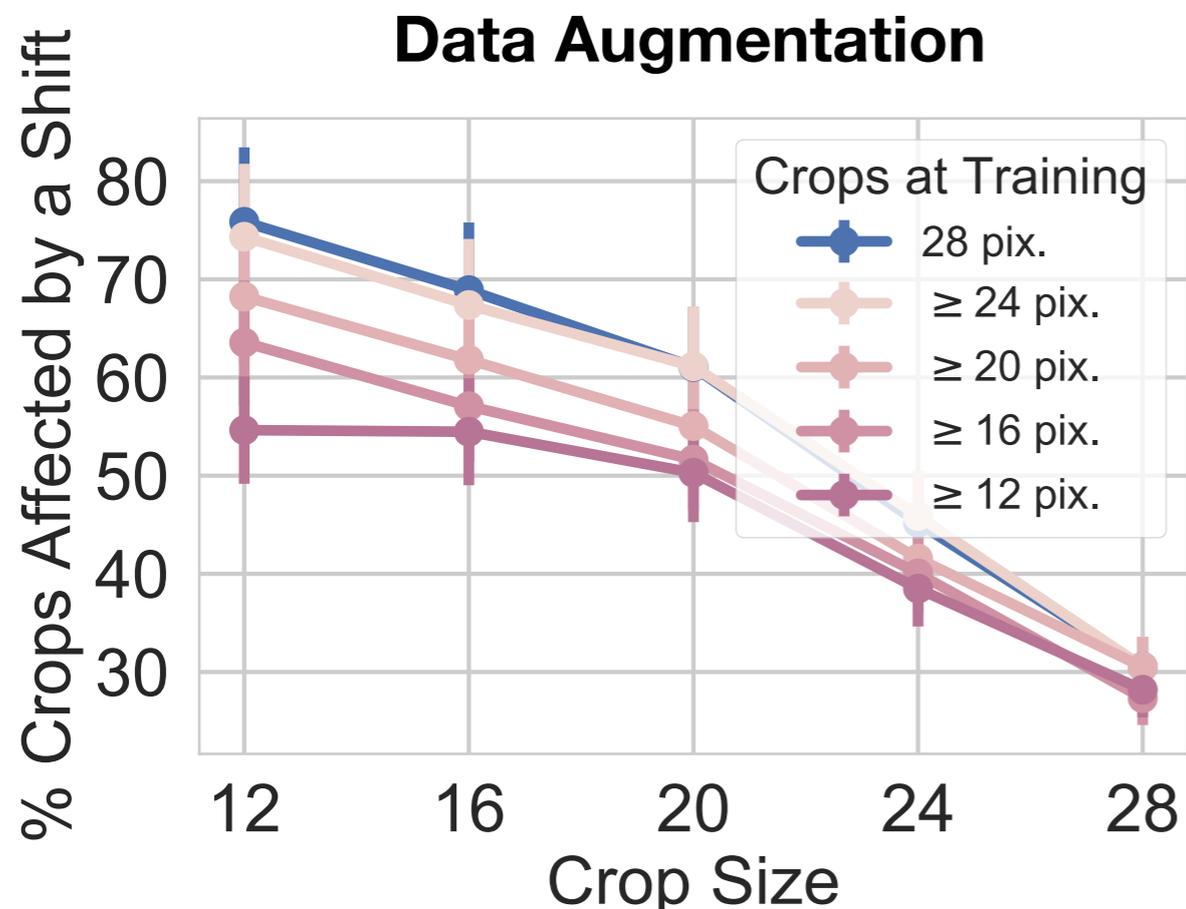
Pooling and Fragile Recognition



Increasing pooling size significantly improves translation invariance at small image sizes, but only somewhat improves minimal images.

FRIs are related to location invariance but also depend on other factors

Can Fragile Recognition be Alleviated?



Data augmentation helps generalization to limited information; regularizers also help

Neither eliminates the problem

VGG16, Resnet, Inception already have data augmentation of relatively large crops

3 Snapshots

1) Failures of DNNs

Minimal images are a common phenomenon among humans and DNNs

How can we make DNNs robust to minimal images as humans?

2) Beyond Object Recognition

3) DNNs to explain the Brain

3 Snapshots

1) Failures of DNNs

Minimal images are a common phenomenon among humans and DNNs

How can we make DNNs robust to minimal images as humans?

2) Beyond Object Recognition

Insideness is solvable but not learnable by state-of-the-art DNNs

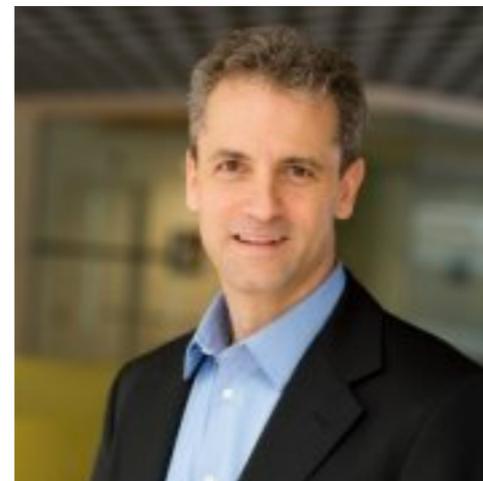
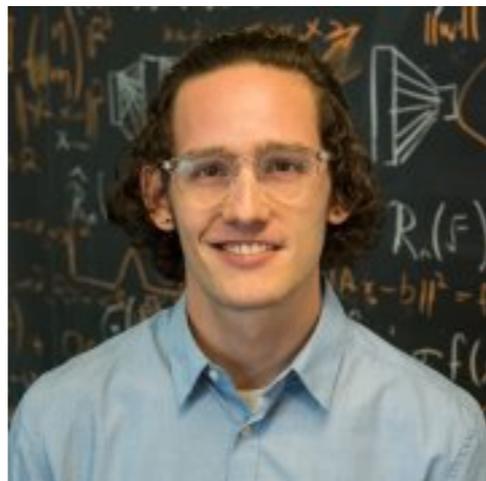
What learning mechanisms could lead to the general solution?

3) DNNs to explain the Brain

Single units in a DNN correspond with neurons in the brain

with

*Luke Arend, Yena Han, Martin Schrimpf, Pouya Bashivan, Kohitij Kar,
Tomaso Poggio, James DiCarlo*

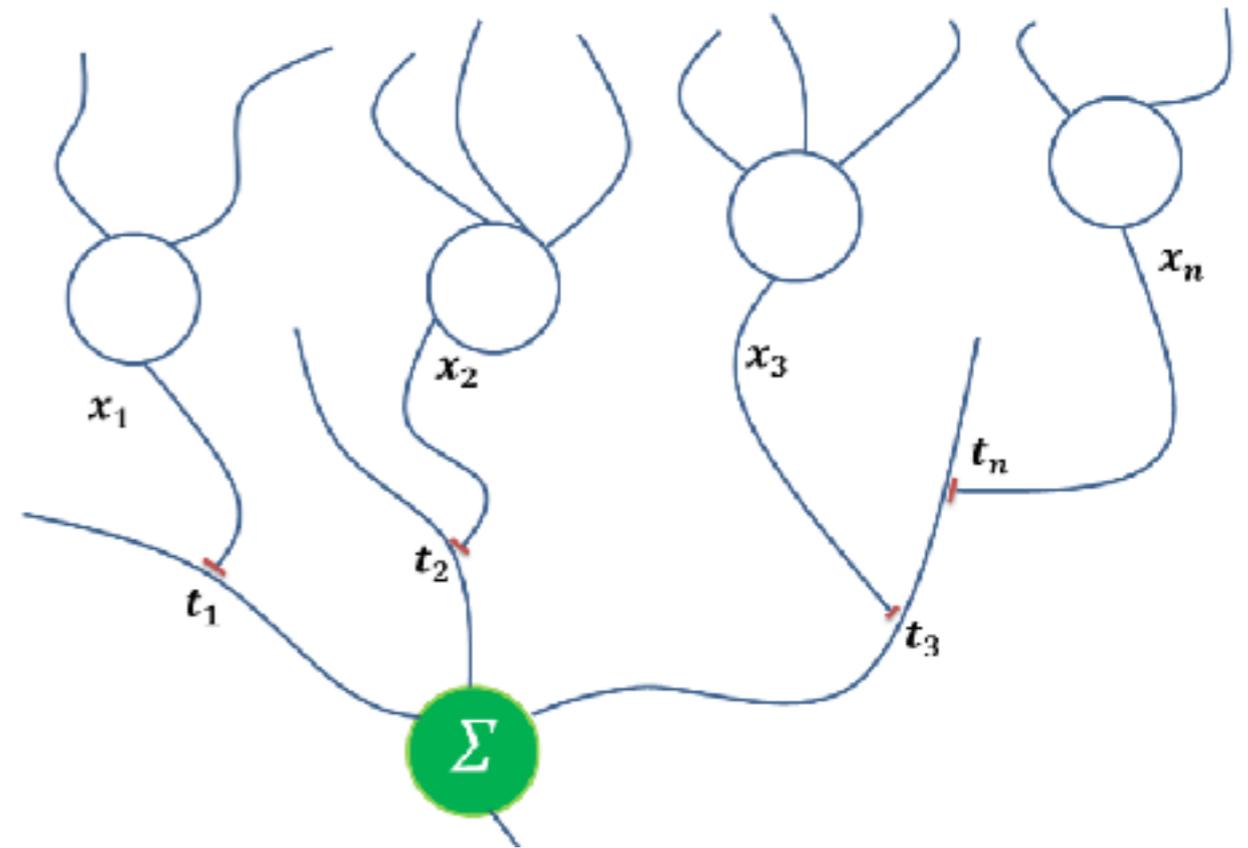


Model of a neuron

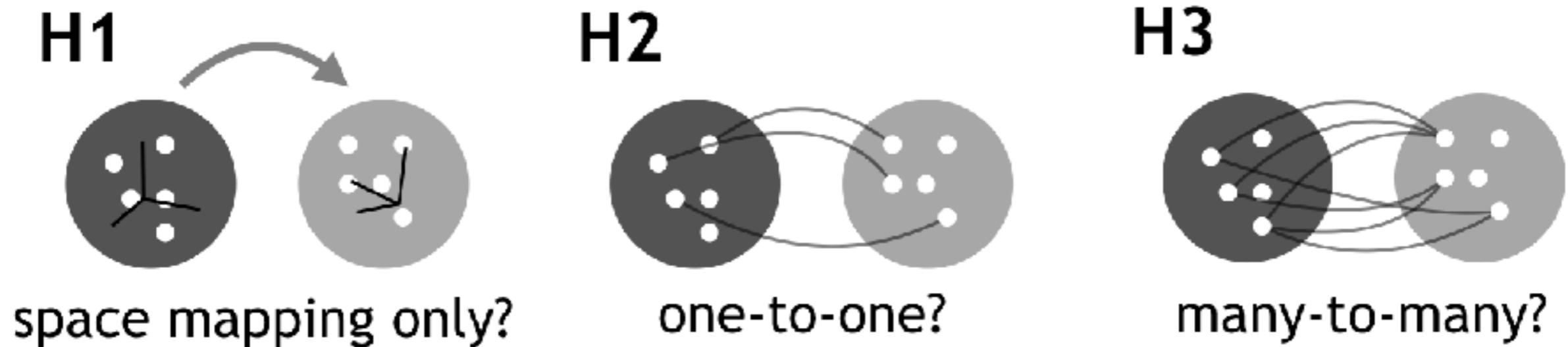
Simplified neuroscience: a neuron computes a dot product between its inputs and the synaptic weights

$$\langle x, t \rangle \longleftrightarrow$$

Neuroscience definition of dot product!



Can we map individual DNN units to individual neurons in the brain?



If yes,

- 1) DNN are better models of the brain than previously thought
- 2) the unit to neuron mapping is more interpretable than a “linear combination of units” to neuron

1) is there a one-to-one mapping?

2) which cells in the model match the neural recordings?

3) which cells in the model don't match the neural recordings?

1) is there a one-to-one mapping?

2) which cells in the model match the neural recordings?

3) which cells in the model don't match the neural recordings?

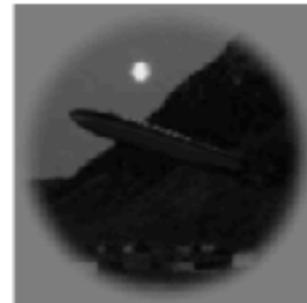
stimuli

- 2560 images
- 8 categories
- high variation in location, pose and size

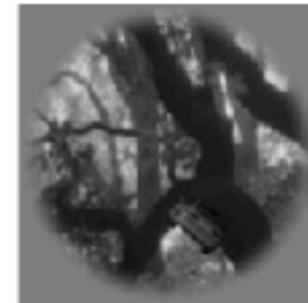
animals



boats



cars



chairs



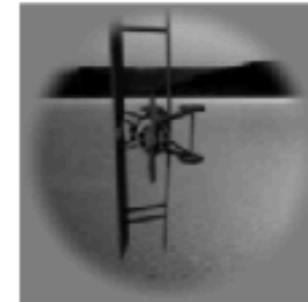
faces



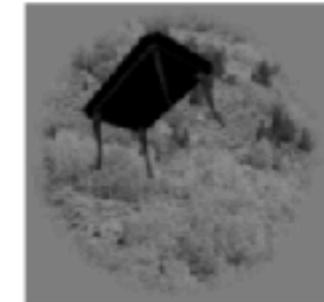
fruits



planes



tables



neural recordings

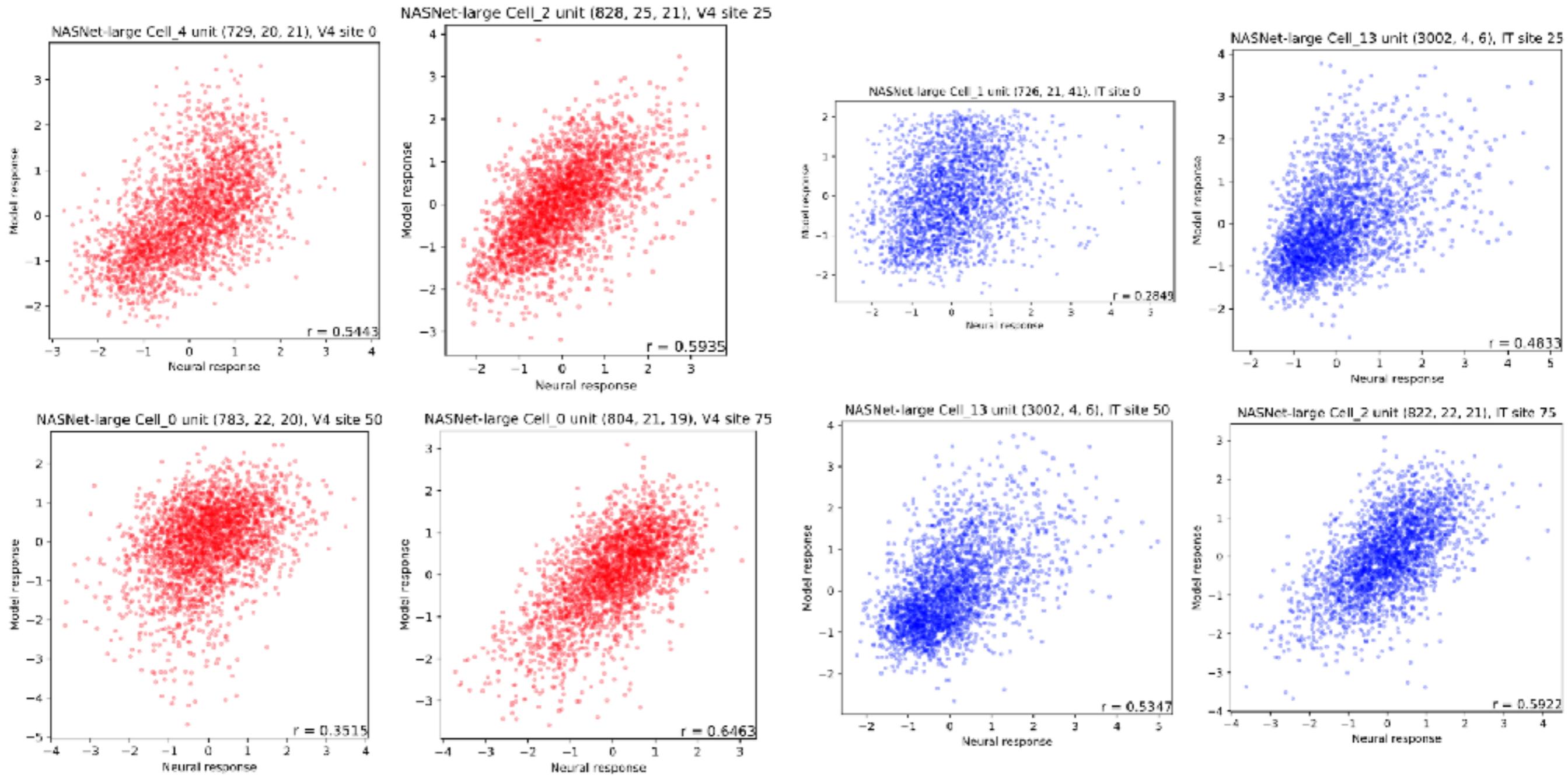
- 88 V4 sites, 168 IT sites
- previously published in [1, 2, 3]

models

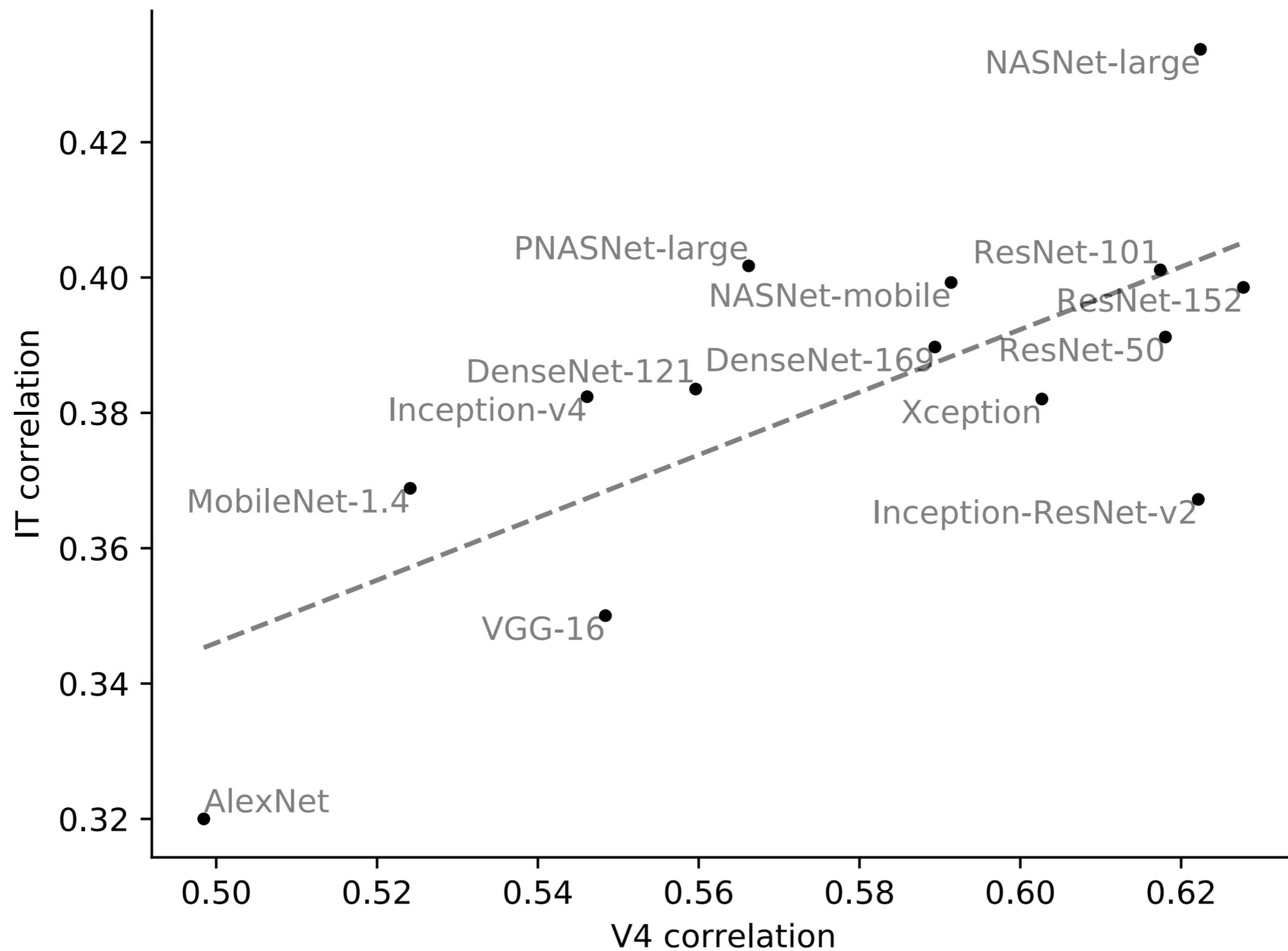
- 14 state-of-the-art deep neural networks
- optimized for object recognition

V4

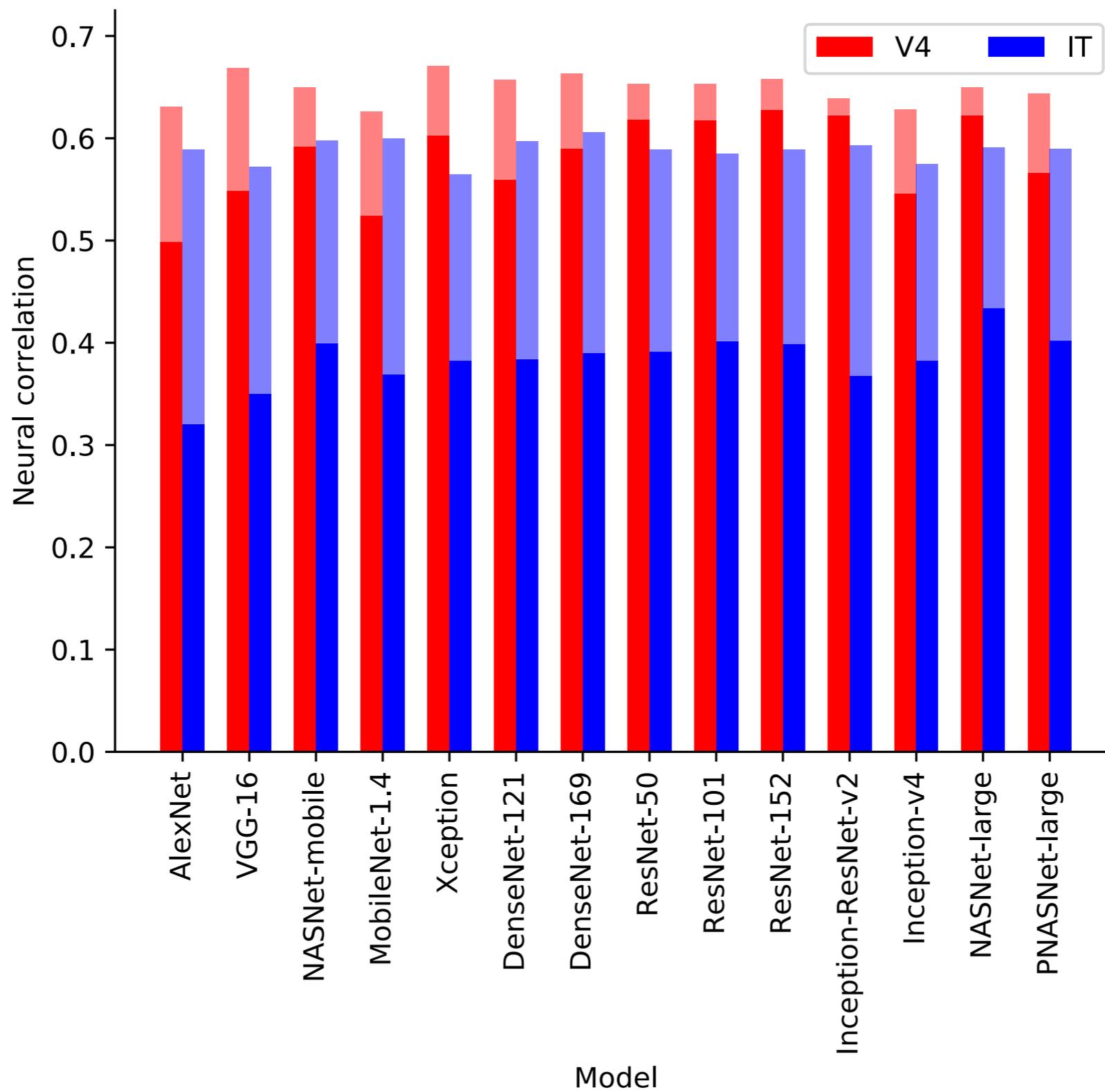
IT



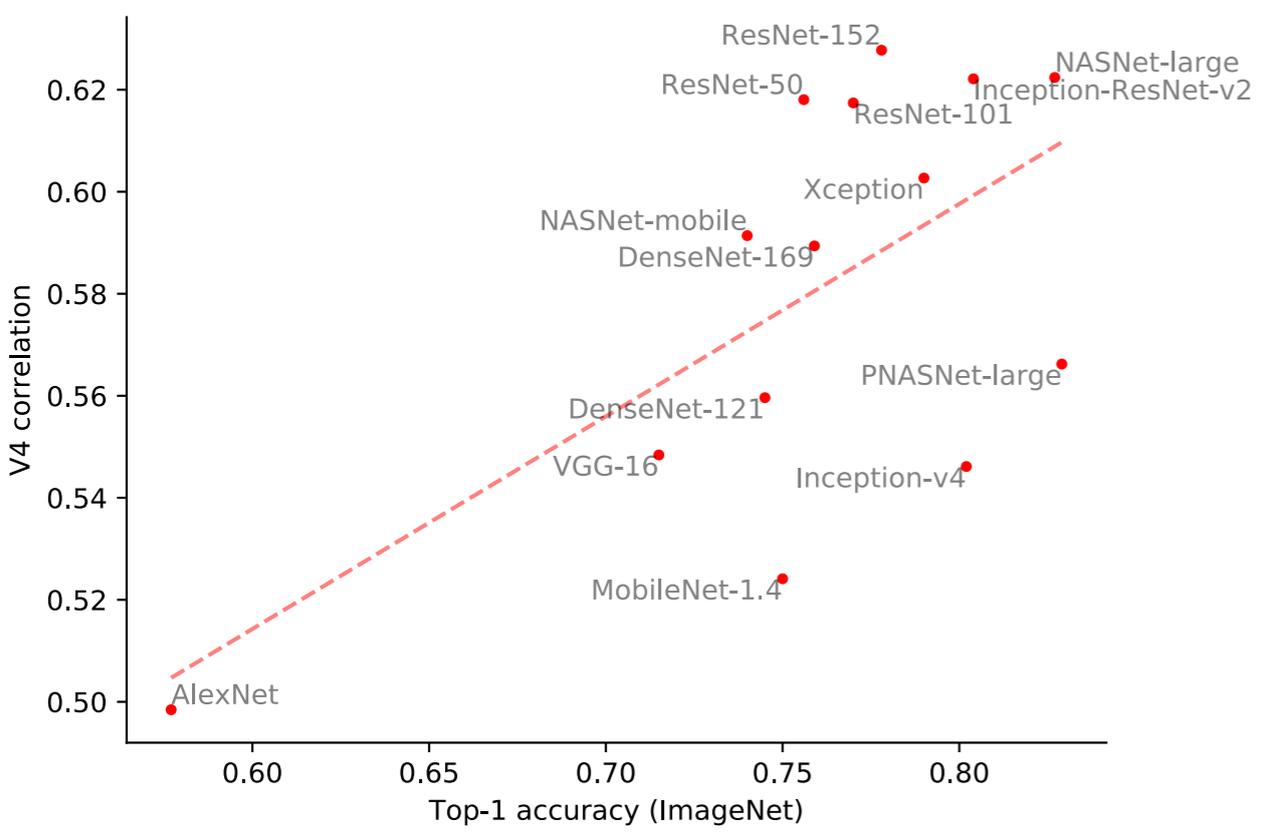
median across population of correlation for best-matched units



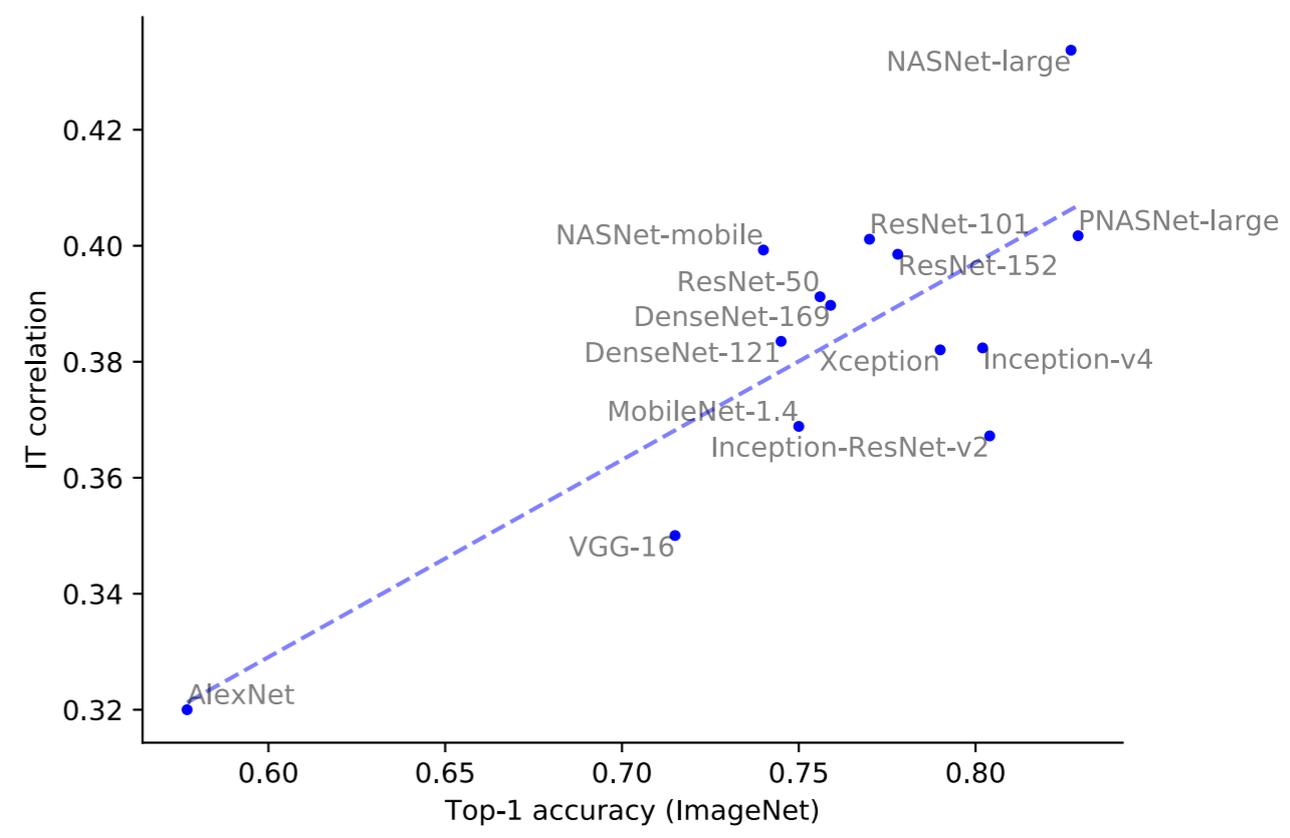
Linear Combination vs Individual Units



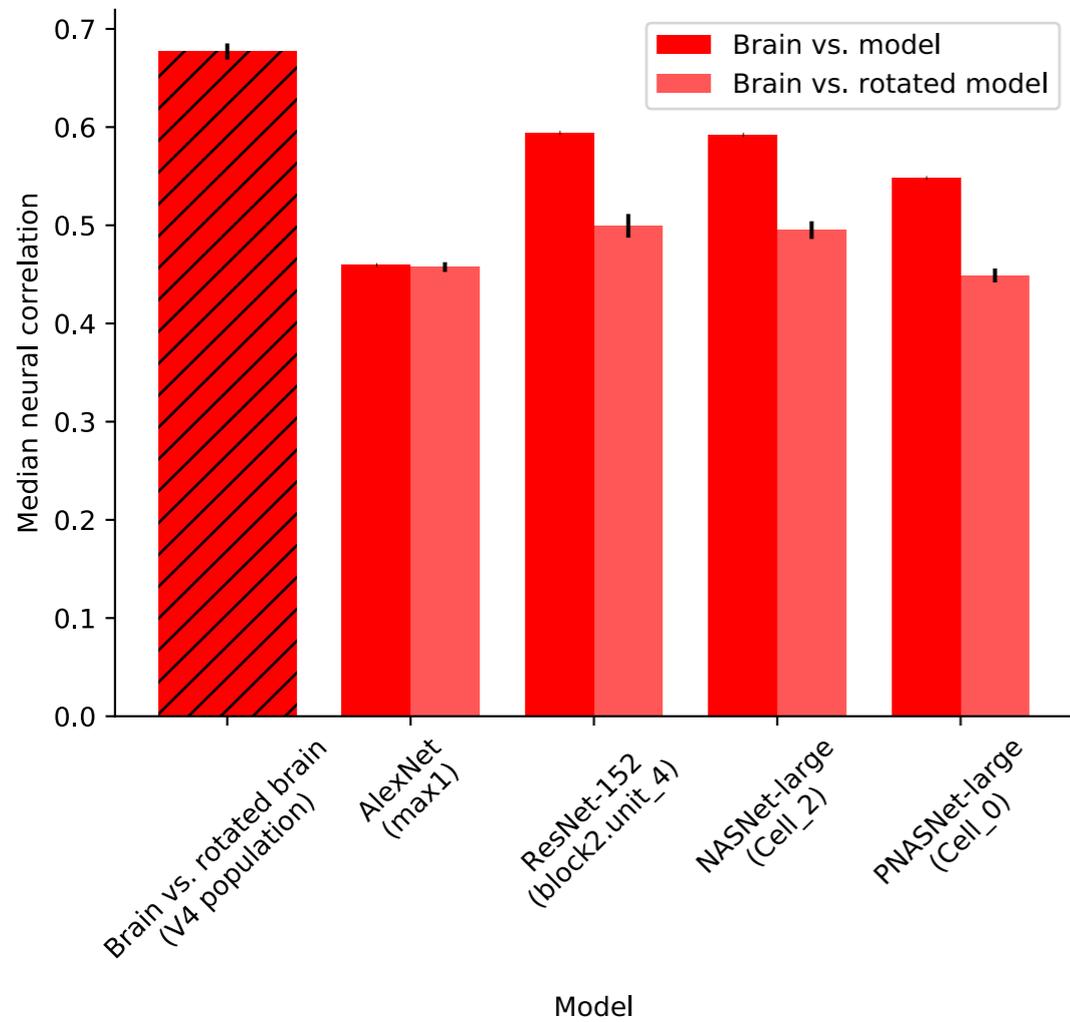
V4



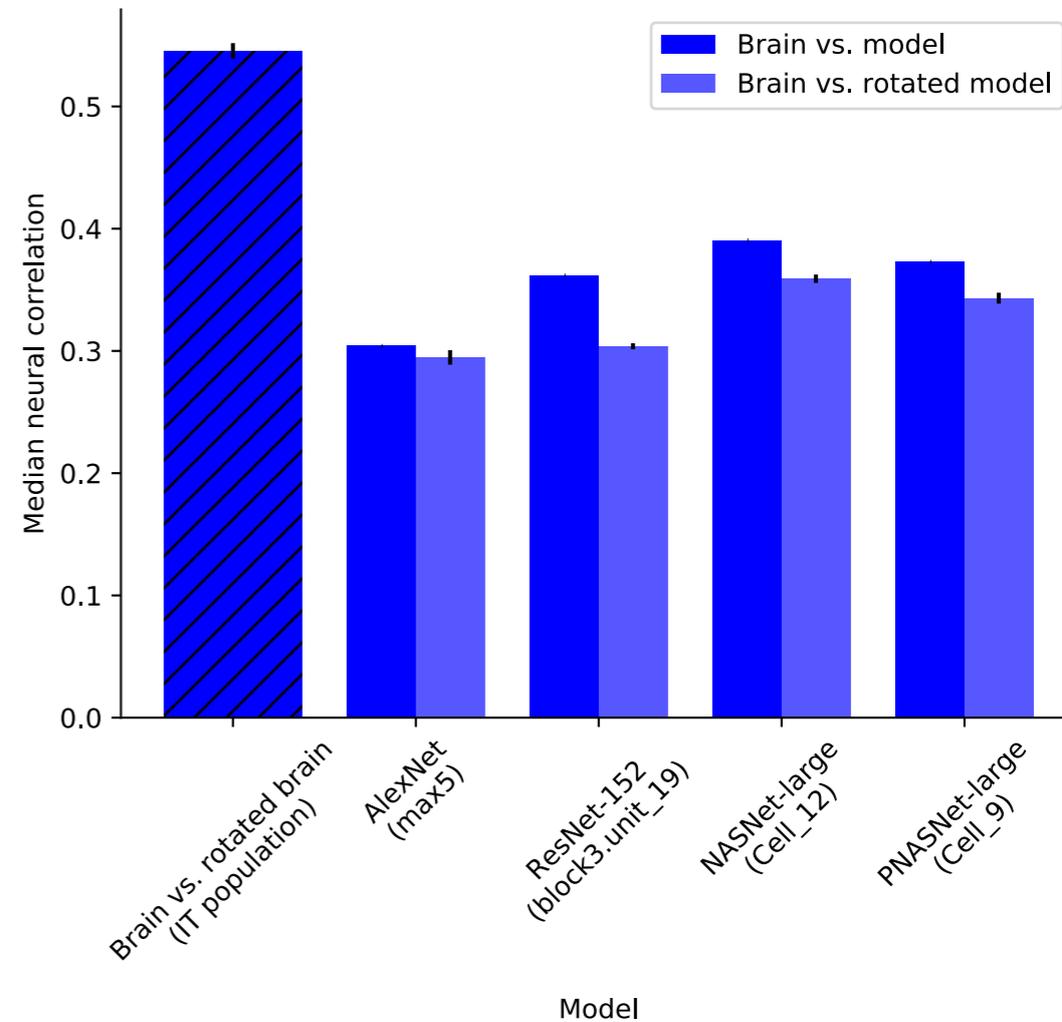
IT



V4



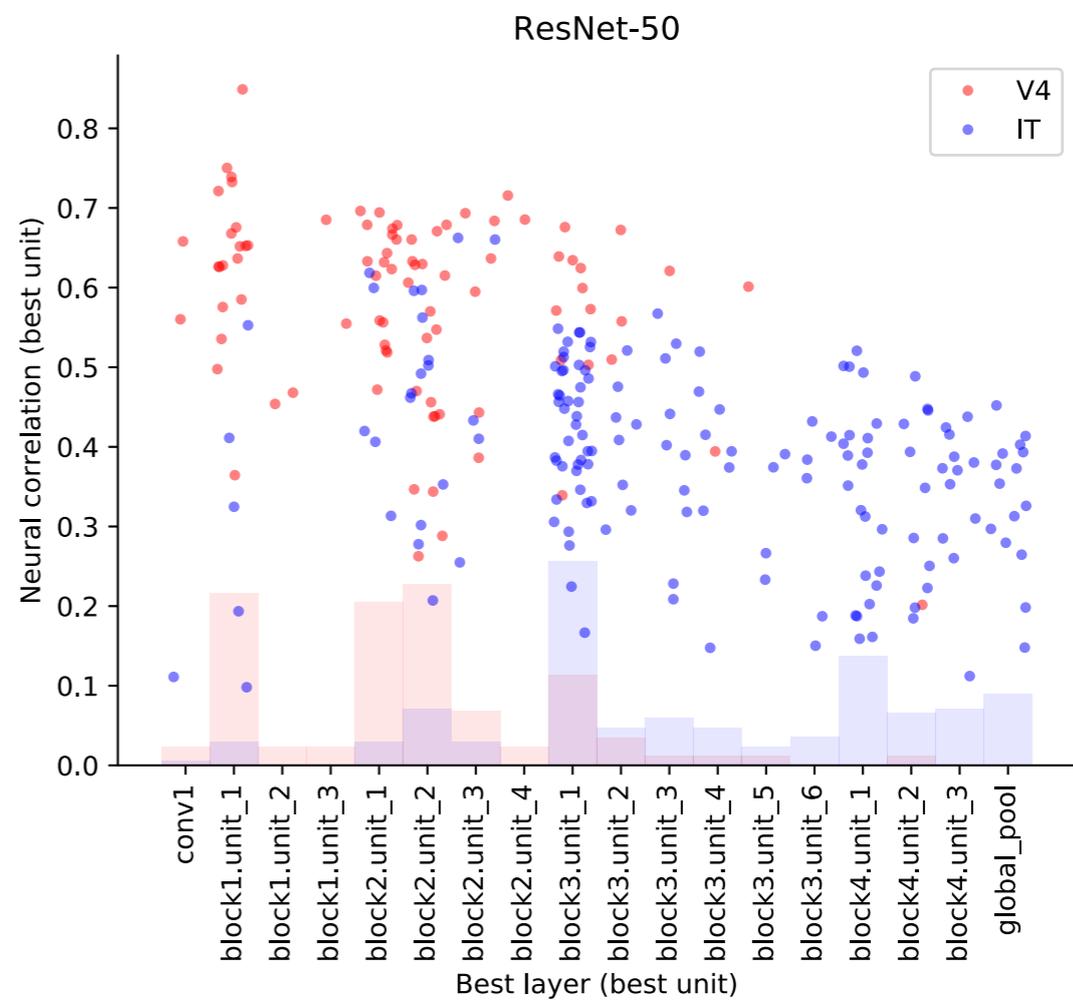
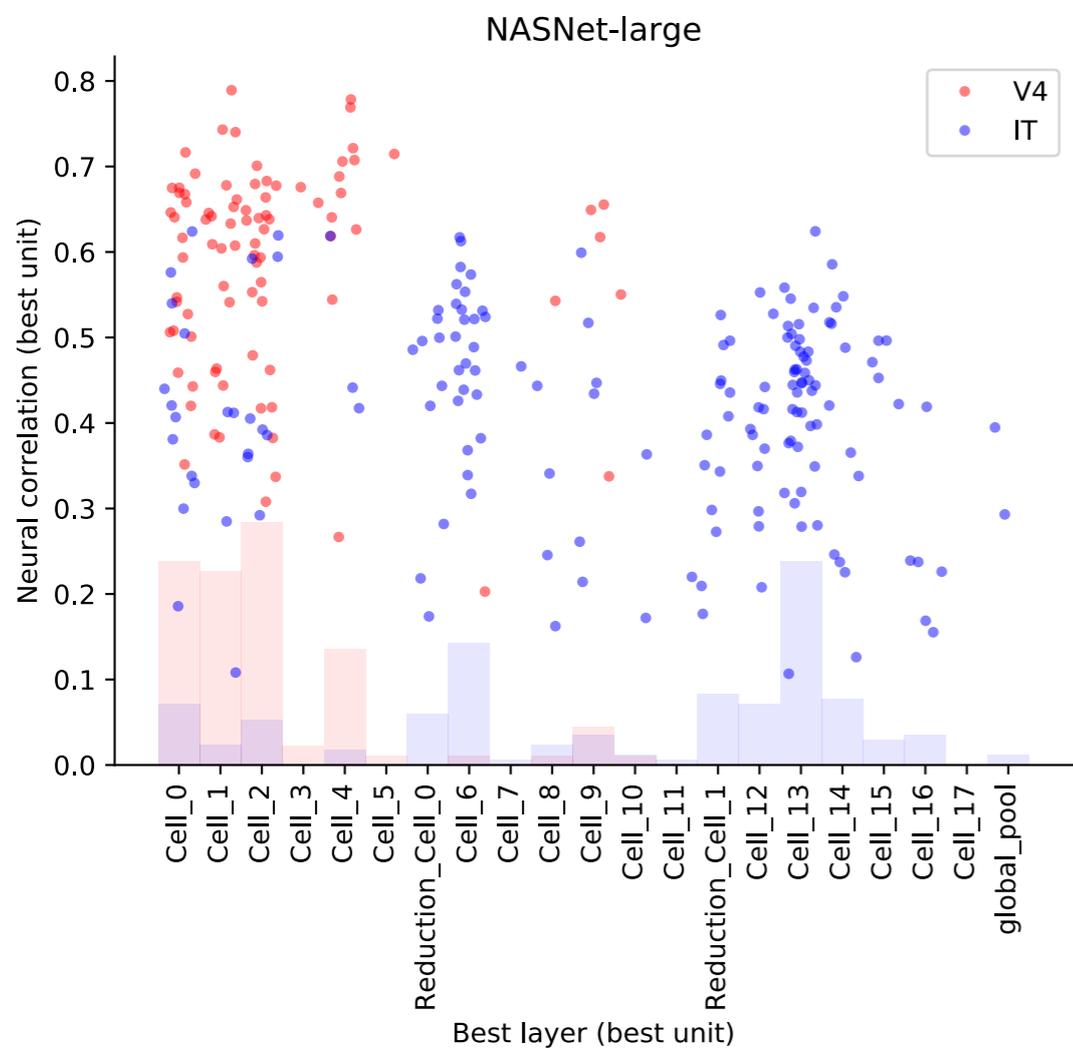
IT

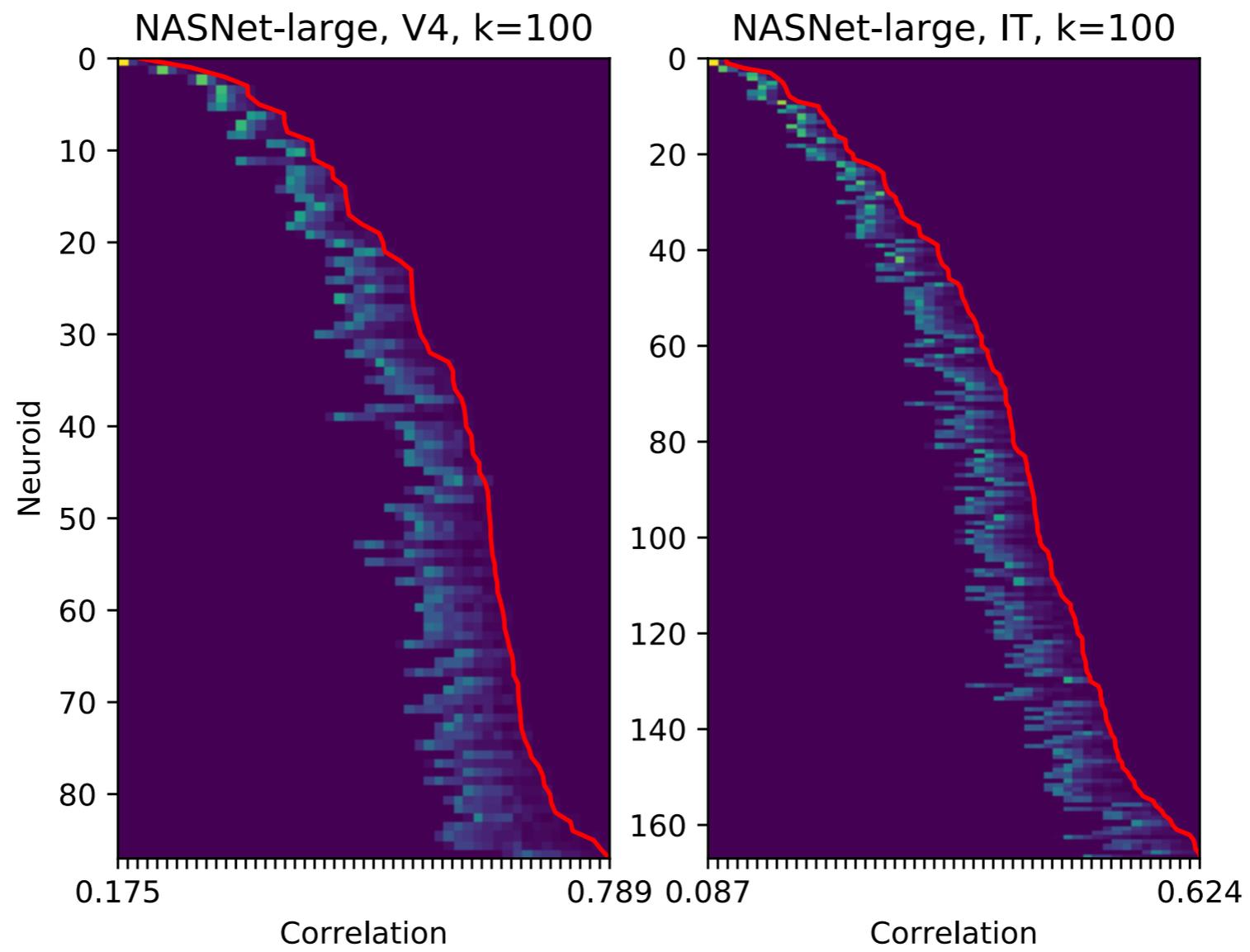


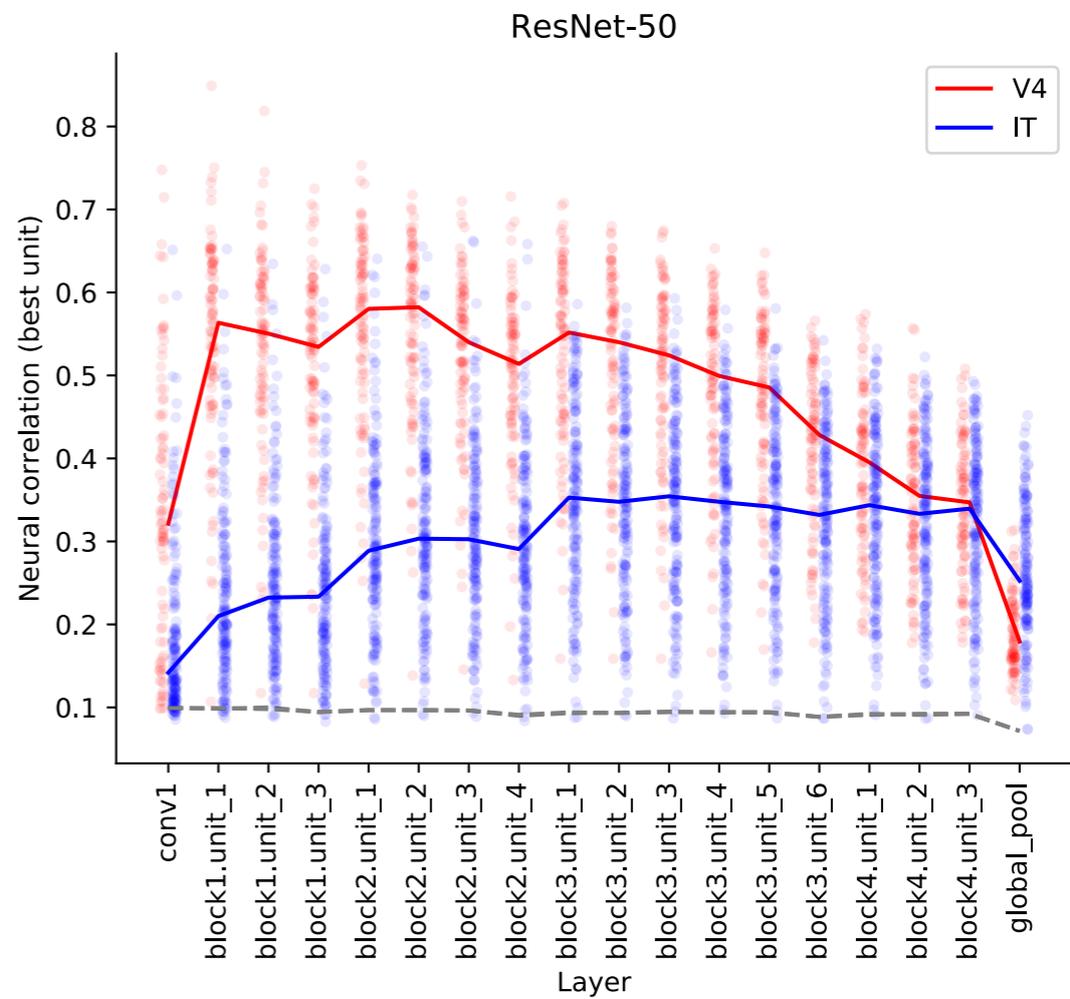
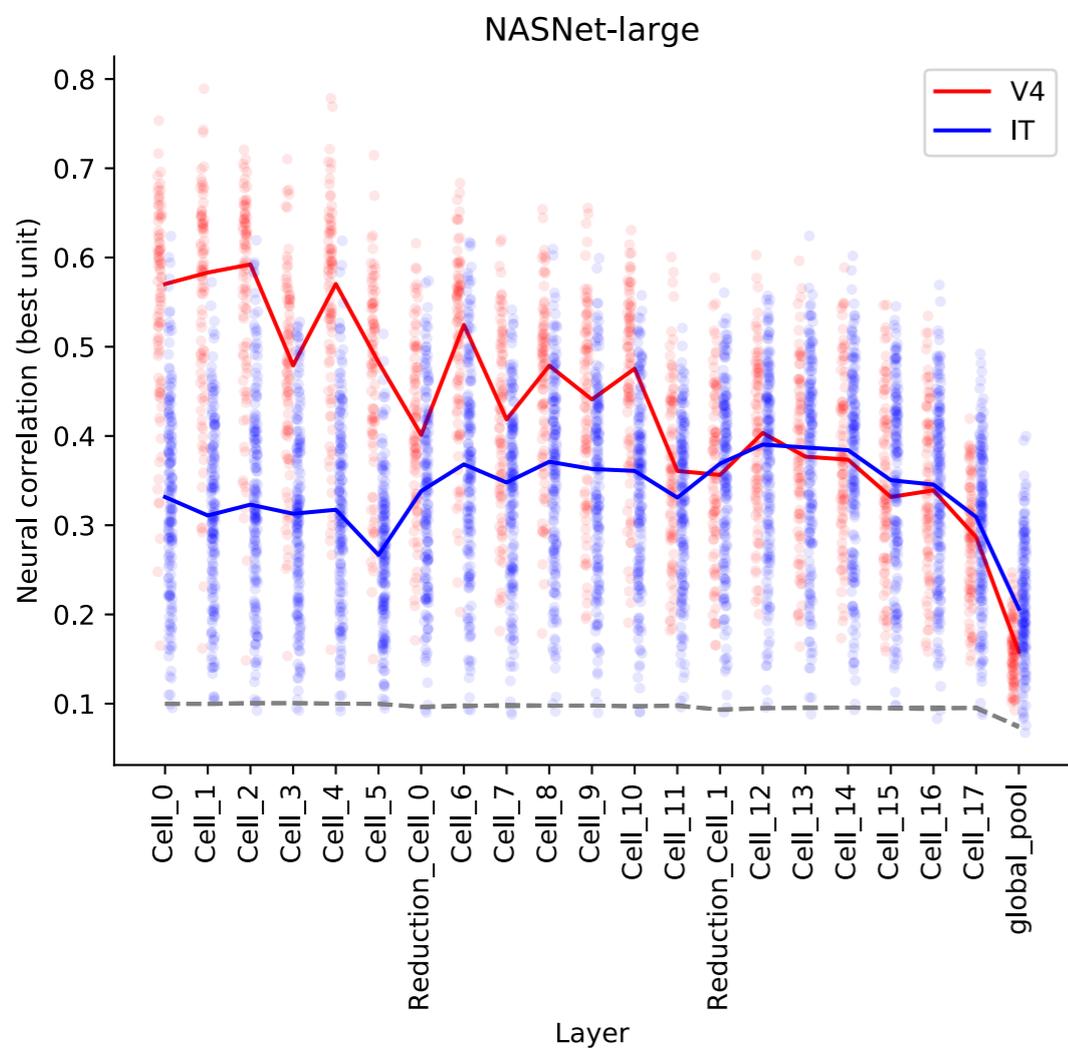
1) is there a one-to-one mapping?

2) which cells in the model match the neural recordings?

3) which cells in the model don't match the neural recordings?





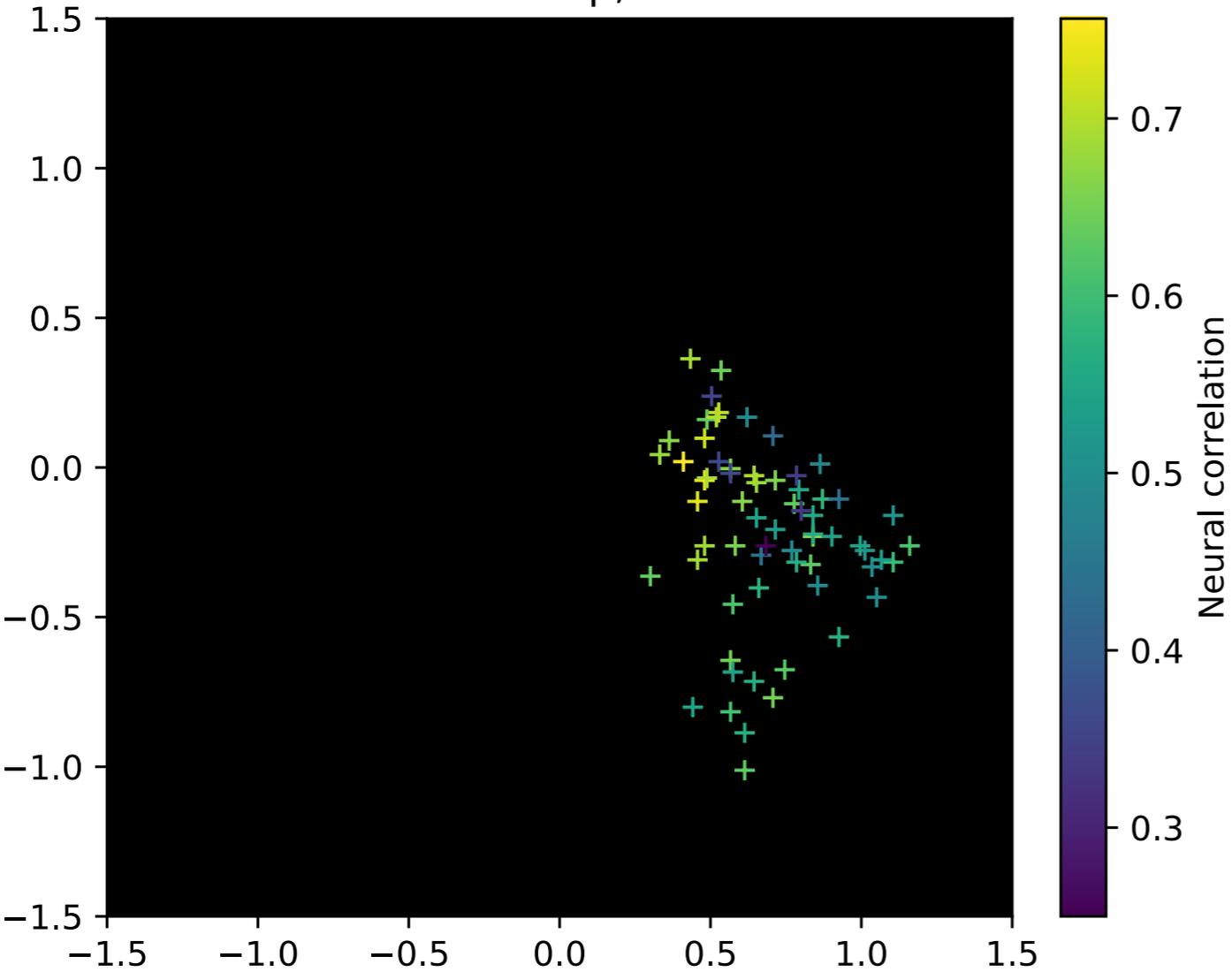


1) is there a one-to-one mapping?

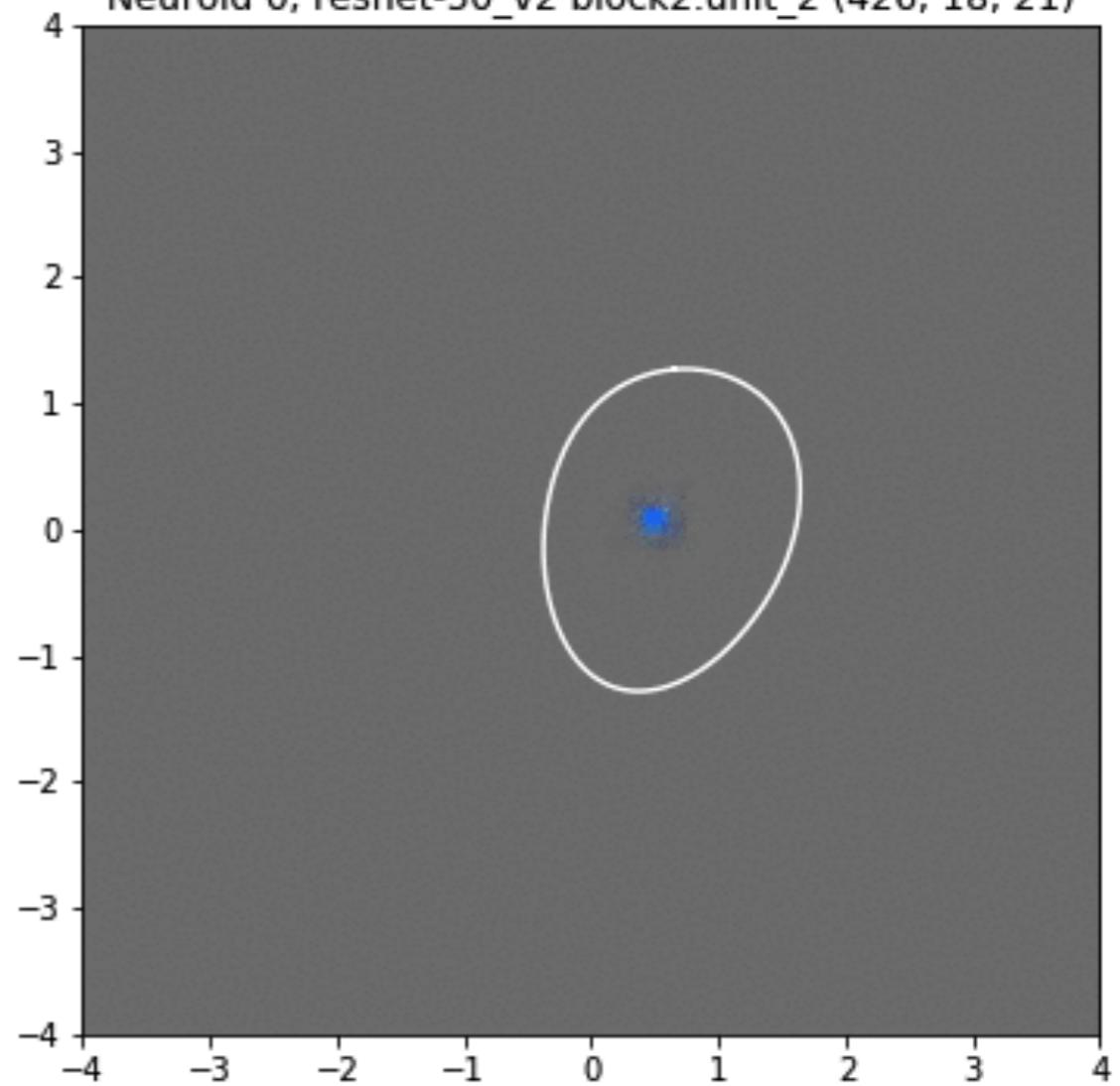
2) which cells in the model match the neural recordings?

3) which cells in the model don't match the neural recordings?

V4 neuroids map, ResNet-50



Neuroid 0, resnet-50_v2 block2.unit_2 (426, 18, 21)



3 Snapshots

1) Failures of DNNs

Minimal images are a common phenomenon among humans and DNNs

How can we make DNNs robust to minimal images as humans?

2) Beyond Object Recognition

Insideness is solvable by state-of-the-art DNNs for image segmentation

What learning mechanisms could lead to the general solution?

3) DNNs to explain the Brain

Single units in DNNs correspond to neurons in V4 and IT (to lesser degree)

What leads to this correspondence?

1) Computational

Same principles should explain
artificial and biological intelligence

(eg. Why the brain is how it is?
eg. Why an algorithm works better than others?)

2) Algorithmic

Computer Vision
(Build better algorithms)

Modelling the Brain
(How does the brain work?)

3) Implementation



CENTER FOR
**Brains
Minds+
Machines**

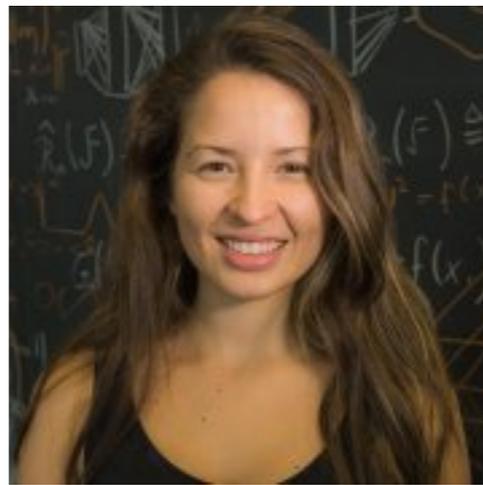
Making a Science from the Computer Vision Zoo

Xavier Boix

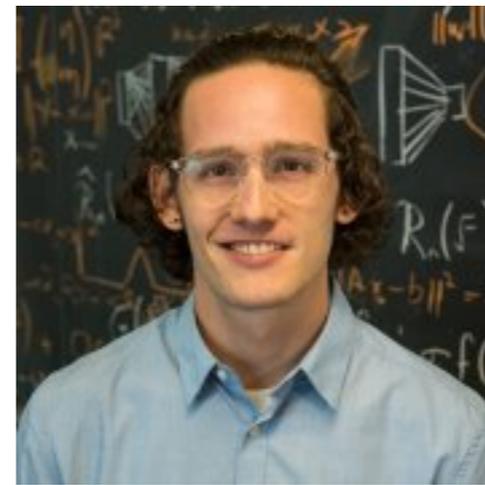
xboix@mit.edu



Sanjana



Kim



Luke



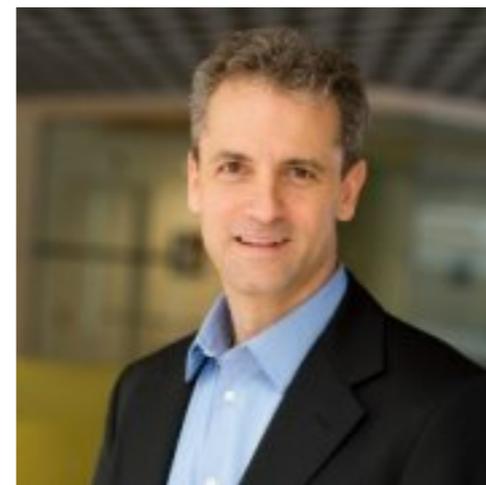
Guy



Tomotake



Fred



Prof. DiCarlo



Prof. Poggio