# Neural mechanisms underlying visual object recognition:
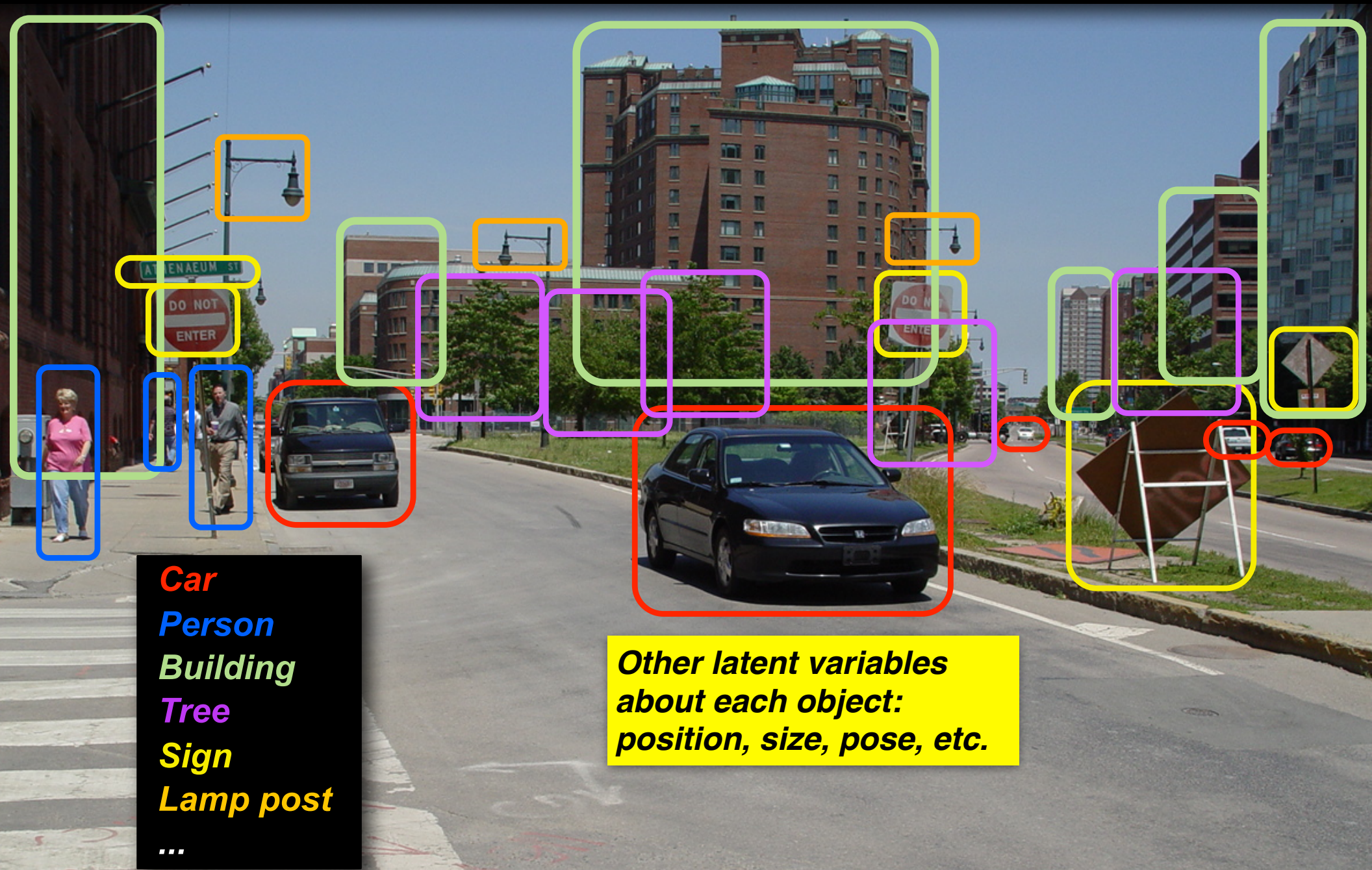## The convergence of computer vision and biological vision

*Center for Brains, Minds, and Machines:  Summer School 2015, Woods Hole, MA*

**James DiCarlo MD, PhD**

*Professor of Neuroscience and Head, Department of Brain and Cognitive Sciences*
*Investigator, The McGovern Institute for Brain Research*
*Massachusetts Institute of Technology, Cambridge MA, USA*

brain+cognitive sciences

MIT

# "Object recognition" (operationalized)



**Car**
**Person**
**Building**
**Tree**
**Sign**
**Lamp post**
**...**

*Other latent variables about each object: position, size, pose, etc.*

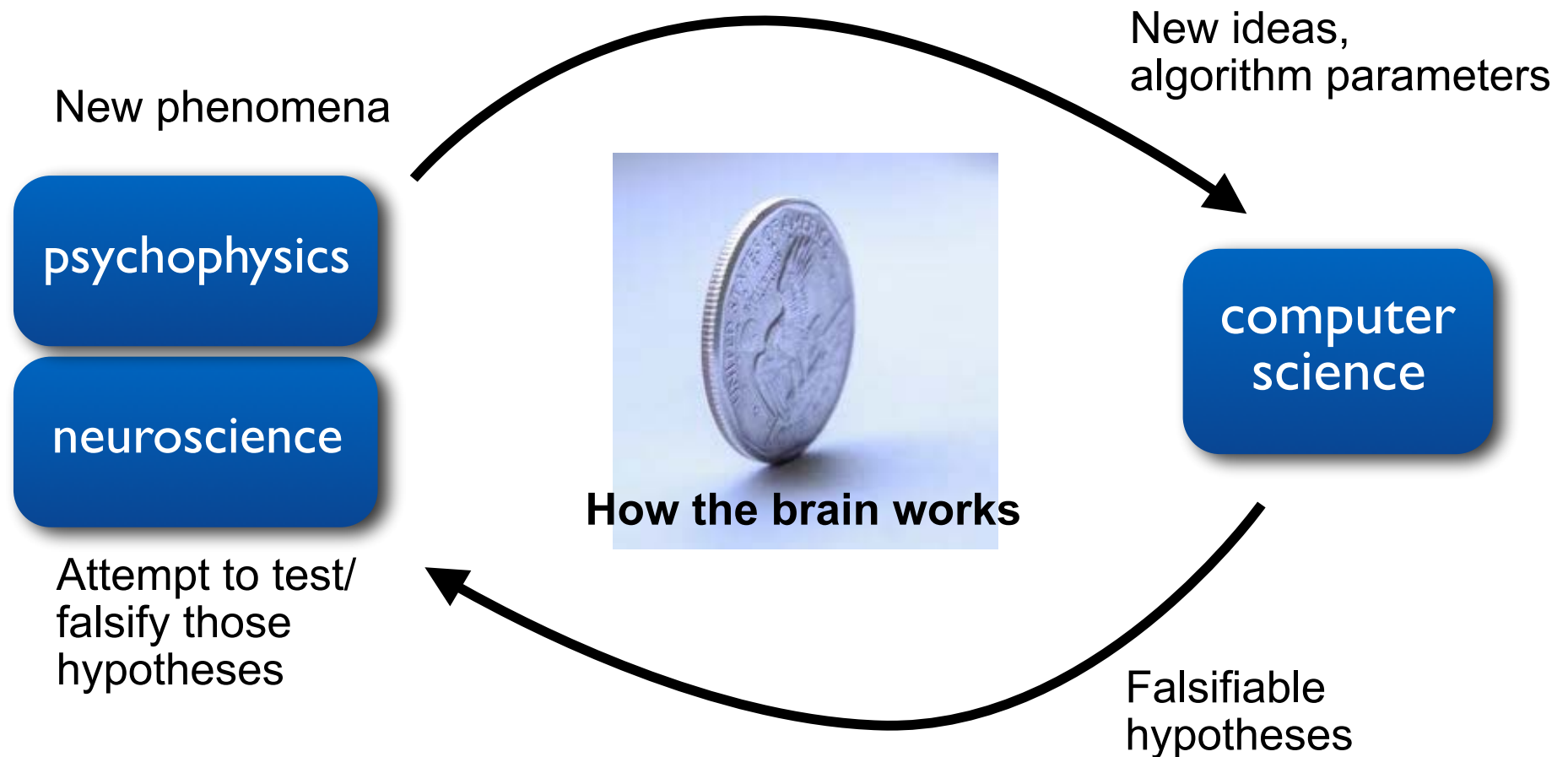*Image adapted from MIT Street Scenes Database (Courtesy of Tommy Poggio)*

# Why study object recognition in the brain?

*The brain's internal representation of objects is the substrate of cognition:*

- *memory*
- *value judgements*
- *decisions*
- *actions*

- *Obstacle avoidance*
- *Navigation*
- *Danger avoidance*
- *Resource detection*
- *Social interactions*
- *Mate selection*
- *Threat detection*
- *Reading*
- *...*

**When biological brains perform better than computers**

New phenomena

New ideas,
algorithm parameters

psychophysics

neuroscience

computer
science

**How the brain works**

Attempt to test/
falsify those
hypotheses

Falsifiable
hypotheses

**When computers perform as well as or better than biological brains**

# A bit of history…

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

PROJECT MAC

Artificial Intelligence Group                                    July 7, 1966
Vision Memo. No. 100.

THE SUMMER VISION PROJECT

The final goal is OBJECT IDENTIFICATION which will actually name objects by matching them with a vocabulary of known objects.

Goals - Specific

We plan to work by getting a simple form of the system going as soon as possible and then elaborating upon it.  To keep the work reasonably coordinated there is a graduated scale of subgoals.

*Courtesy of Mike Tarr*

- *100 billion computing elements*

- *solves problems not soluble by previous machines*
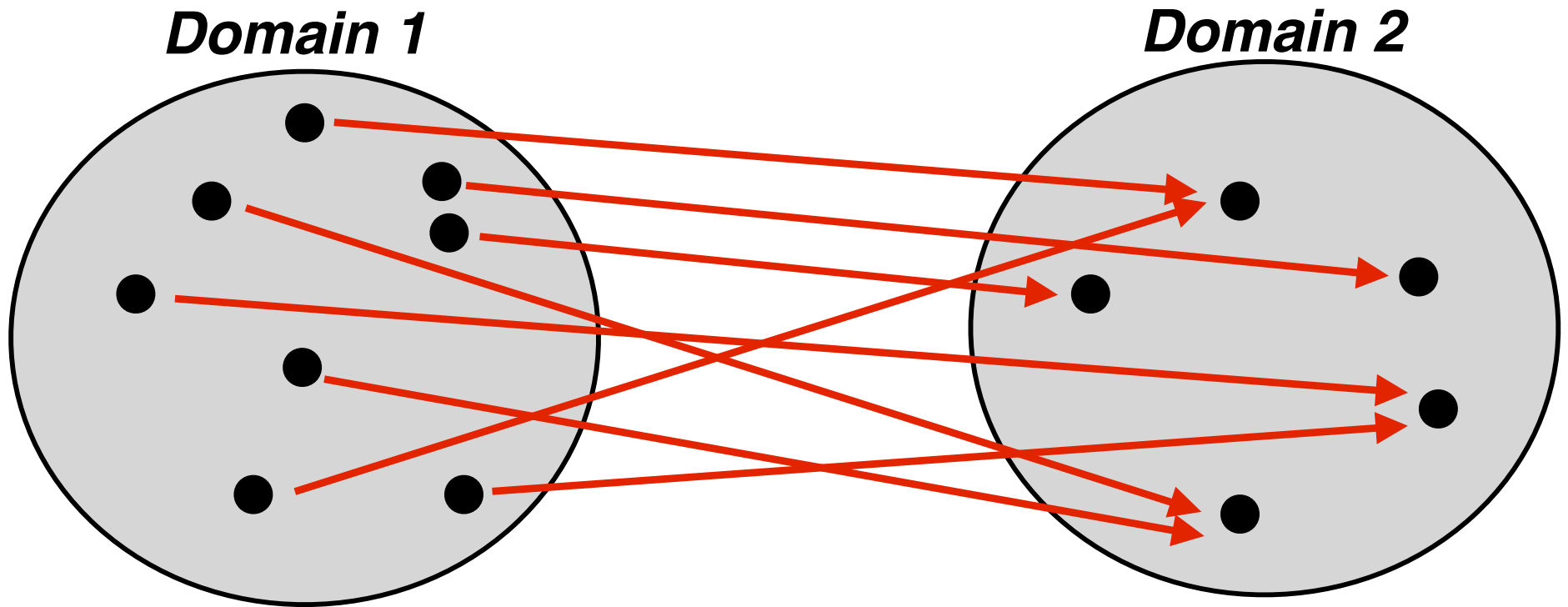
- *requires only 20 watts of power!*

*Key algorithms are **classified***

# An engineer's point of view…

**Which system is better?**

| Problem to solve | Our brain | Machines today (e.g. computers) |
|---|---|---|
| **Calculation** | | **WINNER** |
| **Win at chess** | | **WINNER** |
| **Win at Jeopardy** | | **WINNER** |
| **"Memory"** | Gateway problem (vision, neocortex) | |
| **"Seeing"** | Our goal: Discover how the brain solves object recognition (algorithms) | |
| Pattern matching | | WINNER |
| Object recognition | **WINNER** | |
| Scene "understanding" | **WINNER** | |
| **Walking** | **WINNER** | |

# A scientist's point of view



**Domain 1**     **Domain 2**

*Science:  given state of Domain 1,
<u>predict</u> state of Domain 2*

*The accuracy of this predictive mapping is a
measure of the strength of a scientific field*

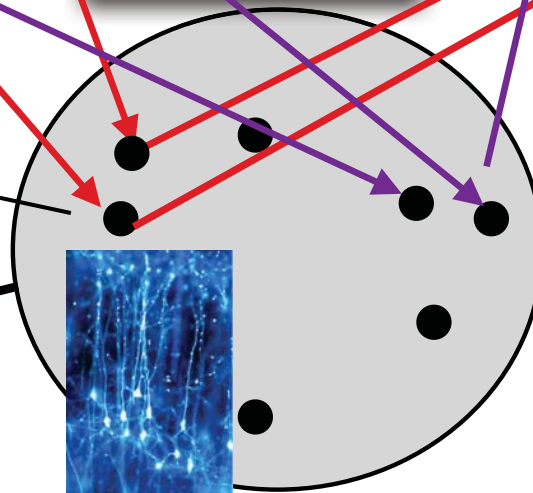**Images**

**Behavioral reports ("perception")**

"clock"

"cat"

"car"

"dog"

"face"

**Neural activity**

spiking pattern of some neural population in response to one image

**"Neural representation"**

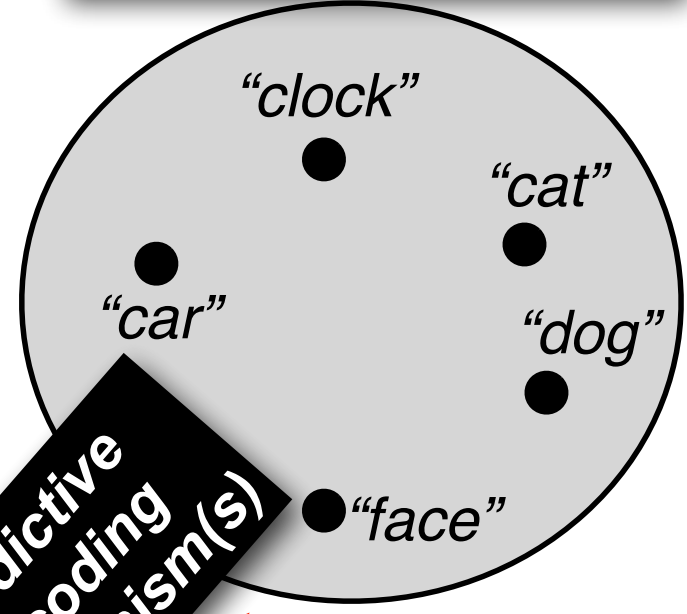**Accurate predictivity** is the core product of science → **Underlies engineer's ability to build, fix, or augment**

**Images**

**Behavioral reports ("perception")**

*"clock"*

*"cat"*

*"car"*

*"dog"*

*"face"*

**Predictive encoding mechanism(s)**

**Neural activ**

**Predictive decoding mechanism(s)**

*spiking pattern of some neural population in response to one image*

***"Neural representation"***

**For visual object perception, this link**

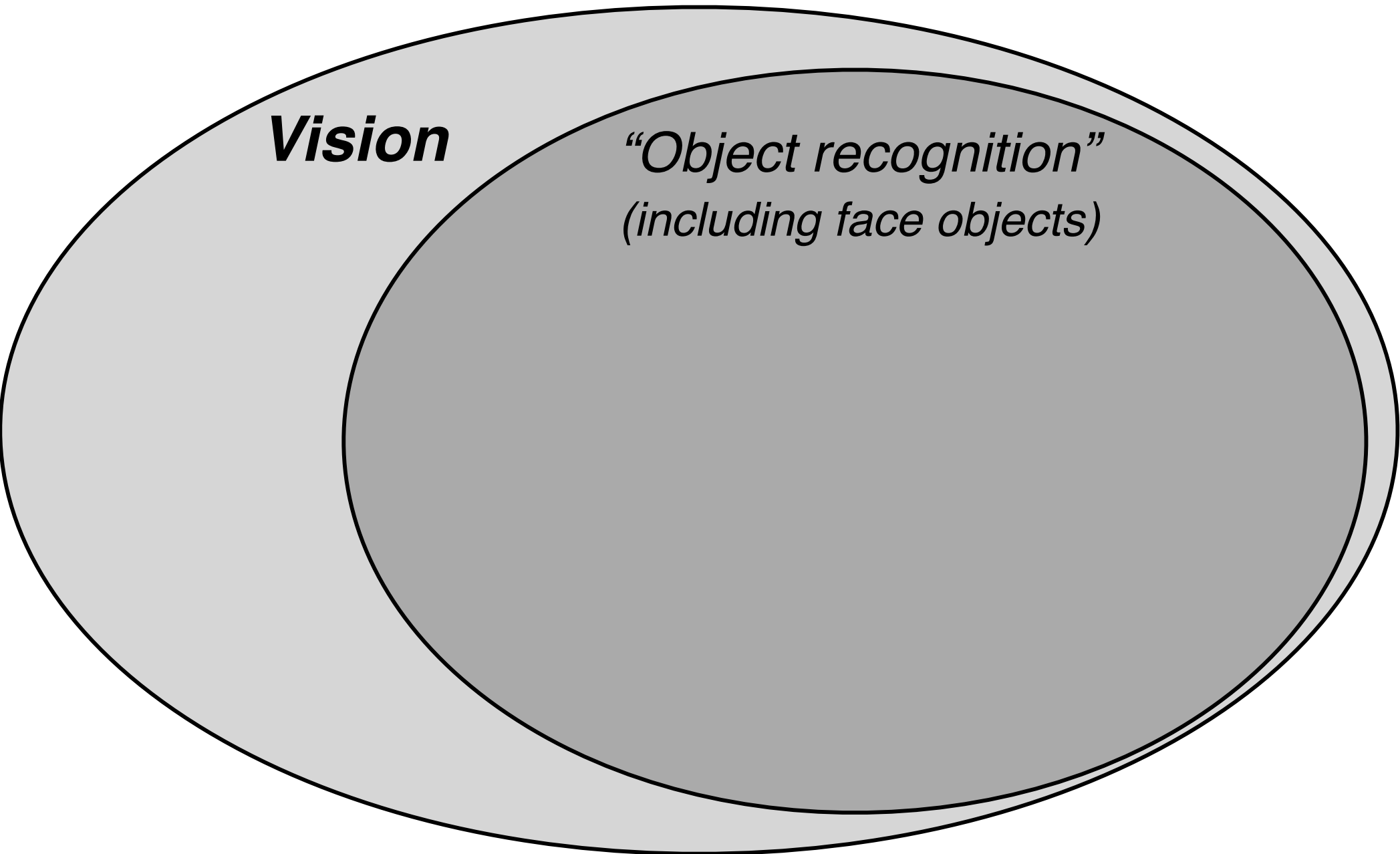**Not doubting the importance of these!**

**word models**

*"IT does object recognition"*

*"Face neurons do face tasks"*

*"Attention solves that"*

10

**Let's try to define a domain of behavior so that we can gauge/make progress in prediction.**

**Vision**

*"Object recognition"*
*(including face objects)*

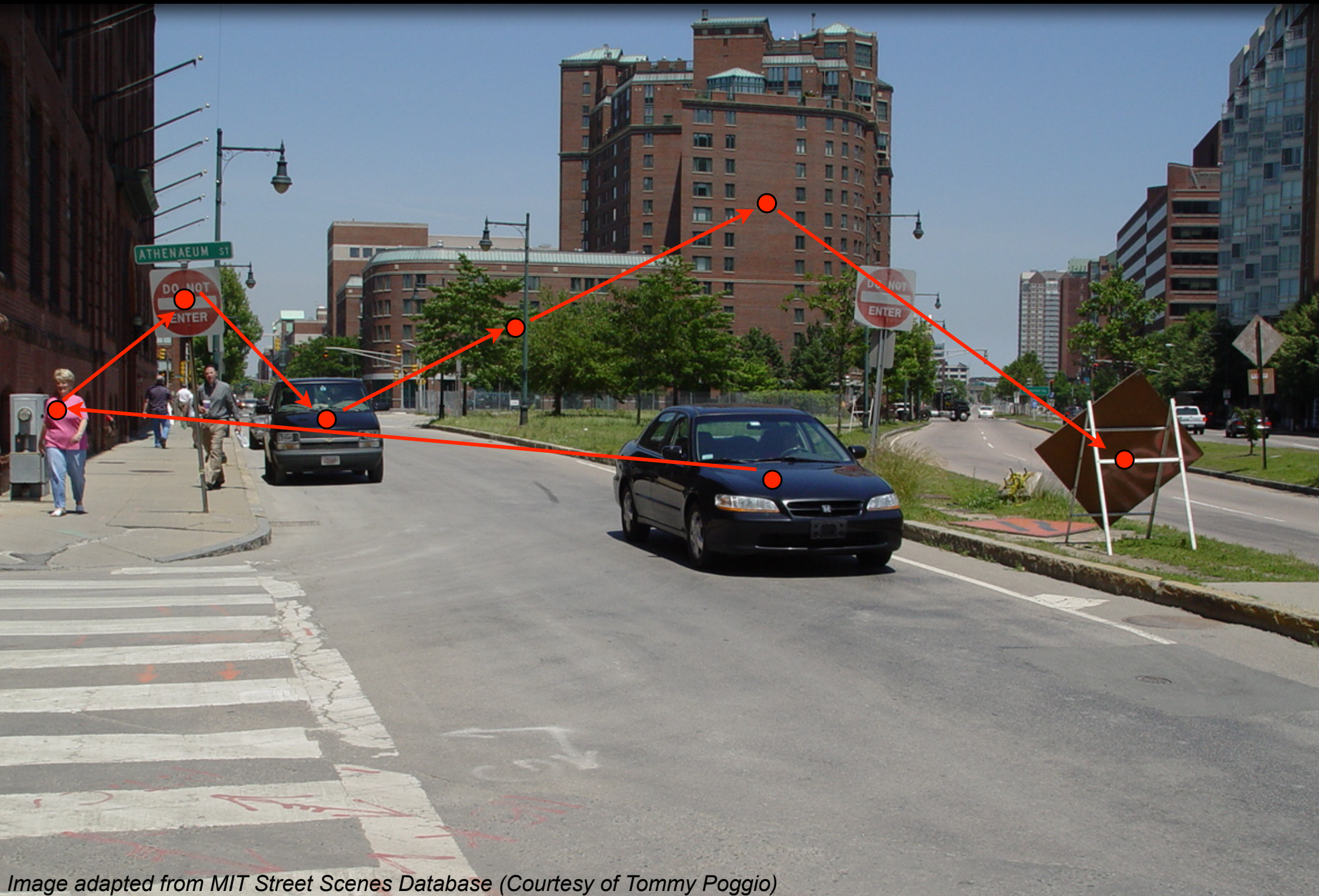**Object recognition as solved by primates**   ~200 ms snapshots

Image adapted from MIT Street Scenes Database (Courtesy of Tommy Poggio)

# Object recognition as solved by primates

**Core object recognition**

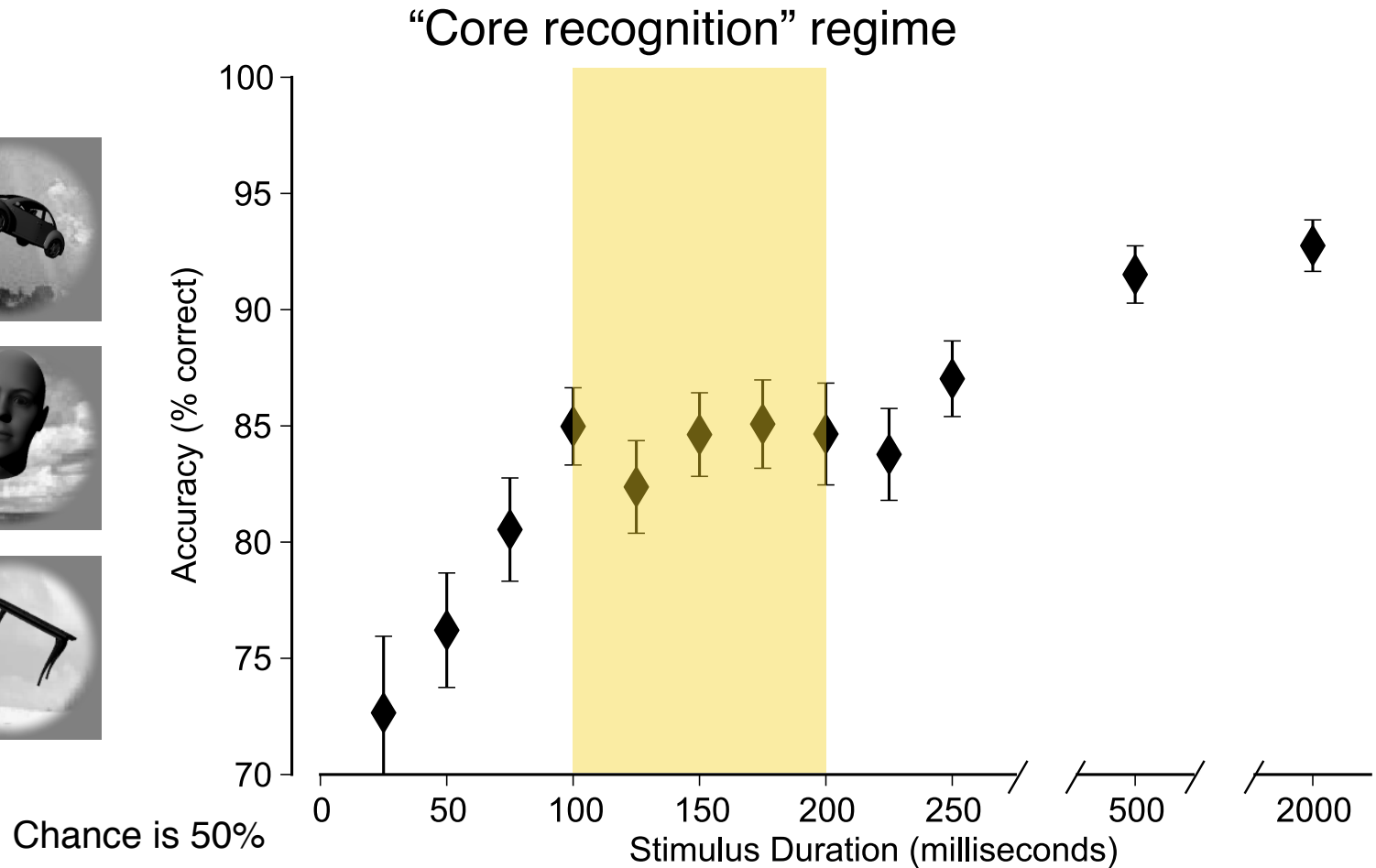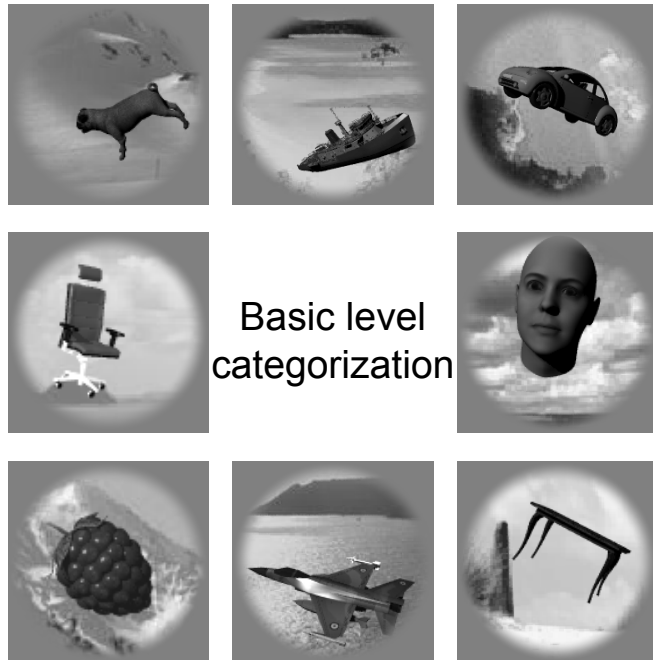central ~10 deg of visual field
100-200 ms viewing duration

# Our visual system excels at core object recognition

**Core object recognition**

central ~10 deg of visual field
100-200 ms viewing duration

# Human object recognition (categorization) accuracy as a function of image viewing time
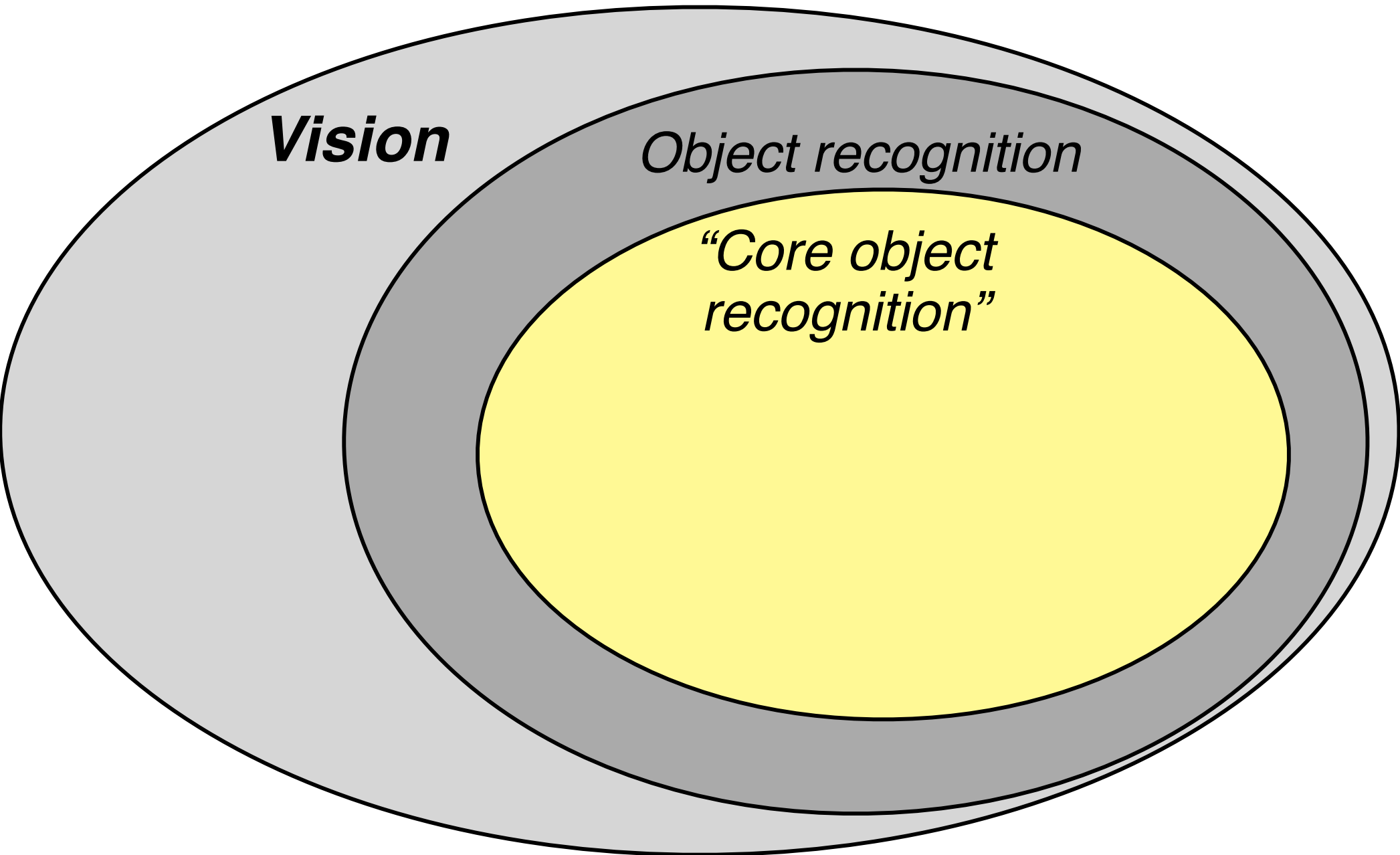


"Core recognition" regime

Basic level categorization

Accuracy (% correct)

Stimulus Duration (milliseconds)

Chance is 50%

All the data I will show you today

Typical primate fixation duration during natural viewing

**Let's try to define a domain of behavior so that we can gauge/make progress in prediction.**
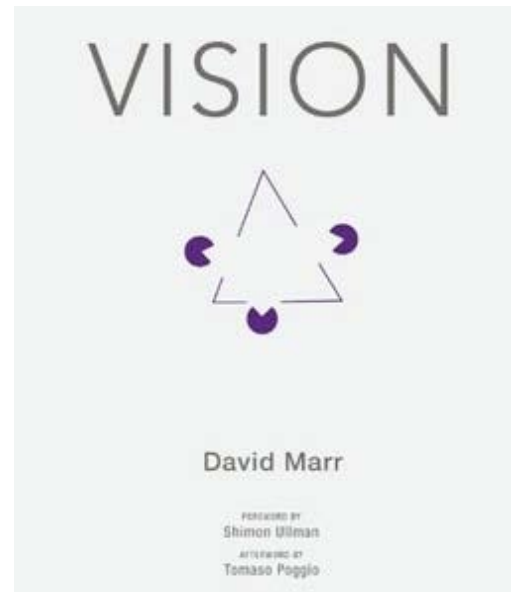
**Vision**

Object recognition

*"Core object recognition"*

| Computational theory | Representation and algorithm | Hardware implementation |
|---|---|---|
| What is the goal of the computation, why is it appropriate, and what is the logic of the strategy by which it can be carried out? | How can this computational theory be implemented? In particular, what is the representation for the input and output, and what is the algorithm for the transformation? | How can the representation and algorithm be realized physically? |

VISION

David Marr

FOREWORD BY
Shimon Ullman

AFTERWORD BY
Tomaso Poggio

David Courtnay Marr
(1946-1980)

Marr, 1982

# Reaching a common language

| | Comp vision, Machine learning | Neuroscience, Cognitive Science |
|---|---|---|
| 1. *What is the problem we are trying to solve?* | Benchmarks Brain solves "it" | "Perception" Behavior Psychophysics |
| 2. *What do good solutions look like?* | Useful image representations ("features") | Explicit neuronal population spiking patterns |
| 3. *How do we instantiate these solutions?* | Algorithms, mechanisms | Neuronal wiring / weighting patterns |
| 4. *How do we construct those instantiations?* | Learning rules, initial conditions, training images | Plasticity, architecture, experience |

**"Identity preserving image variation"**

**View: position, size, pose, illumination**

**Clutter, occlusion**

Pinto, Nicolas, David D. Cox, and James J. Di Carlo. "Why is real-world visual object recognition hard?"
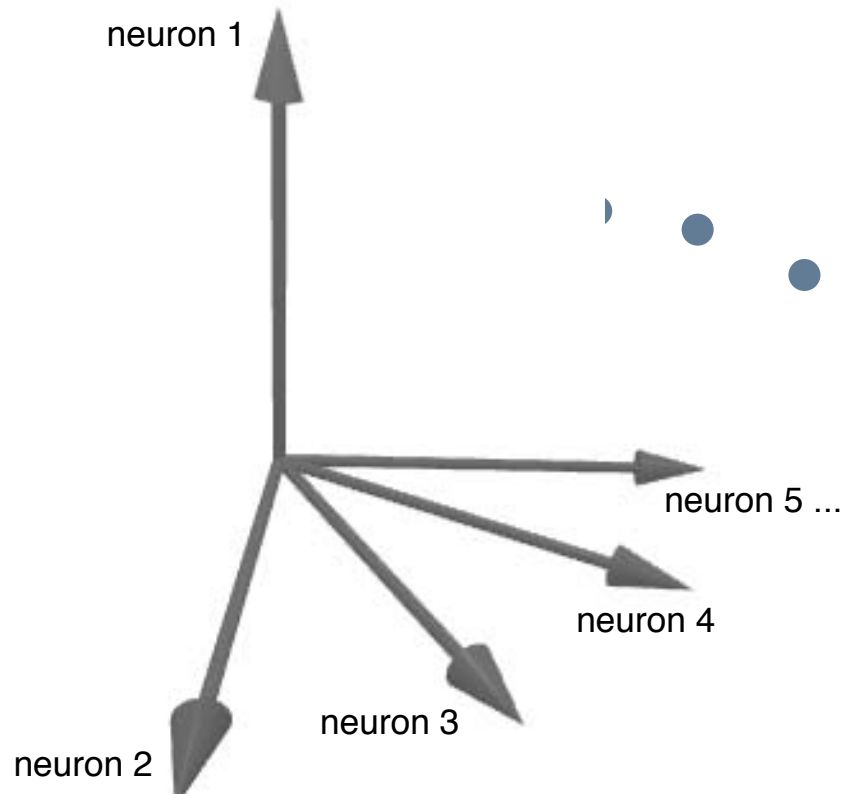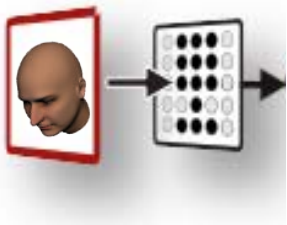PLoS Comput Biol 4, no. 1 (2008): e27. doi: 10.1371/journal.pcbi.0040027. License CC BY.

**subordinate
level variation**

*Poggio, Ullman, Grossberg, Edleman, Biederman, etc.*
*DiCarlo and Cox, **TICS** (2007), Pinto, Cox, and DiCarlo, **PLoS Comp Bio** (2008),*
*DiCarlo, Zoccolan and Rust, **Neuron** (2012)*

# The brain's "camera" represents the image as populations of visually-evoked "features"

"Joe's" identity manifold



neuron 1

neuron 5 ...

neuron 4

neuron 3

neuron 2

"Joe"

"Joe"

pixel    RGC

# The computational crux of object and face recognition

**A "good" set of visual features**

**== "Explicit" representation of object shape**

**We assume: "shape" maps to "identity" and "category"**

individual 2 ("Joe")

"Joe"

**Neural population**

**Should be able to find it with low\* number of training examples**

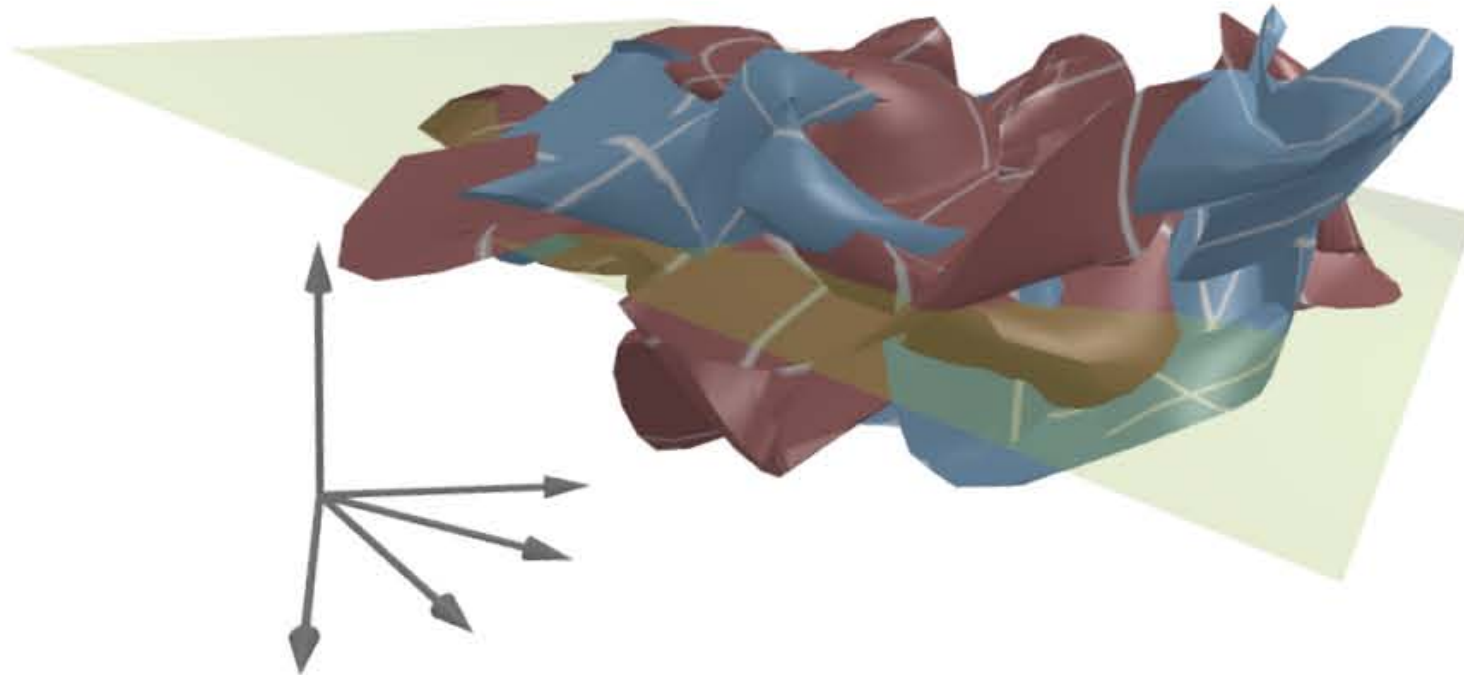separating hyperplane

*linear classifier* ≈ *downstream neuron(s)*

"not Joe"  individual 1 ("Sam")



*DiCarlo and Cox, **TICS** (2007)*

## Pixel population representation
(~ retinal image representation)
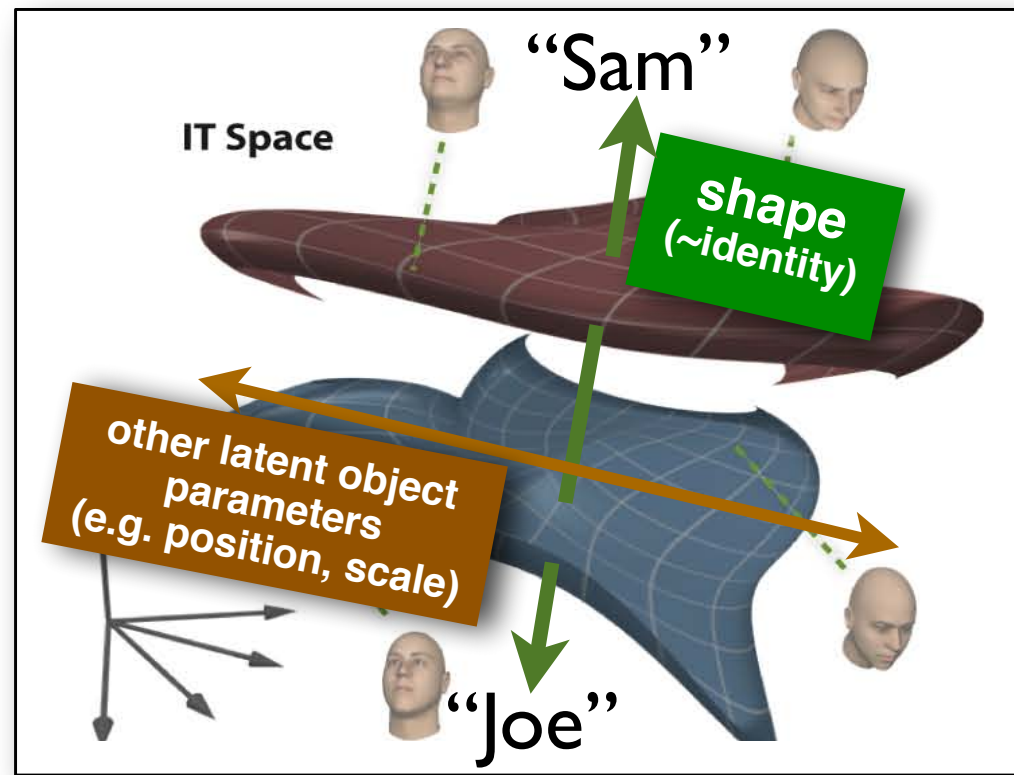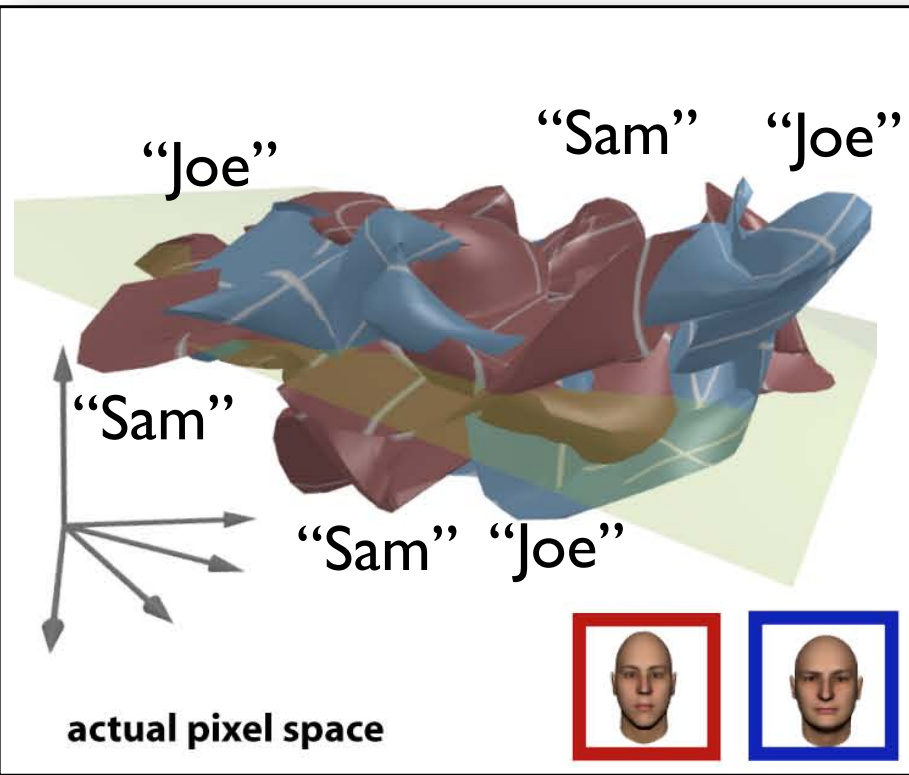


individual 2

ineffective
separating
hyperplane

individual 1

**object manifolds are "tangled"**

*(Due to identity-preserving image <u>variation</u>.)*

*DiCarlo and Cox, **TICS** (2007);  Pinto, Cox, and DiCarlo, **PLoS Comp Bio** (2008)*

**Tangled, implicit object information**

*Transformation* ⟶

*This must be non-linear*

**Untangled, explicit object information**
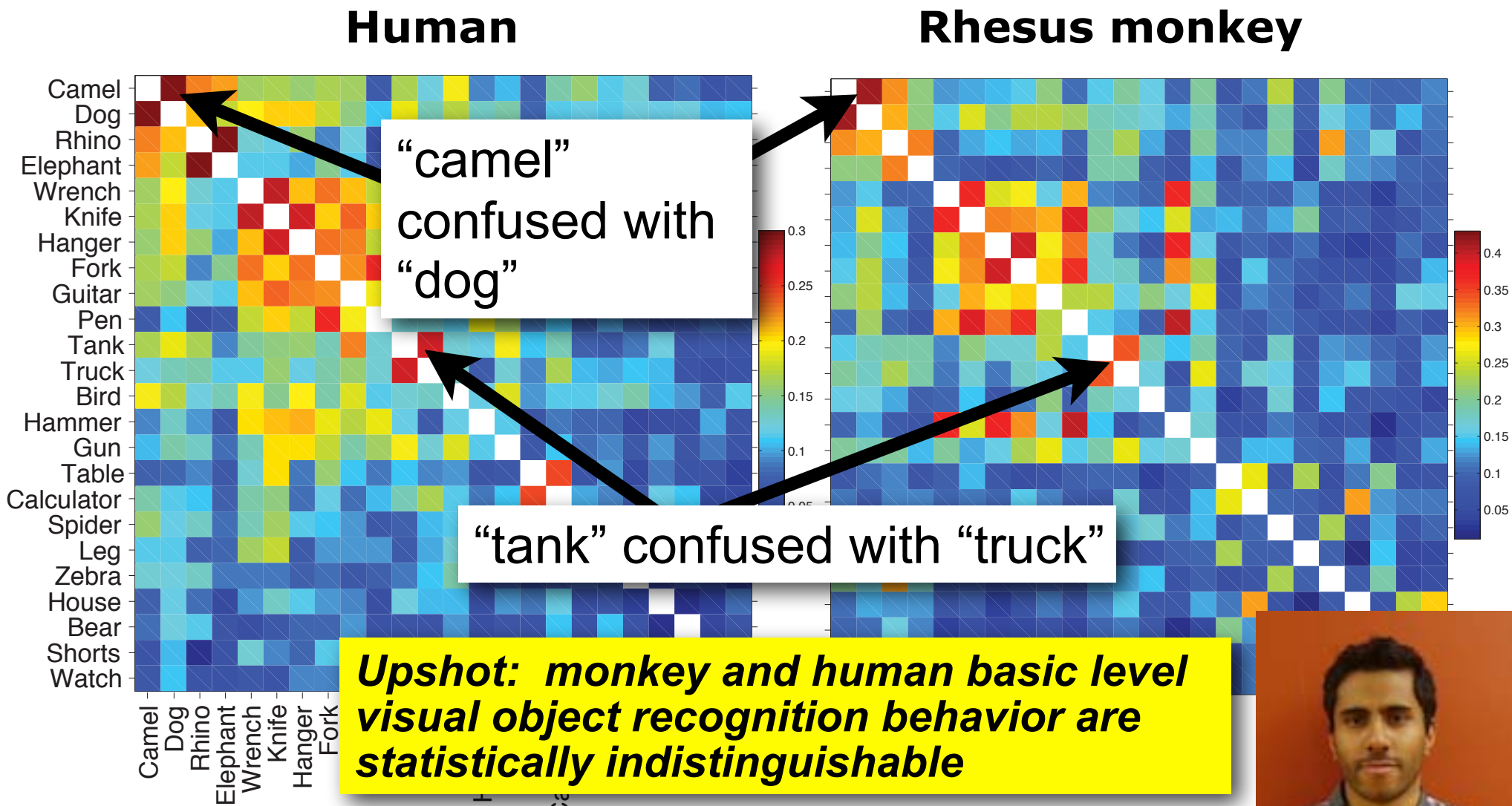
*a poor encoding basis (for this task)*

*a powerful encoding basis somewhere in the brain*

Courtesy of Elsevier, Inc., http://www.sciencedirect.com. Used with permission.
Source: DiCarlo, James J., and David D. Cox. "Untangling invariant object recognition. "Trends in cognitive sciences 11, no. 8 (2007): 333-341;
https://doi.org/10.1016/j.tics.2007.06.010.

# The ventral visual stream



Human · Rhesus monkey

"camel" confused with "dog"

"tank" confused with "truck"

**Upshot: monkey and human basic level visual object recognition behavior are statistically indistinguishable**

*Comparison of Object Recognition Behavior in Human and Monkey*
*R. Rajalingham, K Schmidt, J.J. DiCarlo,* **Vision Sciences Society** *(2014)*
*R. Rajalingham, K Schmidt, J.J. DiCarlo,* **J. Neuroscience** *(in press)*

Adapted from Motter and Mountcastle 1981

***Decision and action***

**V1** **V4**

**V2** **IT**

***Memory***

**Ventral visual stream**

*We think we know where the neural mechanisms and resulting representations that solve core object recognition live in the primate brain.*

*We can measure and manipulate those representations at the level of neuronal spikes.*

Adapted from Motter and Mountcastle 1981

# The ventral visual stream

**IT is believed to be that powerful encoding basis**

**Ventral visual stream**

*Retinotopic map* · *Retinotopic map* · *Retinotopic map* · *Retinotopic map* · *Retinotopic map* · *non-retinotopic*

pixel · RGC · LGN · V1 · V2 · V4 · IT

**Key concept: each area conveys a new neural population representation**

**"IT" (Inferior temporal cortex)**

**A**

V2
V4
V1
PIT
CIT
AIT
Retina
LGN

**B**

| | Latency |
|---|---|
| ~10 M (IT representation) | ~100 ms |
| STP$_a$  AIT ~16 M | |
| STP$_p$  CIT ~17 M | ~90 ms |
| 7a | |
| PIT ~36 M | ~80 ms |
| LIP MST FST | |
| DP  VOT | |
| ~15 M (V4 representation) | ~70 ms |
| MIP PO MT  V4 ~68 M | |
| PIP  V3A | |
| V3  ~29 M (V2 representation) | ~60 ms |
| V2 ~150 M | |
| ~37 M (V1 representation) | ~50 ms |
| V1 ~190 M | |
| LGN ~1 M (LGN representation) | ~40 ms |
| Retina ~1 M (RCG representation) | |

*Adapted from DiCarlo et al. 2012*

29

Courtesy of Elsevier, Inc., http://www.sciencedirect.com. Used with permission.
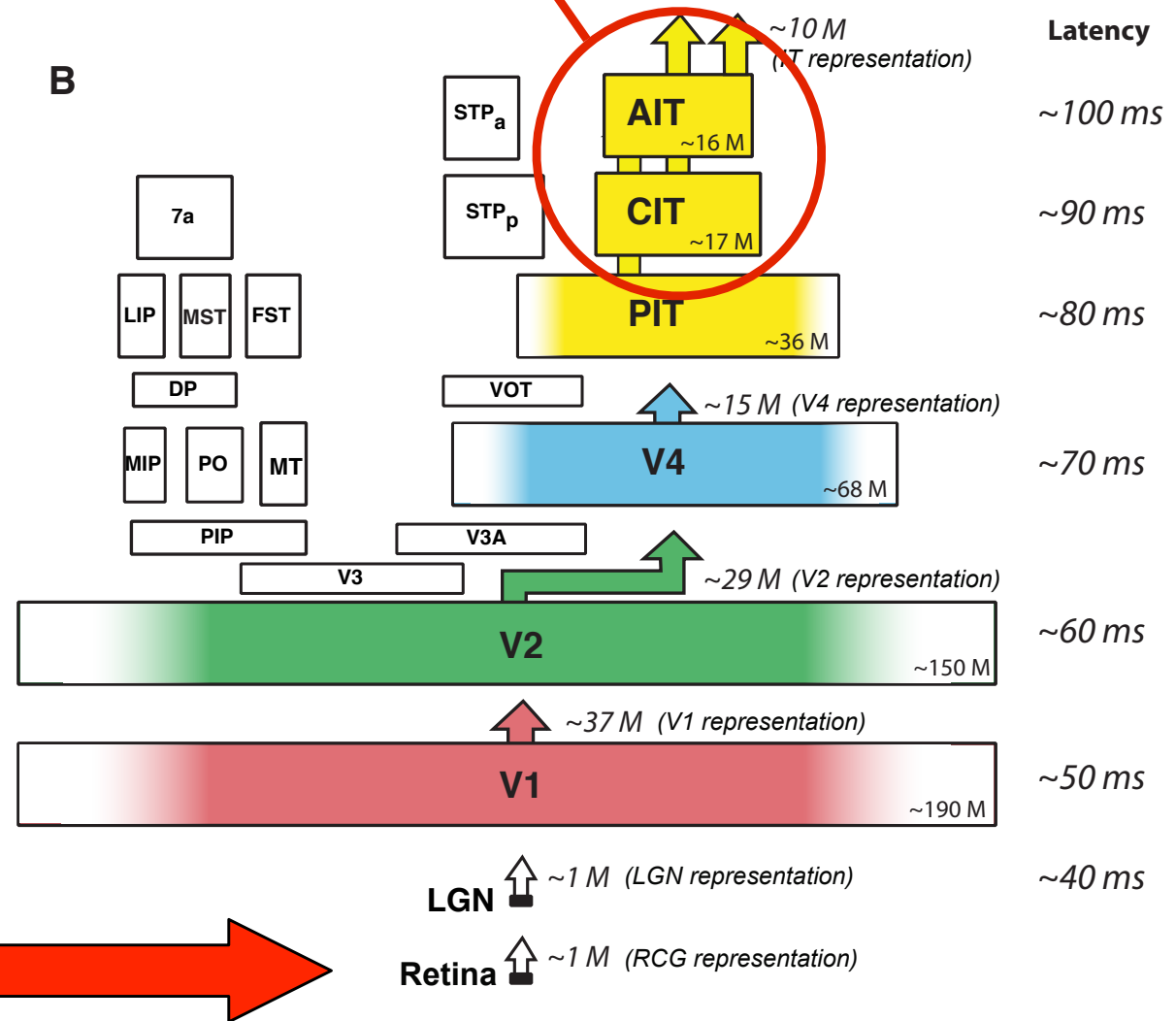Source: DiCarlo, James J., Davide Zoccolan, and Nicole C. Rust. "How does the brain solve visual object recognition?" Neuron 73, no. 3 (2012): 415-434.

*Adapted from DiCarlo et al. 2012*

# Retinal ganglion cell RF structure:

## A Receptive fields of concentric cells of retina and lateral geniculate nucleus

On-center      Off-center

3
Central
illumination

4
Surround
illumination

Source: Siegelbaum, Steven A., and A. James Hudspeth. Principles of neural science. Eds. Eric R.
Kandel, James H. Schwartz, and ThomasM. Jessell. Vol. 4. New York: McGraw-hill, 2000.

*Adapted from Hubel*          *Adapted from Kandel , Schwartz and Jessell*

**"IT" (Inferior temporal cortex)**

You are here.

B

| Region | Latency |
|---|---|
| AIT ~16 M, ~10 M (IT representation) | ~100 ms |
| CIT ~17 M | ~90 ms |
| PIT ~36 M | ~80 ms |
| V4 ~68 M, ~15 M (V4 representation) | ~70 ms |
| V2 ~150 M, ~29 M (V2 representation) | ~60 ms |
| V1 ~190 M, ~37 M (V1 representation) | ~50 ms |
| LGN ~1 M (LGN representation) | ~40 ms |
| Retina ~1 M (RCG representation) | |

Other labeled regions: STPa, STPp, 7a, LIP, MST, FST, DP, VOT, MIP, PO, MT, PIP, V3A, V3

Courtesy of Elsevier, Inc., http://www.sciencedirect.com. Used with permission.
Source: DiCarlo, James J., Davide Zoccolan, and Nicole C. Rust. "How does the brain solve visual object recognition?" Neuron 73, no. 3 (2012): 415-434.

*Adapted from DiCarlo et al. 2012*
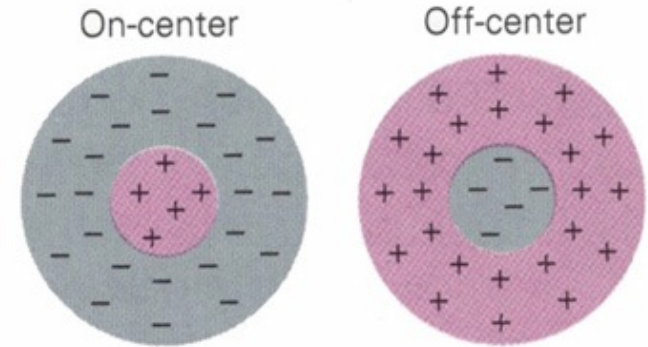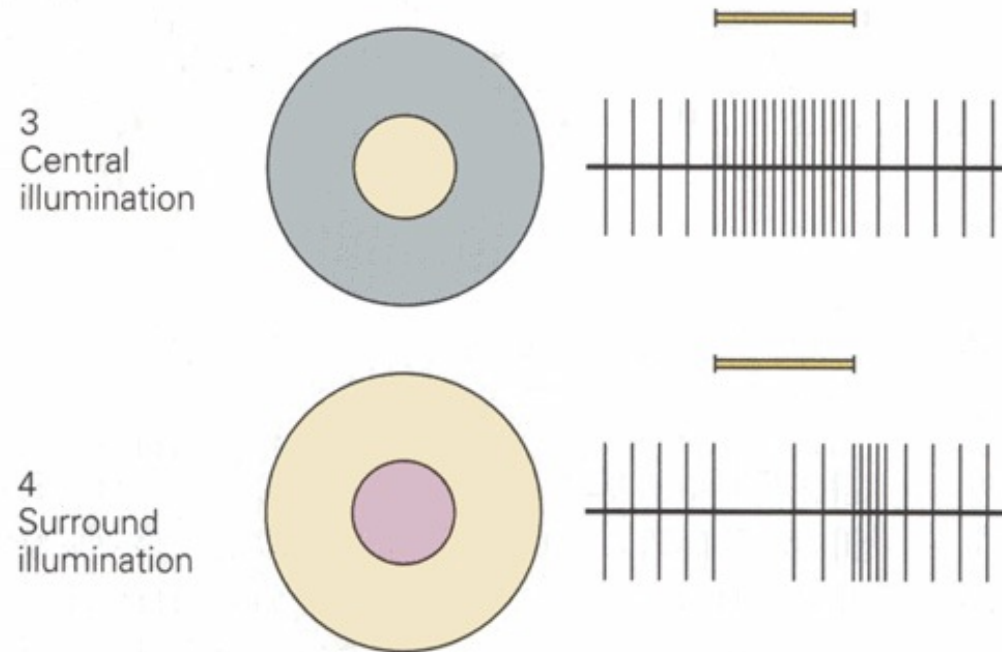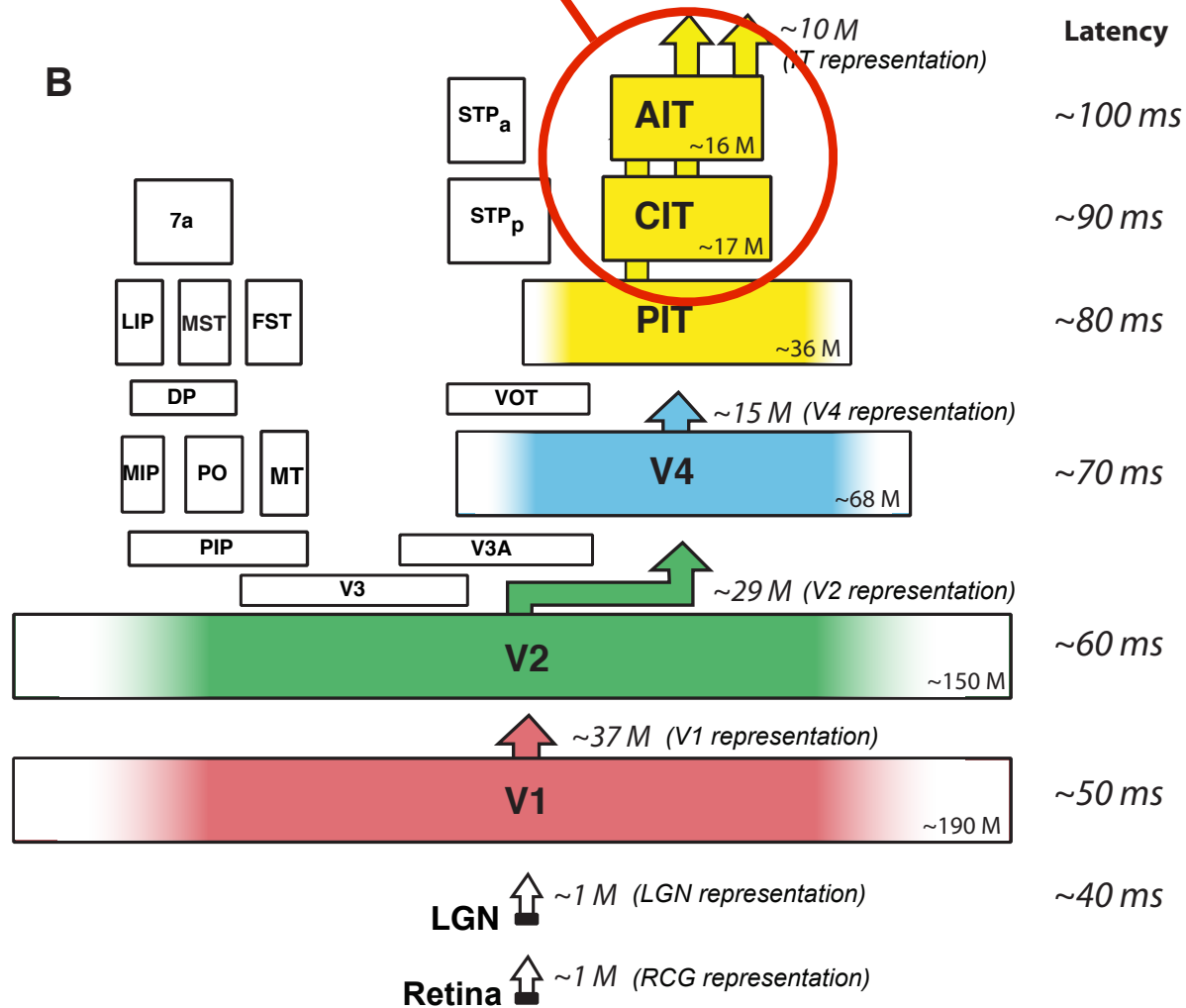
# Primary visual cortex (Area V1):

Figure removed due to copyright restrictions. Please see the video.
Source: Eye, Brain, and Vision. David H. Hubel. New York : Scientific American
Library: Distributed by W.H. Freeman, c1988. ISBN: 0716750201.

Orientation
selectivity

Orientation
selectivity with
some position
tolerance

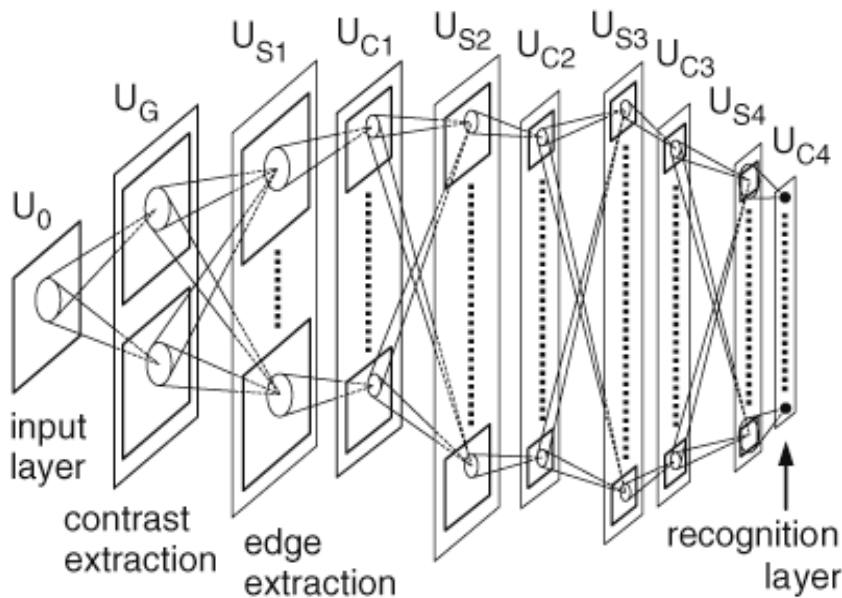*Adapted from Kandel , Schwartz and Jessell*

# Brain-inspired computer algorithms

## Examples:

- *Hubel & Wiesel (1962)*

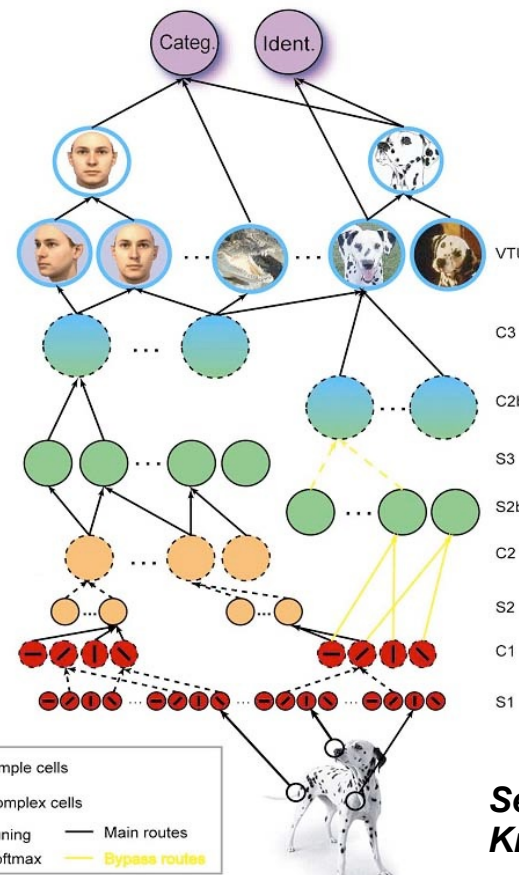Figure removed due to copyright restrictions.
Please see the video. Source: Eye, Brain, and Vision. David H. Hubel. New York: Scientific
American Library: Distributed by W.H. Freeman, c1988. ISBN: 0716750201.

**FROM BIOLOGY:**

- *Hierarchy*
- *Spatially local filters*
- *Convolution*
- *Normalization*
- *Threshold NL*
- *Unsupervised learning*
- *...*

*Serre, Kouh, Cadieu, Knoblich,*
*Kreiman & Poggio 2005*

**"IT" (Inferior temporal cortex)**

Courtesy of Elsevier, Inc., http://www.sciencedirect.com. Used with permission.
Source: DiCarlo, James J., Davide Zoccolan, and Nicole C. Rust. "How does the brain solve visual object recognition?" Neuron 73, no. 3 (2012): 415-434.

*Adapted from DiCarlo et al. 2012*

# Area V2 (first cortical area after V1):
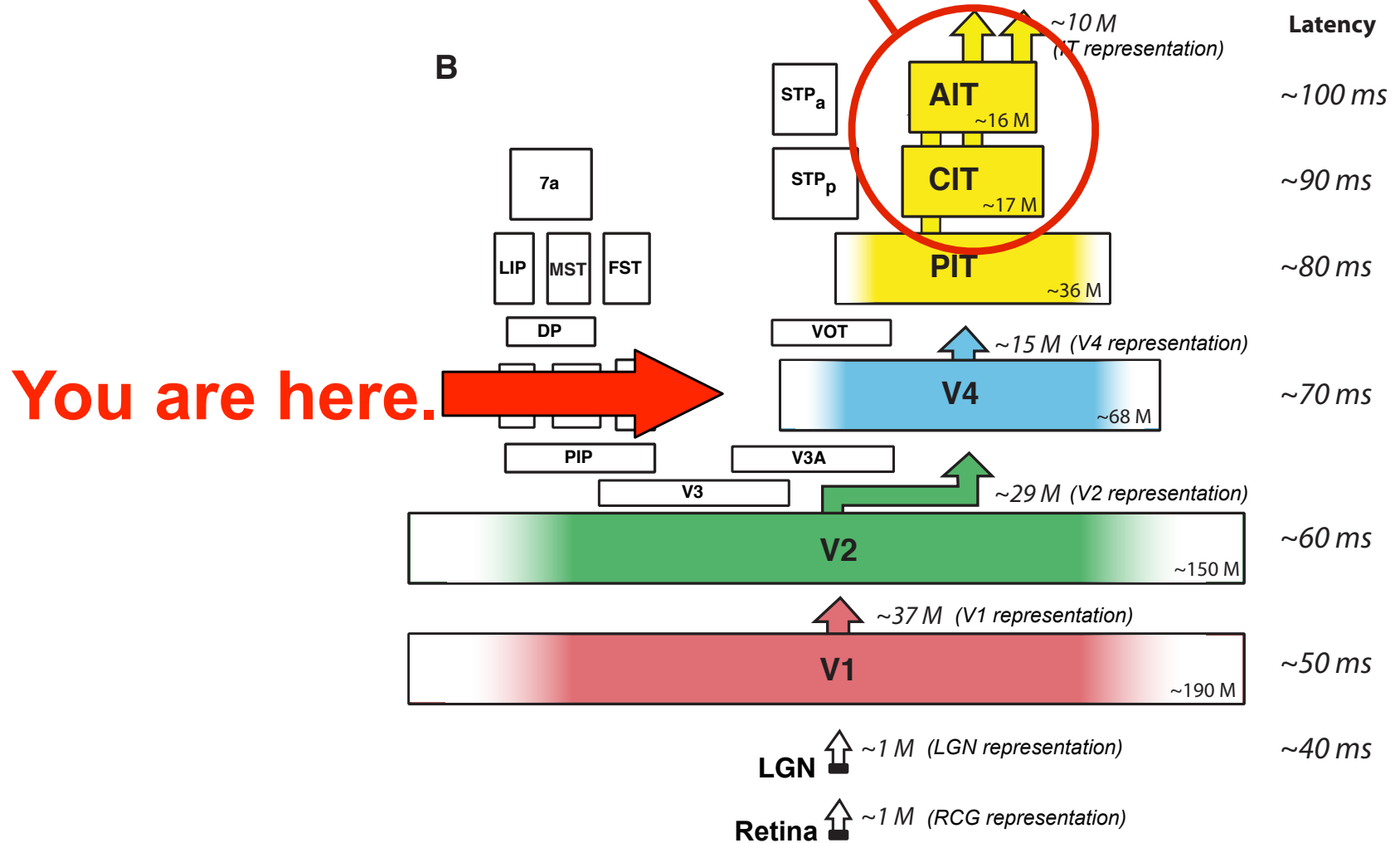


Original photographs

V1-like filters matched: spectrally matched noise

Correlations matched: naturalistic texture

**b**

V1
n = 102

V2
n = 103

Naturalistic

Noise

Null

Normalized firing rate

Time from stimulus onset (ms)

**Interpretation:**

- **V2 neurons apply "and-like" operators on V1 outputs**

- **those "ands" are tuned toward natural co-occurring V1 statistics**

Reprinted by permission from Macmillan Publishers Ltd: Nature Neuroscience.
Source: Freeman, Jeremy, Corey M. Ziemba, David J. Heeger, Eero P. Simoncelli, and J. Anthony Movshon. "A functional and perceptual signature of the second visual area in primates. "Nature neuroscience 16, no. 7 (2013): 974-981.

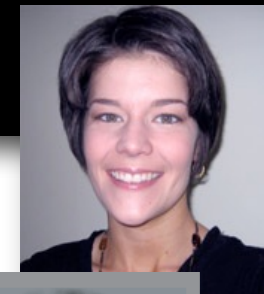*Adapted from Freeman, Ziemba, Heeger, Simoncelli, & Movshon, Nature Neuro (2013)*

36

**"IT"** (Inferior temporal cortex)

**You are here.**

Latency

~10 M (IT representation) — ~100 ms

~90 ms

~80 ms

~15 M (V4 representation) — ~70 ms

~29 M (V2 representation) — ~60 ms

~37 M (V1 representation) — ~50 ms

~1 M (LGN representation) — ~40 ms

~1 M (RCG representation)

*Adapted from DiCarlo et al. 2012*

# What is V4 doing?

**Increased selectivity for conjunction of features that tend to co-occur in natural images**

Same animal, task, stimuli.

*Rust & DiCarlo **J Neurosci** (2010)*

*Rust & DiCarlo **J Neurosci** (2012)*

**Easier to read-out object identity in IT**
*(per neuron, matched for information)*

**Tangled, implicit object information**

**Explicit, untangled object representation**

?

pixel    RGC    LGN    V1    V2    V4    IT

DOG   S(·)

V4 Responses to Non-Cartesian Gratings
Gallant et al. 1996

Courtesy of Journal of Neurophysiology. Used with permission.
Source: Gallant, Jack L., Charles E. Connor, Subrata Rakshit, James W. Lewis, and DAVID C. Van Essen. "Neural responses to polar, hyperbolic, and Cartesian gratings in area V4 of the macaque monkey." Journal of neurophysiology 76, no. 4 (1996): 2718-2739.

# What shape features drive V4 responses?

*Adapted from C.E. Connor*

Make a basis for shapes:
each shape = set of curved elements
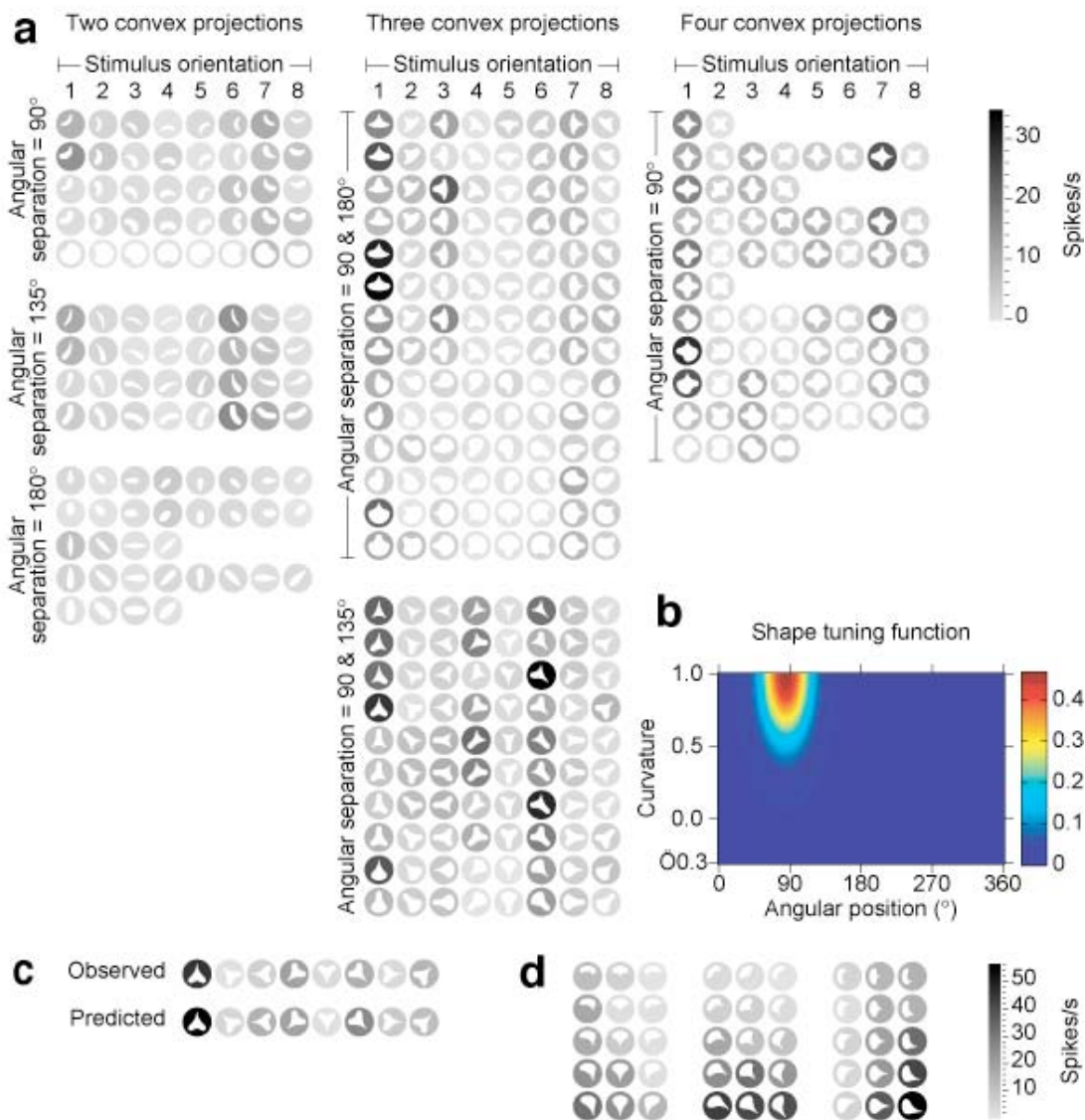each element = (ang position, curvature)

Hypothesis:
V4 neurons are tuned in this basis

Figure removed due to copyright restrictions. Please see the video.
Source: "Shapes Dimensions and Object Primitives" from Chalupa,
Leo M., and John Simon Werner. The visual neurosciences. [Vol. 2].
MIT Press, 2004. Harvard.

*Adapted from C.E. Connor*

Make a basis for shapes:
each shape = set of curved elements
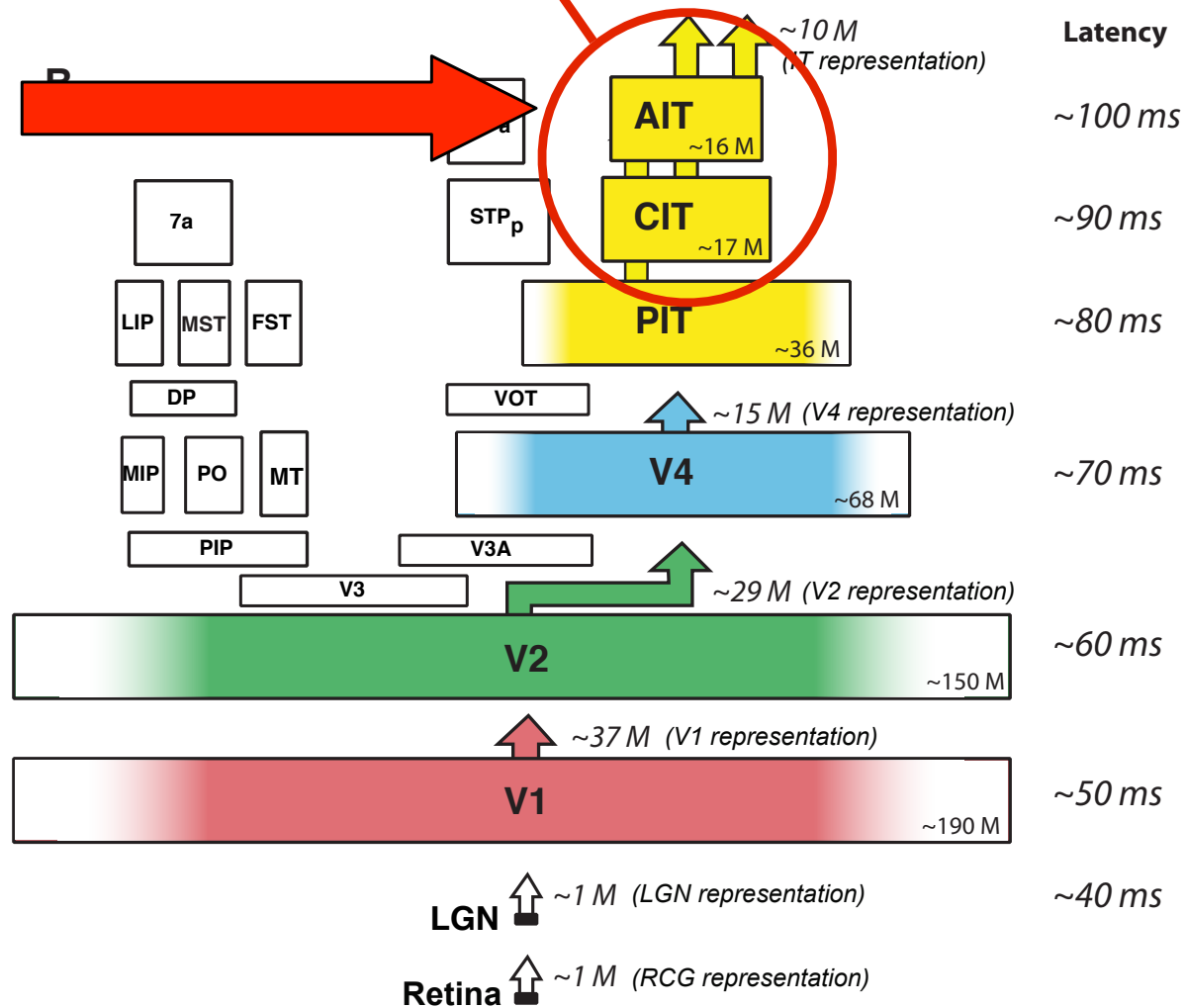each element = (ang position, curvature)

Hypothesis:
V4 neurons are tuned in this basis



Experimental result:
Hypothesis explains ~50% of the explainable response variance

*Pasupathy and Connor (V4)*
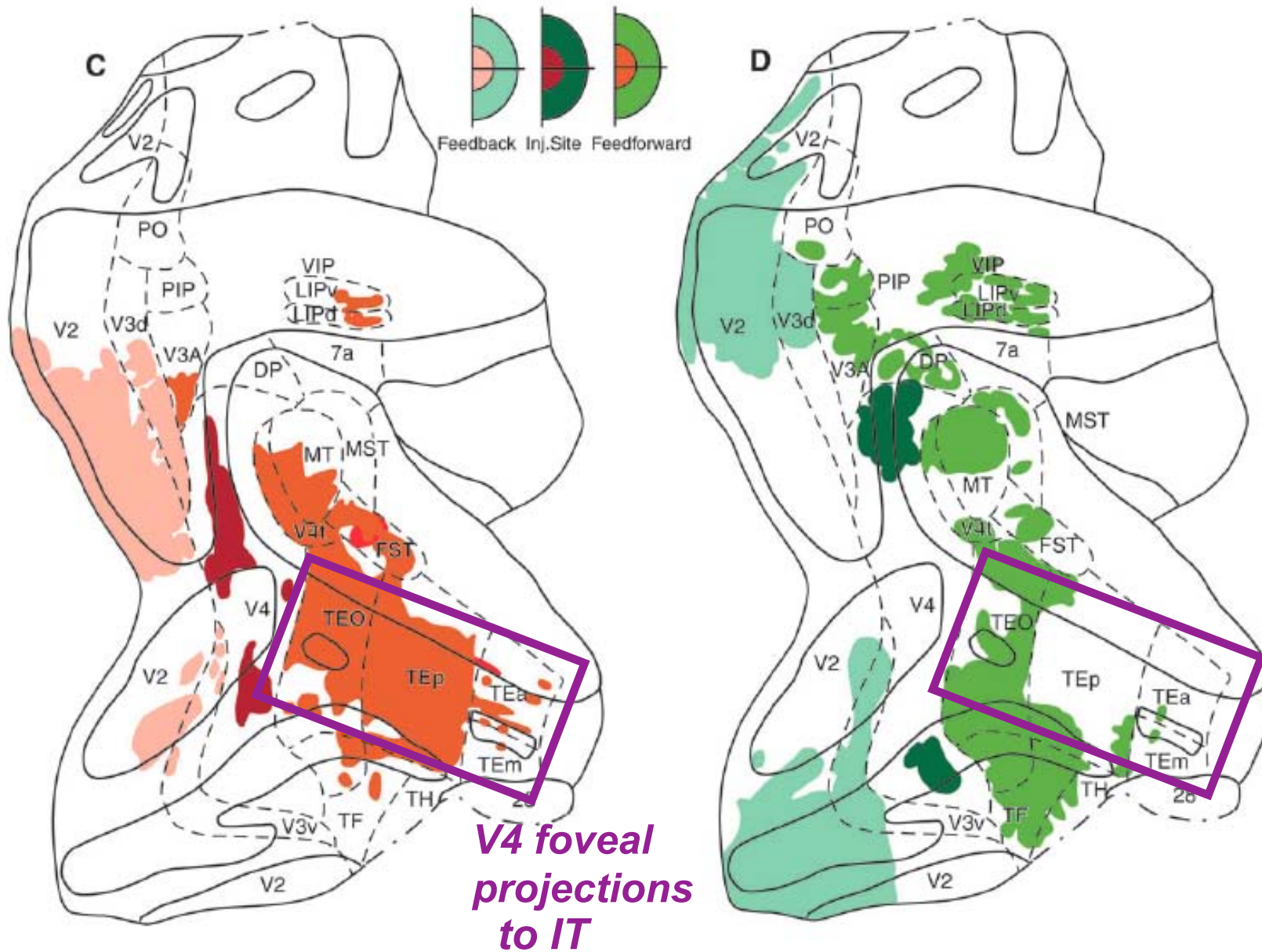*Brincat and Connor (PIT)*

Reprinted by permission from Macmillan Publishers Ltd: Nature Neuroscience.
Source: Pasupathy, Anitha, and Charles E. Connor. "Population coding of shape in area V4." Nature neuroscience 5, no. 12 (2002): 1332-1338.

41

**"IT"** (Inferior temporal cortex)

**You are here.**

| Region | | Latency |
|---|---|---|
| | ~10 M (IT representation) | |
| AIT ~16 M | | ~100 ms |
| CIT ~17 M | STP_p | ~90 ms |
| PIT ~36 M | | ~80 ms |
| VOT | ~15 M (V4 representation) | |
| V4 ~68 M | | ~70 ms |
| V3A | | |
| V3 | ~29 M (V2 representation) | |
| V2 ~150 M | | ~60 ms |
| | ~37 M (V1 representation) | |
| V1 ~190 M | | ~50 ms |
| LGN ~1 M (LGN representation) | | ~40 ms |
| Retina ~1 M (RCG representation) | | |

Other labeled regions: 7a, LIP, MST, FST, DP, MIP, PO, MT, PIP

*Adapted from DiCarlo et al. 2012*

# IT is about central vision



*V4 foveal projections to IT*

Source: Ungerleider, Leslie G., Thelma W. Galkin, Robert Desimone, and Ricardo Gattass. "Cortical connections of area V4 in the macaque." Cerebral Cortex 18, no. 3 (2008): 477-499.

*Ungerleider, L. G. et al. Cereb. Cortex 2007*  43

# Stimulus selectivity in inferotemporal cortex
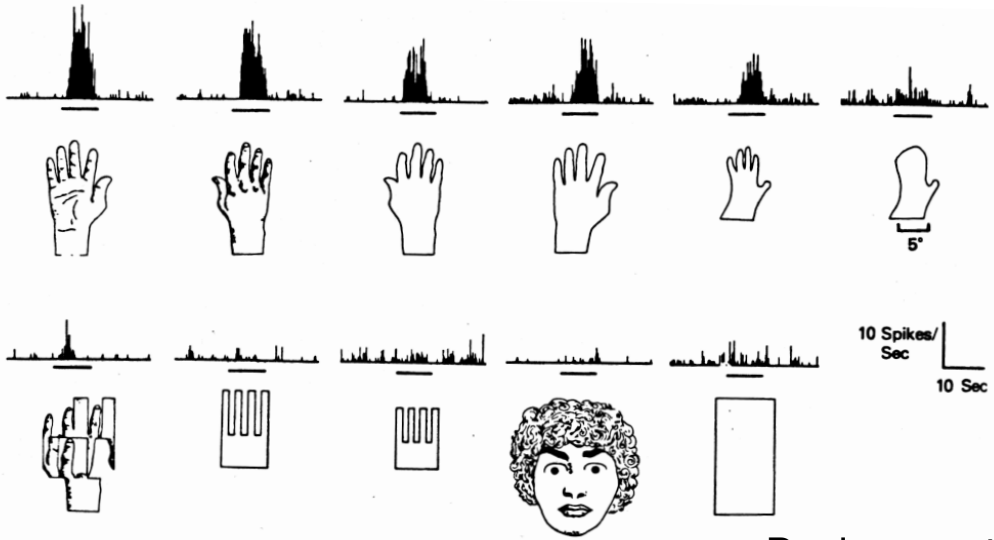## Gross, Rocha-Miranda & Bender 1972

Figure removed due to copyright restrictions. Please see the video.
Source: Gross, Charles G., Carlos Eduardo de Rocha-Miranda, and
David B. Bender. "Visual properties of neurons in inferotemporal cortex
of the Macaque." Journal of neurophysiology 35, no. 1 (1972): 96-111.

*The use of [these] stimuli was begun one day when, having failed to drive a unit with any light stimulus, we waved a hand at the stimulus screen and elicited a very vigorous response from the previously unresponsive neuron...*

*We then spent the next 12 hr testing various paper cutouts in an attempt to find the trigger feature for this unit. When the entire set of stimuli used were ranked according to the strength of the response that they produced, we could not find a simple physical dimension that correlated with this rank order. However, the rank order of adequate stimuli did correlate with similarity (for us) to the shadow of a monkey hand" (Gross et al., 1972).*
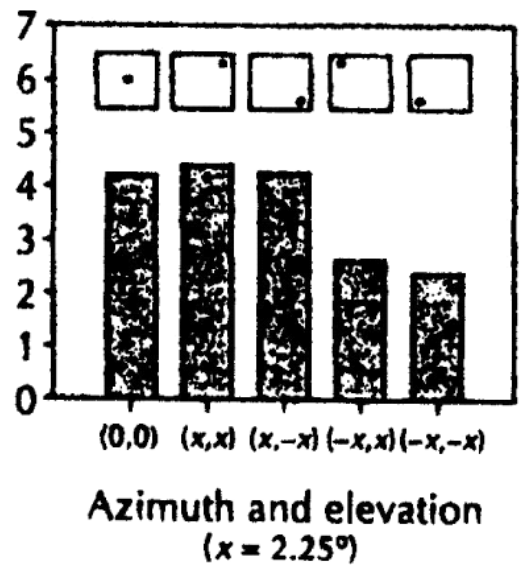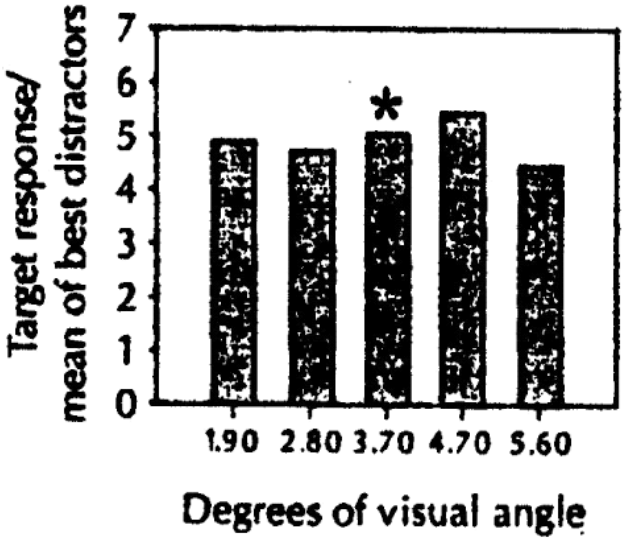
## The ventral stream and object recognition



IT neurons can be tuned to specific combinations of features (high "selectivity")

Desimone et al. (1984)

That selectivity is tolerant to changes in position and size

Logothetis et al. (1995)

# Primary visual cortex:

Orientation
selectivity

Figure removed due to copyright restrictions. Please see the video.
Source: Eye, Brain, and Vision. David H. Hubel. New York : Scientific American
Library: Distributed by W.H. Freeman, c1988. ISBN: 0716750201.

Orientation
selectivity with
some position
tolerance

*Adapted from Kandel , Schwartz and Jessell*

# What stimulus feature are IT neurons actually "tuned" to?

Figure removed due to copyright restrictions. Please see the video.
Source: Tanaka, Keiji. "Neuronal mechanisms of object recognition."
Science-New York Then Washington 262 (1993): 685-685.

Figure removed due to copyright restrictions. Please see the video.
Source: Tanaka, Keiji. "Columns for complex visual object features in
the inferotemporal cortex: Clustering of cells with similar but slightly
different stimulus selectivities." Cerebral cortex 13, no. 1 (2003): 90-99.
doi: 10.1093/cercor/13.1.90.

# IT has spatial organization at 500 um - 1 mm scale

Figure removed due to copyright restrictions. Please see the video.
Source: Tanaka, Keiji. "Columns for complex visual object features
in the inferotemporal cortex: Clustering of cells with similar but slightly
different stimulus selectivities." Cerebral cortex 13, no. 1 (2003): 90-99.

Figure removed due to copyright restrictions. Please see the video.
Source: Tanaka, Keiji. "Columns for complex visual object features in
the inferotemporal cortex: Clustering of cells with similar but slightly
different stimulus selectivities." Cerebral cortex 13, no. 1 (2003): 90-99.
doi: 10.1093/cercor/13.1.90.

# Larger scale (2-6 mm) organization for some image contrasts

ML



Figure removed due to copyright restrictions. Please see the video.

Tsao, Freiwald, and Livingstone used fMRI to reveal a set of face selective regions in IT (aka "face patches")
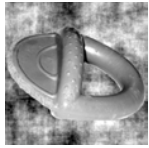
Most of the single neurons in these regions showed a preference for frontal faces

Tsao et al., *Science* 2006

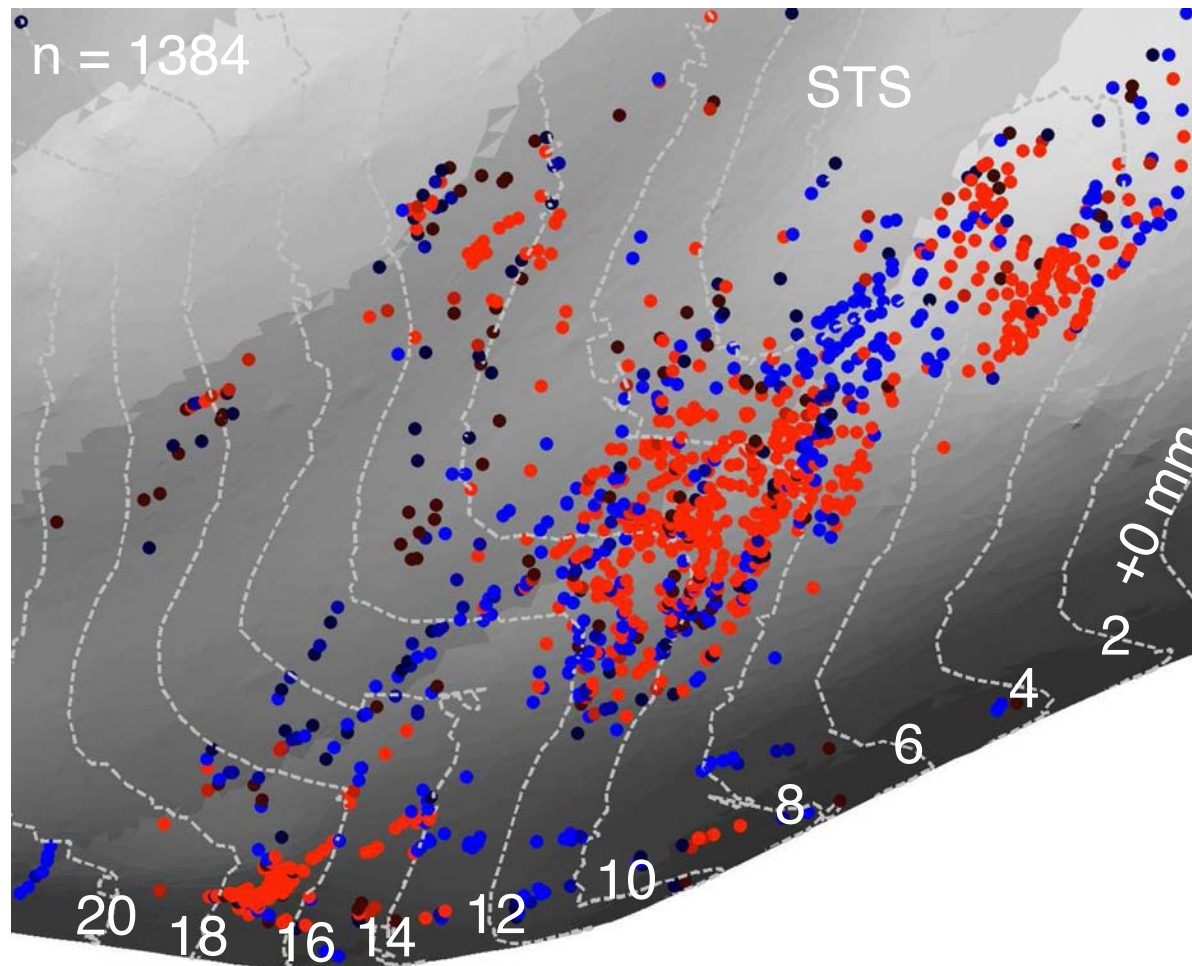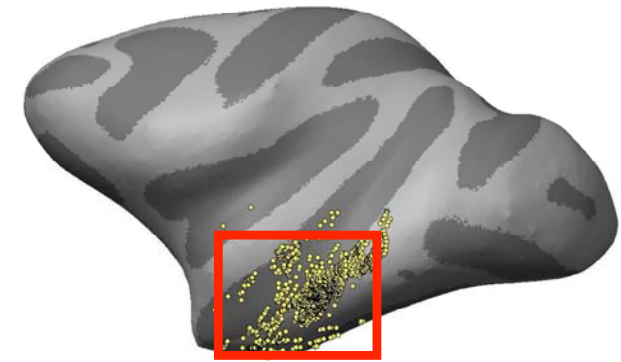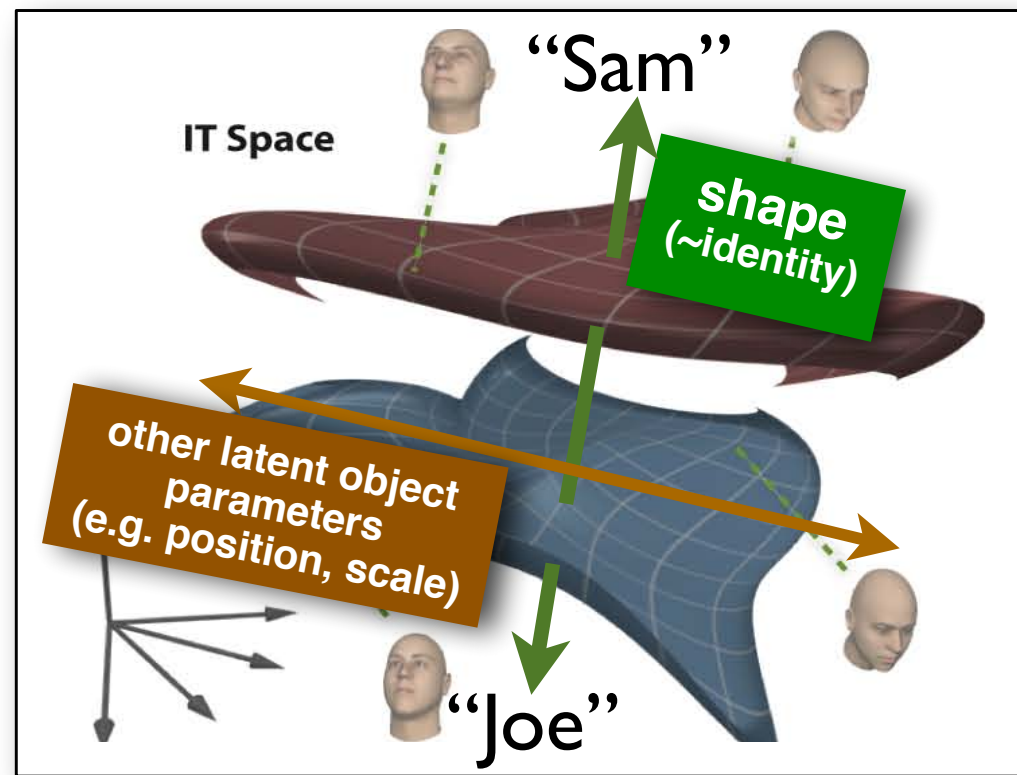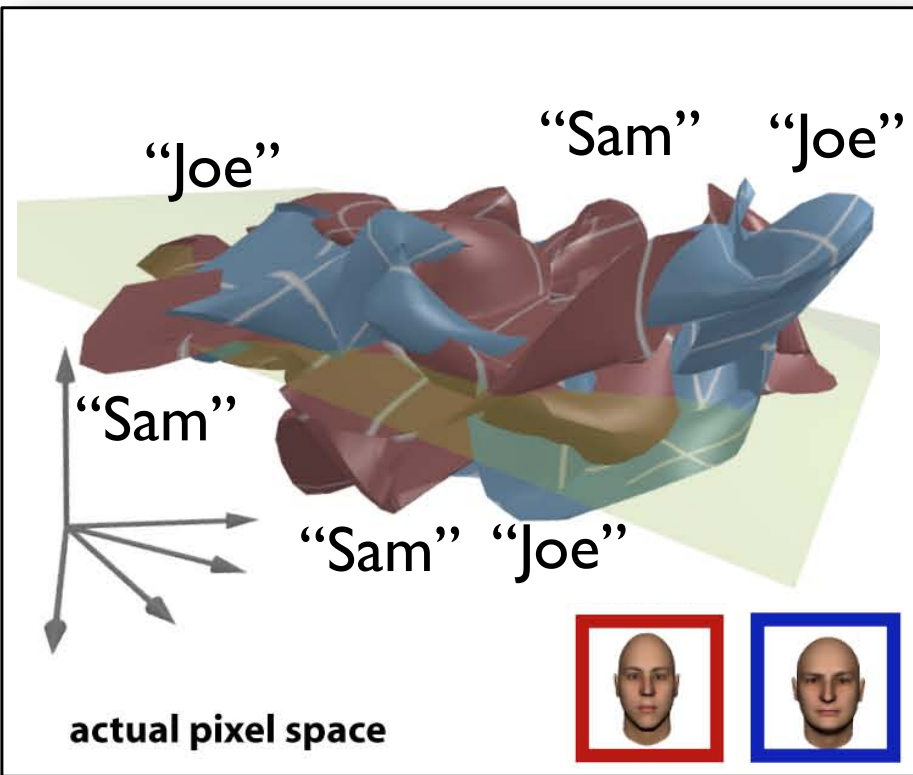# IT selectivity is particularly clustered for some image contrasts

vs

**face objects**

**non-face objects**

MUA

n = 1384

STS

+0 mm

2

4

6

8

10

20  18  16  14  12

50

"Joe"   "Sam"   "Joe"

"Sam"

"Sam"   "Joe"

actual pixel space

IT Space

**shape (~identity)**

**other latent object parameters (e.g. position, scale)**

"Sam"

"Joe"

**Tangled, implicit object information**

*Transformation* ⟶

**Untangled, explicit object information**

*a poor encoding basis (for this task)*

*a powerful encoding basis somewhere in the brain*

pixel   RGC   LGN   V1   V2   V4   IT

T(·)   T(·)   T(·)   T(·)   T(·)   T(·)
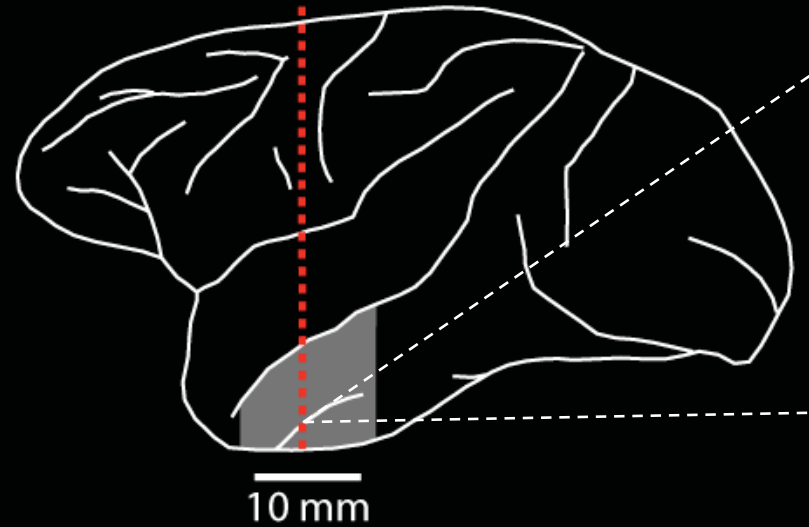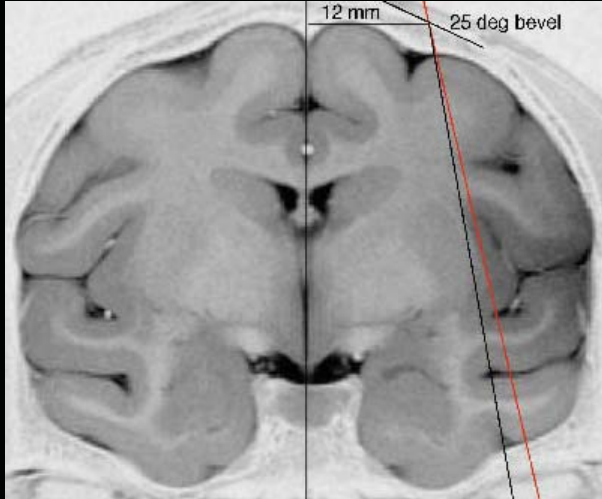
51

# Example spiking activity in IT

Figure removed due to copyright restrictions. Please see the video. Source: Eye, Brain, and Vision. David H. Hubel. New York: Scientific by W.H. Freeman, c1988. ISBN: 0716750201.

Site 1

0  100
ms

# An early test of the IT population

**A broad set of 78 test objects from eight categories …**

For each, test changes
in position and scale

*0.5x*

*2 deg*  *4 deg*

*2x*

100 ms  100 ms  100 ms

- *fixation task*
- *15 images per trial*
- *10 repetitions per image*
- *randomized and counter-balanced*

*time* →  100 ms

*Hung*, Kreiman*, Poggio and DiCarlo, **Science** (2005)*

# The "mean" IT population

(n ~ 350 IT sites)

Recording Site

63

1

Object image

1

78

# How do we test if the population image is "good"?

**Implicit representation**

neuron 2

neuron 1

*"inaccessible"*
*object information*

*BAD*

**Explicit representation**

neuron 2

NOT object "A"

Object "A"

neuron 1

*Linearly separable*

*"accessible"*
*object information*

*GOOD*

# How explicit ("good") is object information in IT?



Hung*, Kreiman*, Poggio and DiCarlo, *Science* (2005)

**Explicit object information in IT ?**

Does not work in earlier visual areas
e.g. V1 vs. IT  or  V4 vs. IT

• Consistent with other IT work
(e.g. Rolls, Tanaka, Miyashita, Yamane, Sugase, Logothetis, Vogels, Connor, ...)

*Rapid, explicit object representation in IT*

Hung*, Kreiman*, Poggio and DiCarlo, **Science** (2005)

57

**Summary so far:**

the problem of visual object recognition

a tour of the ventral stream

IT population seems to have solved a key problem

Over the last 40 years. we (the field) have largely described important phenomenology

Next phase of this field:  developing and testing predictive models

*Images*

**Behavioral reports /
perception ("mind")**

"clock"

"cat"

"car"

"dog"

"face"

**Neural activity**

**???**

**Decoding
algorithm ?**

*e.g. spiking pattern of
a neural population*

**"Neural representation"**

**Goal is accurate
predictivity**

**(Domain: core object recognition)**

# Goal:  end-to-end understanding

**1. Can we infer the precise <span style="color:red">decoding</span> mechanism(s) that the brain uses to support perceptual reports about visually presented objects?**

**2. Can we infer the <span style="color:red">encoding</span> mechanism(s) that accurately predicts the <span style="color:yellow">relevant</span> ventral stream population patterns of neural activity from each image?**

**Images**

**Behavioral reports ("perception")**

**Specific task domain**

*(nouns)*

"clock"

"cat"

"car"

"dog"

"face"

**Generative image domain**
*(single foreground object)*

**???**

**Neural activity**

*a specific spiking pattern over the IT neural population in response to a specific image*

**???**

**"IT Neural representation"**

61

# 3-d object Models
## (e.g. "car")

# experimenter-chosen view parameters

+

Position

Size

Pose

# ray-trace render

# place on a randomly-chosen background image

- generative space of images, each with a single foreground object and experimenter-known viewing parameters.

- uncorrelated, new background every image
  ==> challenging for computer vision, doable by humans

# 8 deg image at center of gaze, 100 ms viewing time

One example core object recognition test:

**"face"**    **not "face"**



**n>100**                    **n>700**

Another example core object recognition test:

## "Beetle"  |  ## Not "Beetle"



n>100

n>700

## (Domain: core object recognition)

# Goal:  end-to-end understanding

**1. Can we infer the <span style="color:red">decoding</span> mechanism(s) that the brain uses to support perceptual reports about visually presented objects?**

> **Note: this must <span style="color:yellow">predict</span> behavioral report and it must include a falsifiable statement of the <span style="color:yellow">relevant</span> aspects of neural activity (aka "neural code")**

2. Can we infer the encoding mechanism(s) that accurately predicts the relevant ventral stream population patterns of neural activity from each image?

**(Domain: core object recognition)**

# Goal:  end-to-end understanding

**1. Can we infer the <span style="color:red">decoding</span> mechanism(s) that the brain uses to support perceptual reports about visually presented objects?**

> **Note: this must <span style="color:yellow">predict</span> behavioral report and it must include a falsifiable statement of the <span style="color:yellow">relevant</span> aspects of neural activity (aka "neural code")**
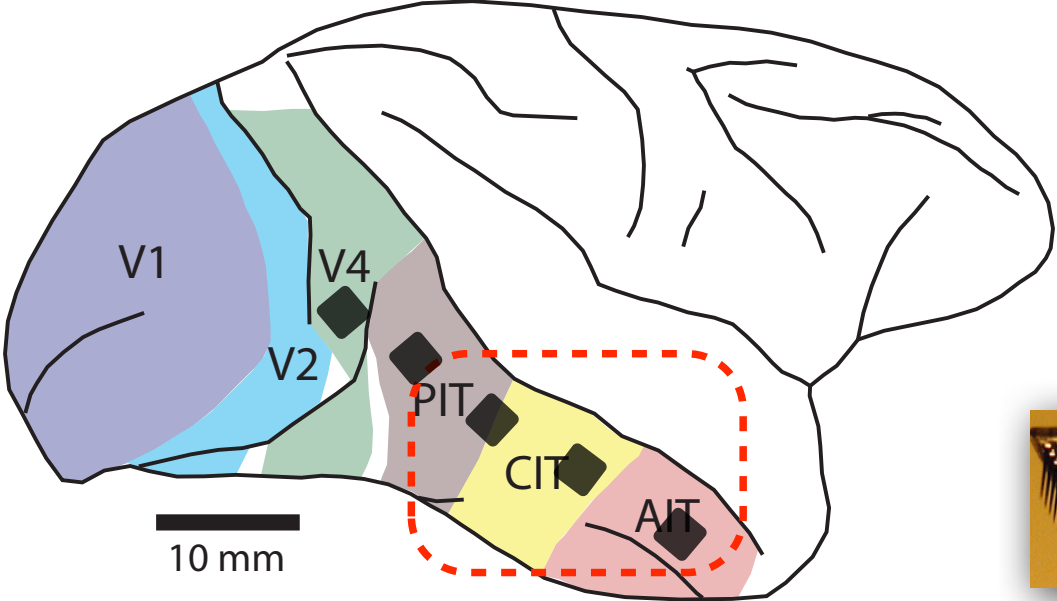
*Adapted from Kelly et al. J. Neurosci (2007)*

V1

V4

V2

PIT

CIT

AIT

10 mm

## Three, 96-electrode arrays

Courtesy of Society for Neuroscience. License CC BY-NC-SA. Source: Kelly, Ryan C., Matthew A. Smith, Jason M. Samonds, Adam Kohn, A. B. Bonds, J. Anthony Movshon, and Tai Sing Lee. "Comparison of recordings from microelectrode arrays and single electrodes in the visual cortex." Journal of Neuroscience 27, no. 2 (2007): 261-264.

Array 1 location

Array 2 location

Array 3 (in place)

1

2

3

*e.g. "response" = mean firing rate 70-170 ms after image onset*
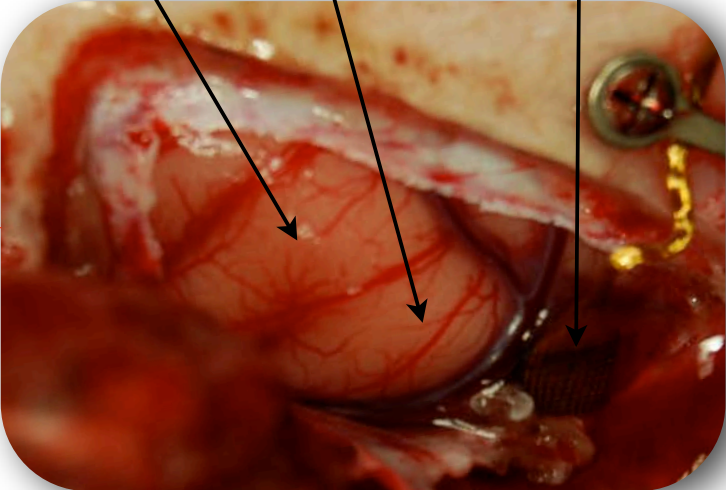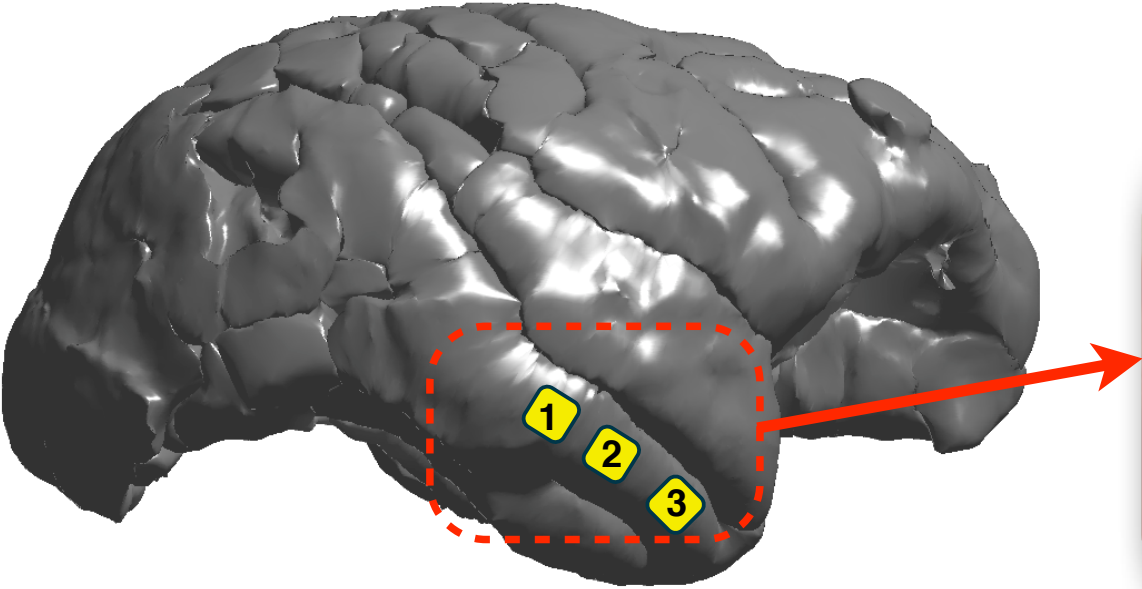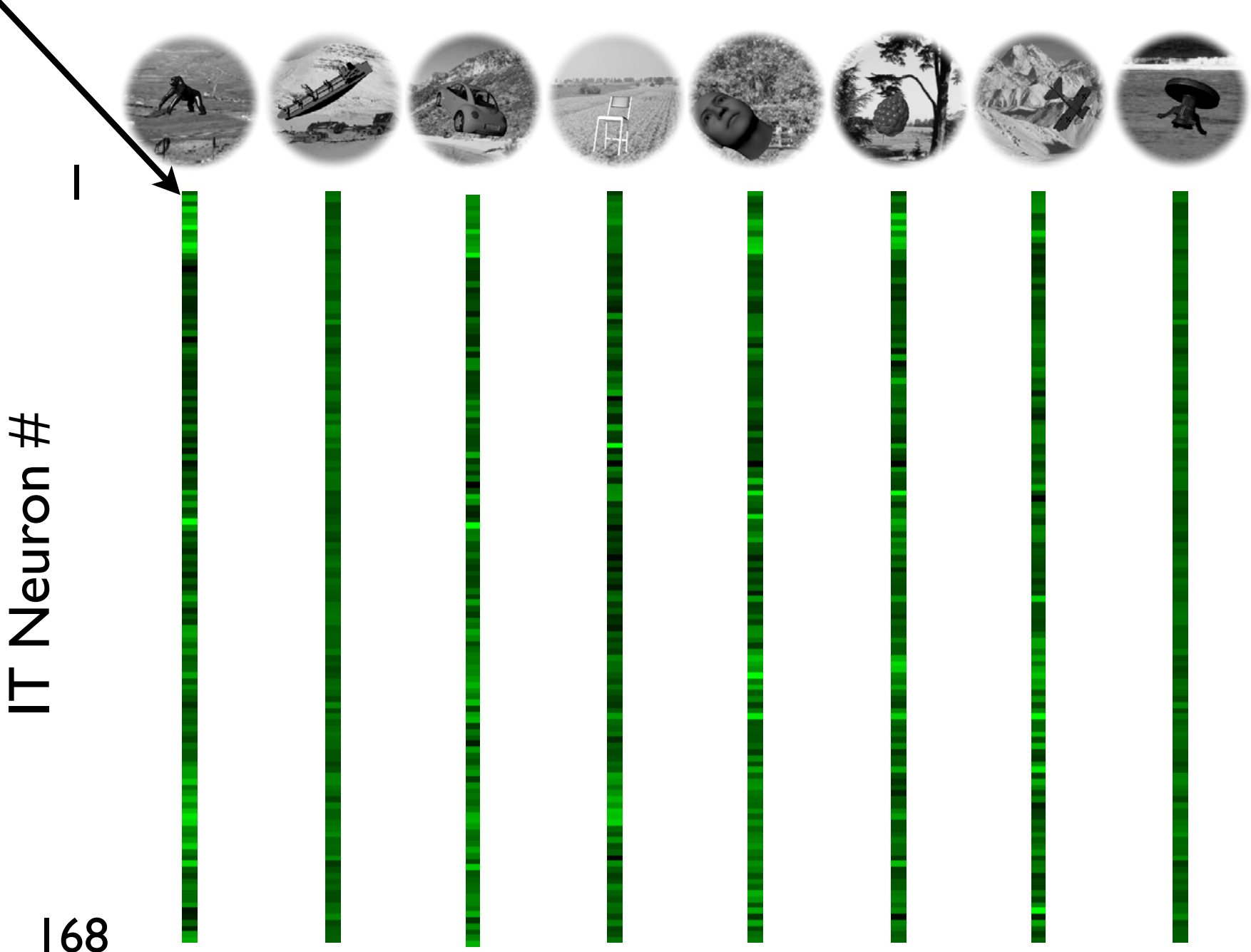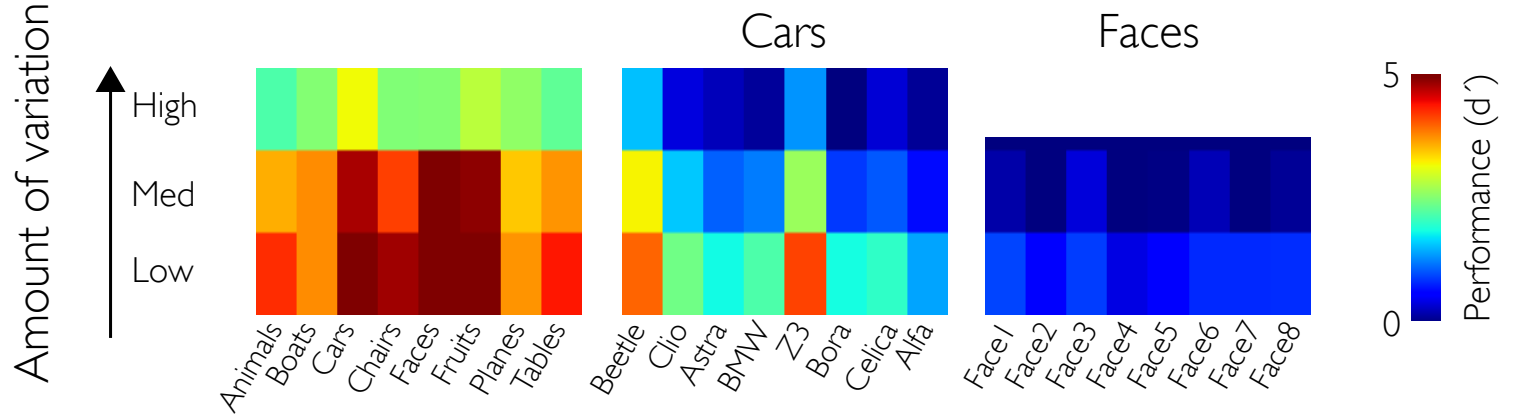


I

IT Neuron #

168

# BEHAVIOR
(64 object recognition tests using same images)

Amount of variation: High / Med / Low

Cars

Faces

Performance (d'): 0 – 5

Animals, Boats, Cars, Chairs, Faces, Fruits, Planes, Tables

Beetle, Clio, Astra, BMW, Z3, Bora, Celica, Alfa

Face1, Face2, Face3, Face4, Face5, Face6, Face7, Face8

## NEURAL ACTIVITY

IT Neuron # (1 – 168)

Image # (1 – 2560)

**Humans and monkeys find some object recognition tests more difficult than others.**

**This pattern of difficulty is very reliable across observers.**

**This pattern is not explained by "low level" visual features.**

**Which, if any, part of the IT population neural activity pattern predicts the observed behavioral performance over _all_ 64 object recognition tests?**

# We had previously shown that simple weighted sums of IT population responses have high performance in recognition tasks
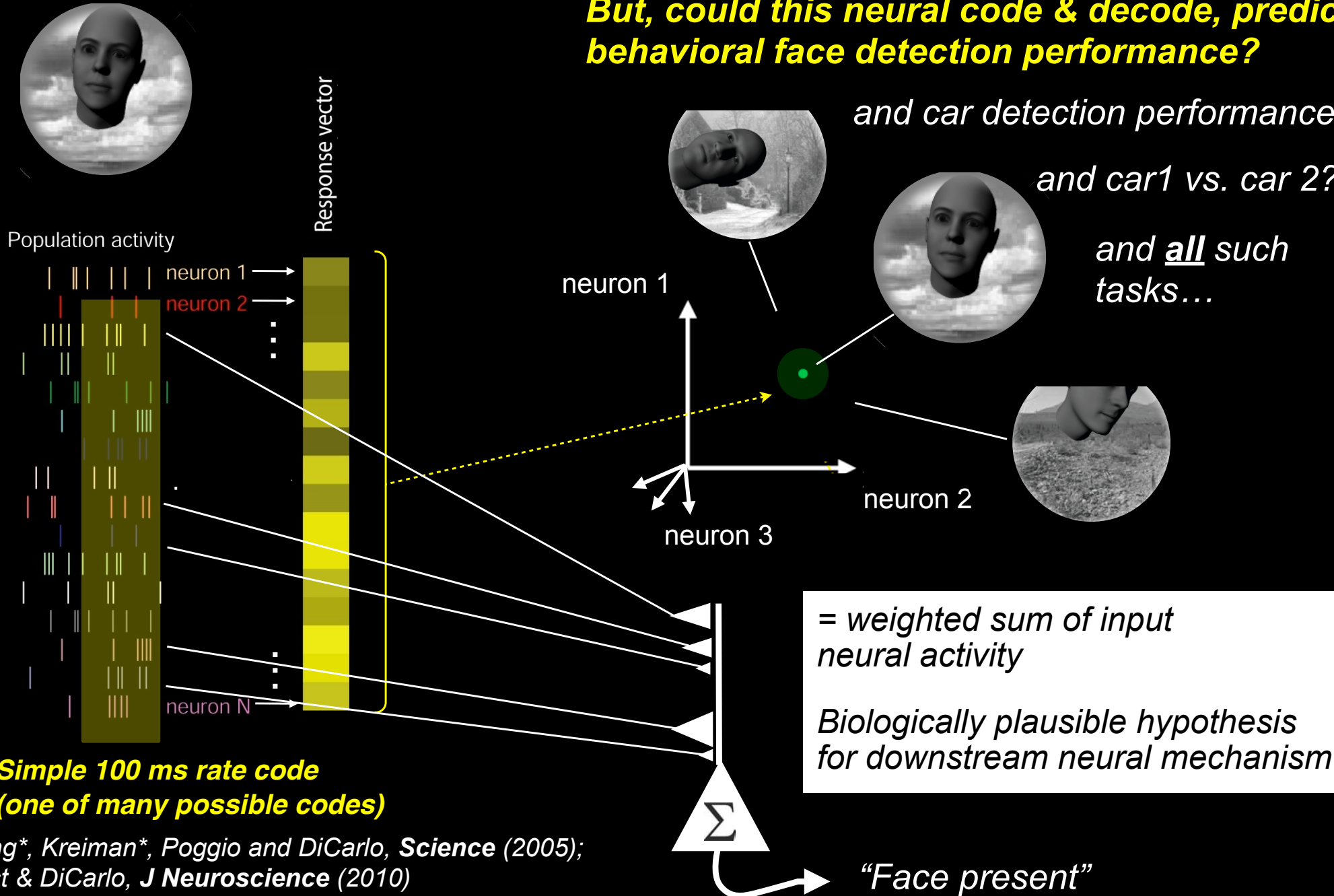
*But, could this neural code & decode, predict behavioral face detection performance?*

*and car detection performance?*

*and car1 vs. car 2?*

*and **all** such tasks…*

Response vector

Population activity

neuron 1
neuron 2

neuron 1

neuron 3

neuron 2

.
.
.

neuron N

**Simple 100 ms rate code (one of many possible codes)**

= *weighted sum of input neural activity*

*Biologically plausible hypothesis for downstream neural mechanism*

$\Sigma$

*"Face present"*

*Hung\*, Kreiman\*, Poggio and DiCarlo, **Science** (2005);*
*Rust & DiCarlo, **J Neuroscience** (2010)*

# What code & decoding mechanism explains object recognition?

## Our working hypothesis from previous work:

**Passively-evoked** spike **rate codes (using a single, fixed time scale)** that are **spatially distributed** over a **single, fixed number** of non-human primate **IT cortex** neurons and **learned from a reasonable number of examples.**

> If correct, this code/decode should predict monkey **and human** reports about object category and object identity for all tasks.

## Other possibilities:

Attentional and/or arousal mechanisms are needed to "activate" IT

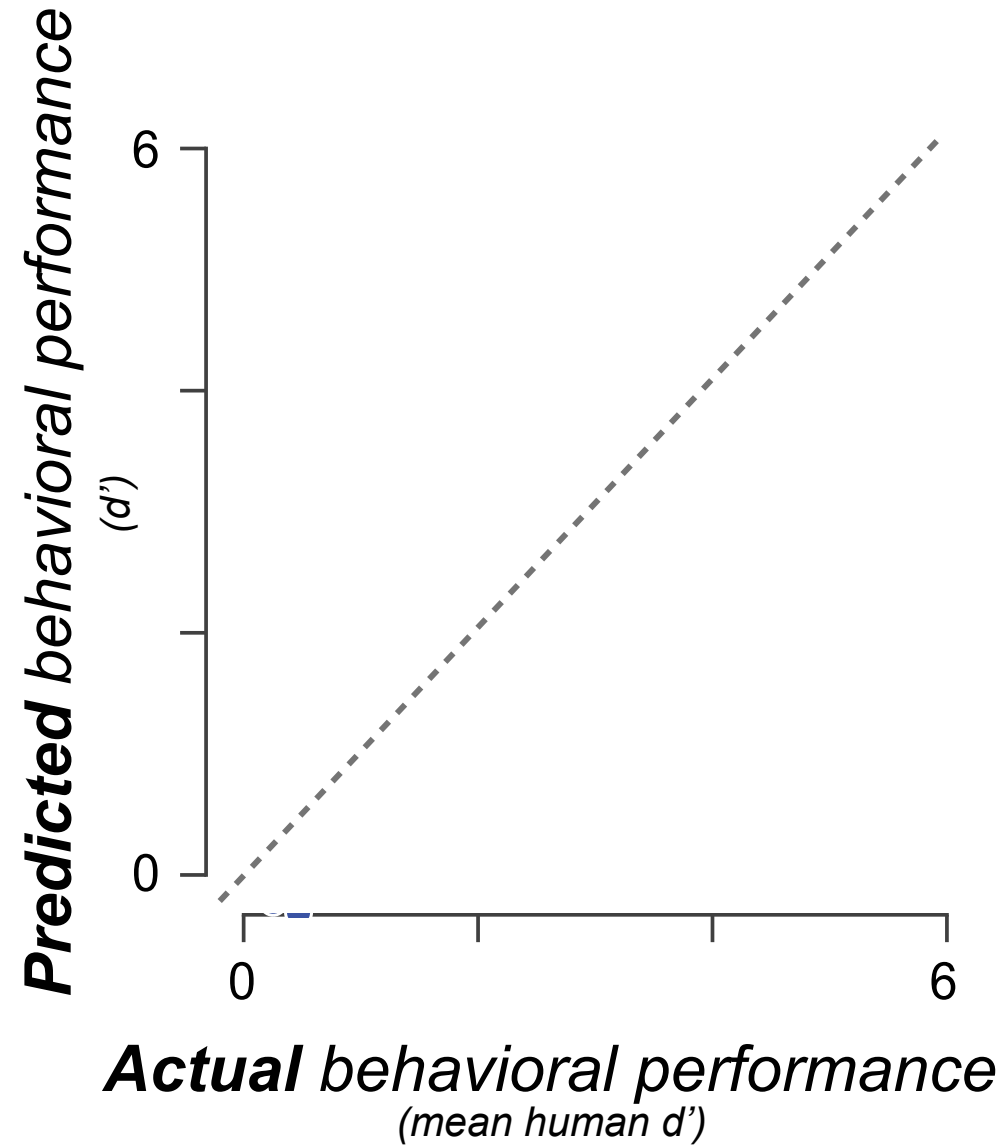Trial-by-trial coordinated spike timing patterns are crucial

Compartments within IT must be carefully considered
(e.g. tasks related to faces handled exclusively by "face patch" network)

IT does not directly underlie object recognition

Performance requires too many training examples

Monkey neuronal codes cannot explain human behavior

**Our first decoder (based on previous work), with number of neurons chosen (once) to match human performance**

**Predicted** *behavioral performance* *(d')*

6

0

**Actual** *behavioral performance* *(mean human d')*
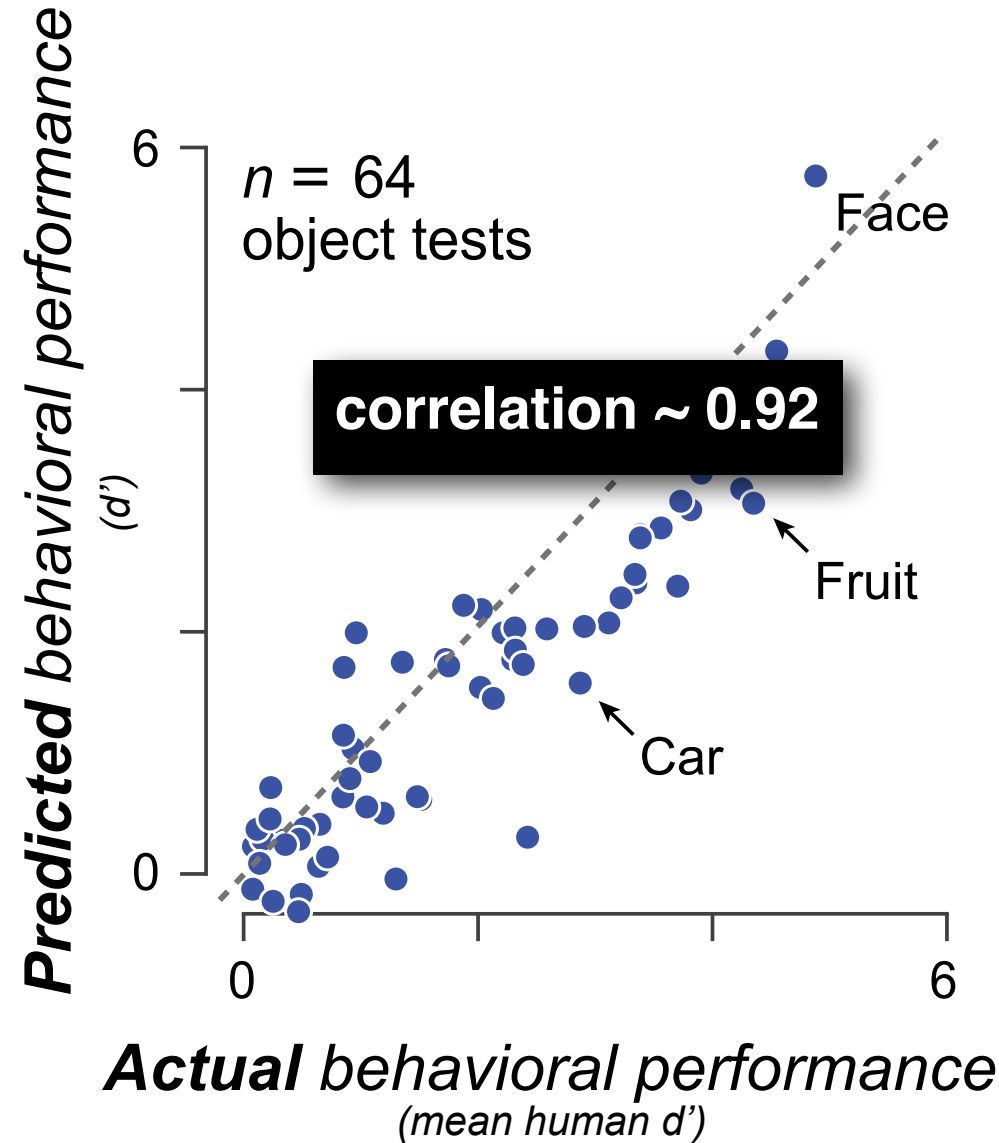
0          6

**Take home: simple, learned weighted sums of IT firing rates accurately predict the pattern of PERFORMANCE over all object recognition tests**

**Parameters of inferred neural code/decoding mechanism:**

*- for each new object, <u>randomly</u> sample ~50,000 single neurons spatially <u>distributed</u> over <u>IT</u>*

*- "listen" to each IT site's <u>average</u> spiking response (ave over <u>100 ms</u>)*

*- <u>learn</u> an appropriately <u>weighted sum</u> of those IT spiking outputs, and then use ~10% of them to judge the likelihood of the object being present*

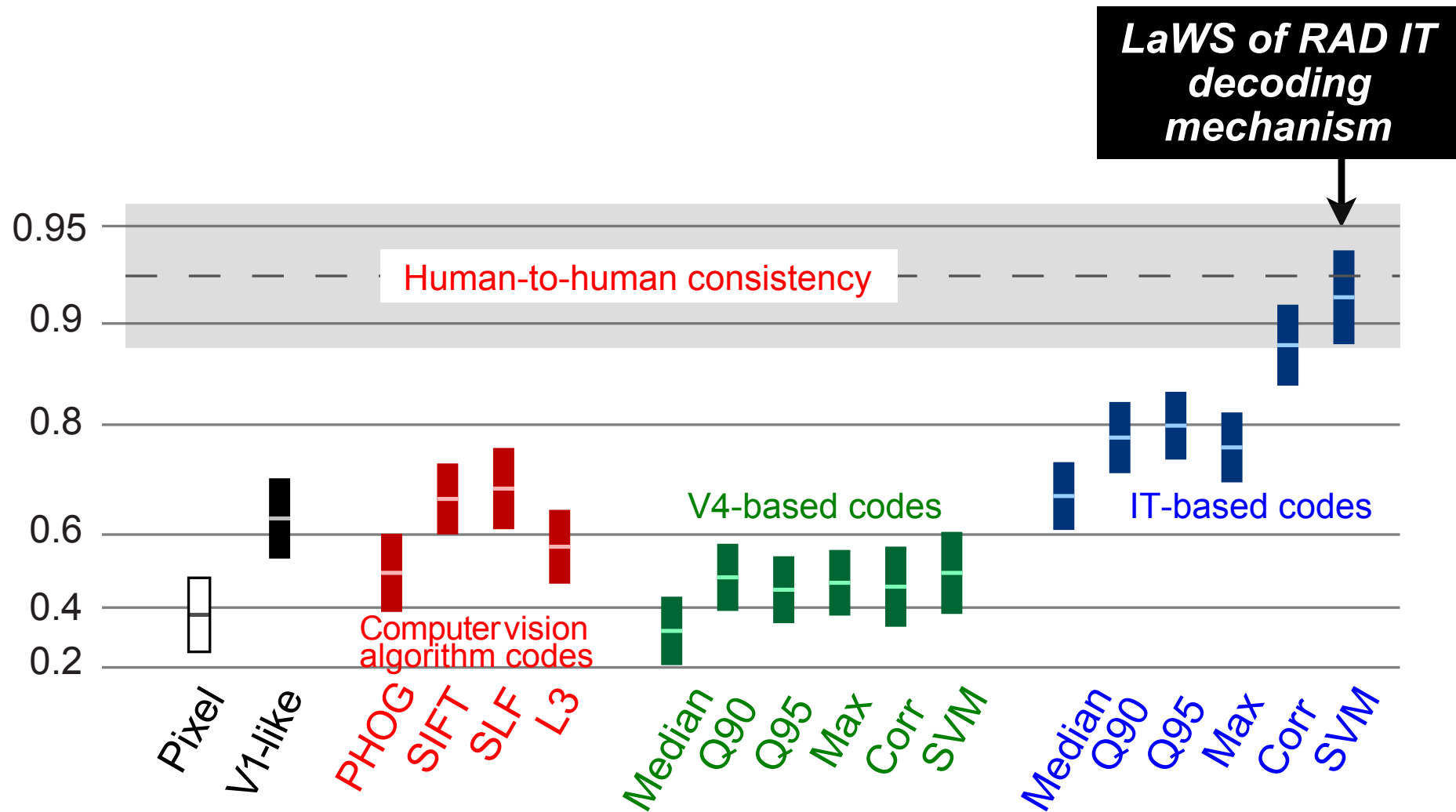**L**earned **W**eighted **S**ums of (~50,000) **R**andom **A**verage (100 ms) single unit responses **D**istributed over **IT**

**"LaWS of RAD IT" decoding mechanism**

**Predicted** *behavioral performance* (d')

**Actual** *behavioral performance* (mean human d')

n = 64 object tests

**correlation ~ 0.92**

Face

Fruit

Car

6

0

0        6

*Majaj, Hong, Solomon, and DiCarlo, **Cosyne 2012***
*Majaj, Hong, Solomon, and DiCarlo, **Under Review***

# Some controls…
# Most alternative codes/decoding mechanisms are not even close.



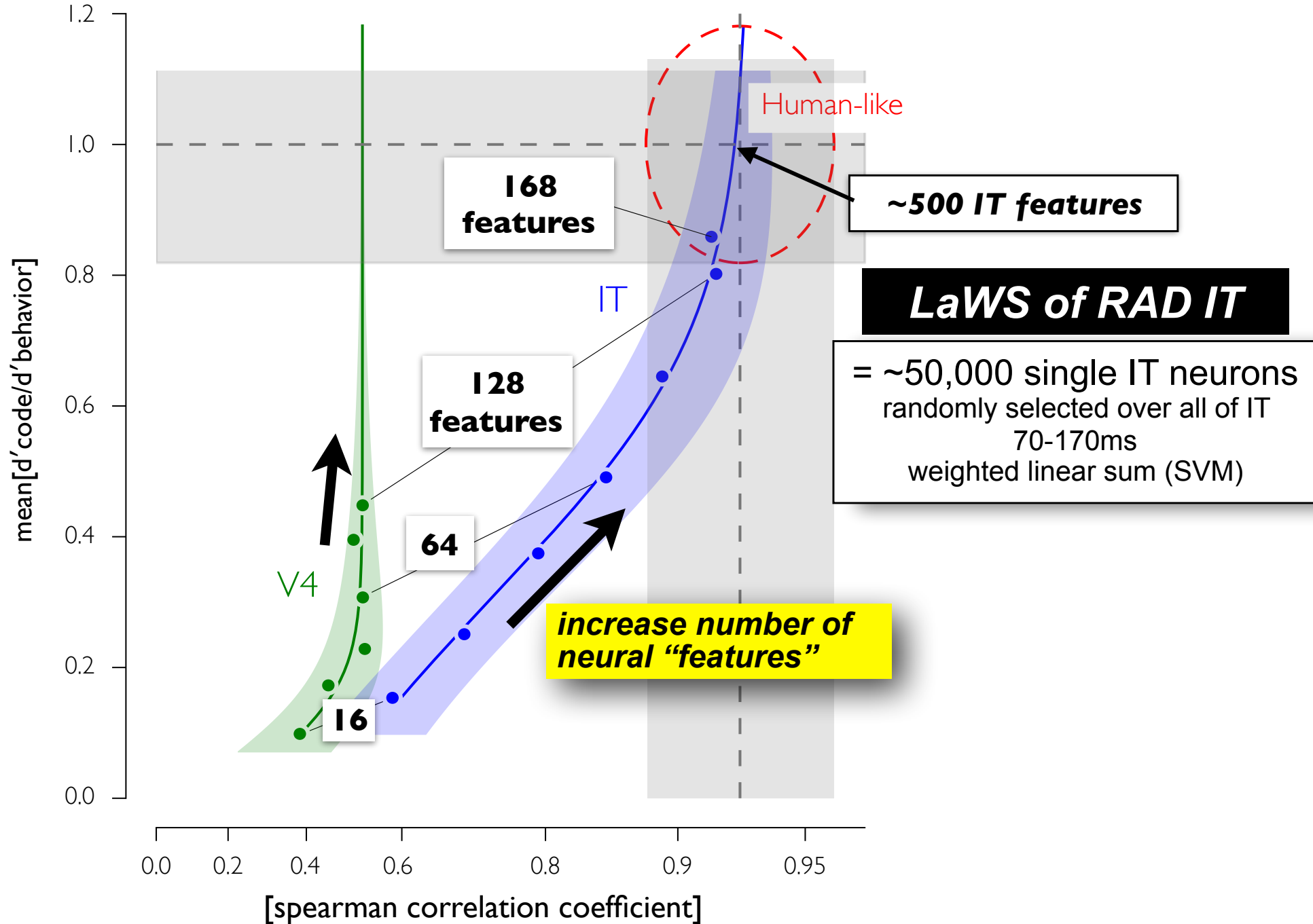Correlation of model performance predictions with human performance

LaWS of RAD IT decoding mechanism

Human-to-human consistency

Computer vision algorithm codes

V4-based codes

IT-based codes

Pixel · V1-like · PHOG · SIFT · SLF · L3 · Median · Q90 · Q95 · Max · Corr · SVM · Median · Q90 · Q95 · Max · Corr · SVM

0.95 · 0.9 · 0.8 · 0.6 · 0.4 · 0.2

*Majaj, Hong, Solomon, and DiCarlo,* **Cosyne 2012**
*Majaj, Hong, Solomon, and DiCarlo,* **Under Review**

79

Consistency with humans

Number of single units needed to support single-trial performance

Number of neural "features" (multi-unit, trial averaged)

Number of output neurons in IT
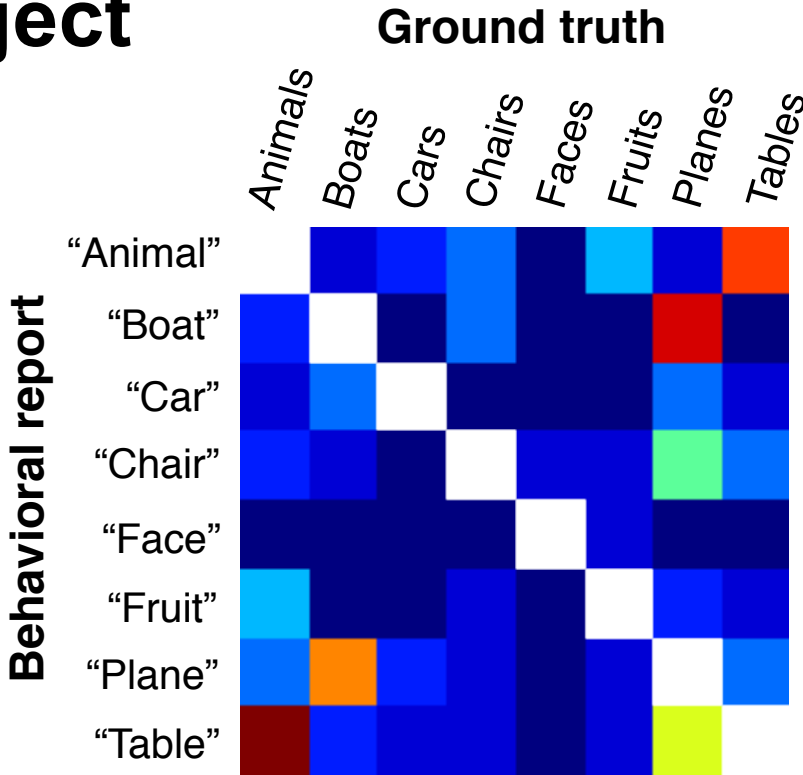
~10 M

~50,000 single IT neurons

$10^6$
$10^5$
$10^4$
$10^3$
~500
$10^2$
$10^1$
$10^0$

*

~100

A family of IT codes/decodes that each accurately predict pattern of behavioral performance

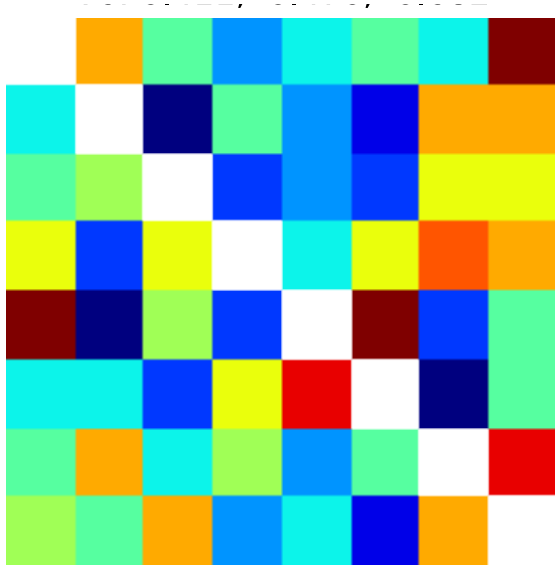$10^0$ $10^2$ $10^4$ $10^6$ $10^8$ $10^{10}$ $10^{12}$ $10^{14}$ $10^{16}$

Number of training examples per object

# Behavioral object confusions

**Ground truth**

*Predicted:*

**LaWS of RAD IT decoding mechanism**



Animals · Boats · Cars · Chairs · Faces · Fruits · Planes · Tables

Behavioral report: "Animal" · "Boat" · "Car" · "Chair" · "Face" · "Fruit" · "Plane" · "Table"
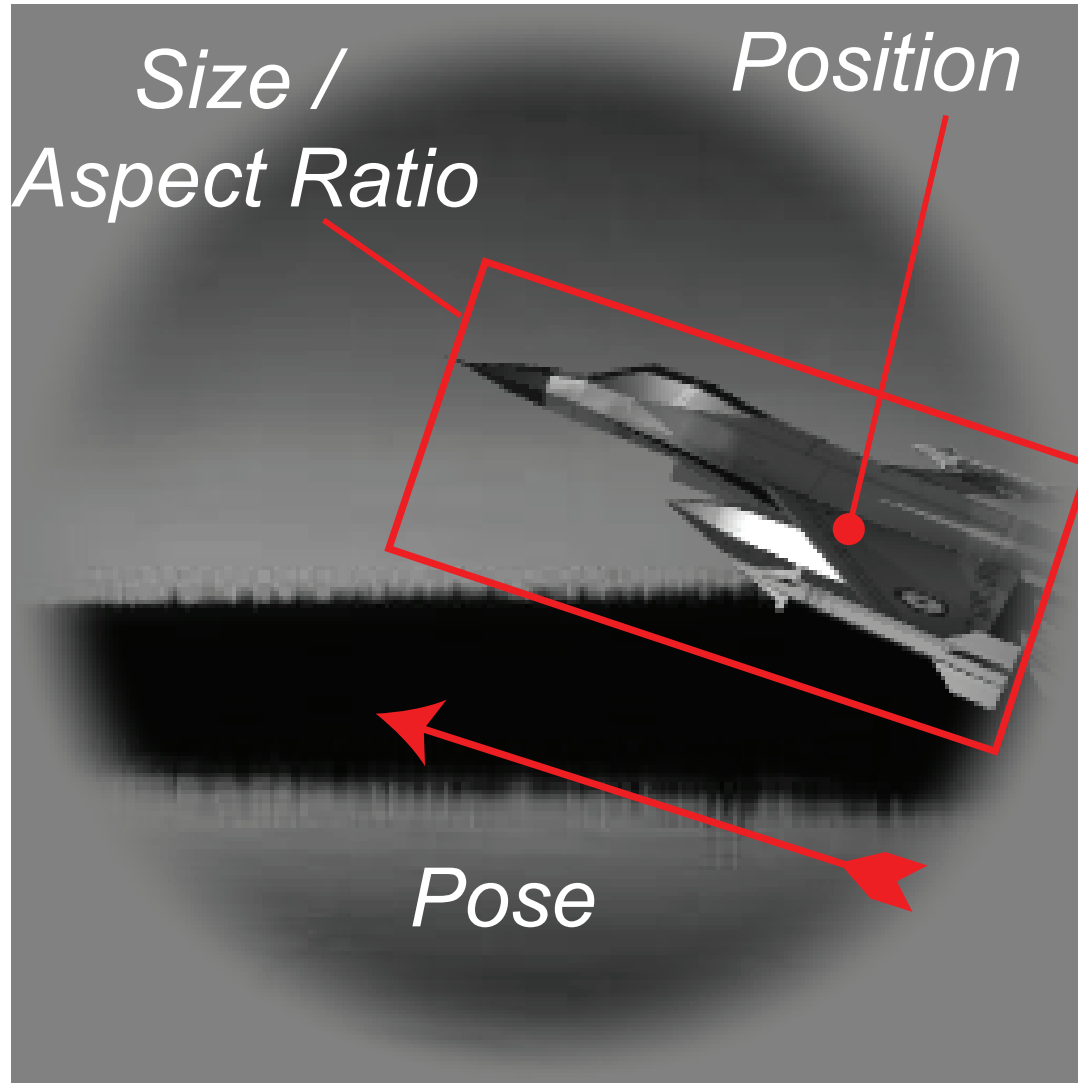
*Noise-corrected correlation:* **0.91** ↕

**0.68** ↕

## This is an opportunity to push forward: image grain predictions to distinguish among alternative IT codes

*High variation*

# Other object latent variables



Category: plane
Identity: f16

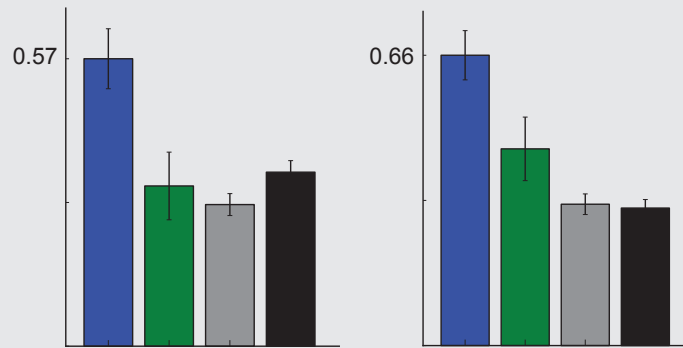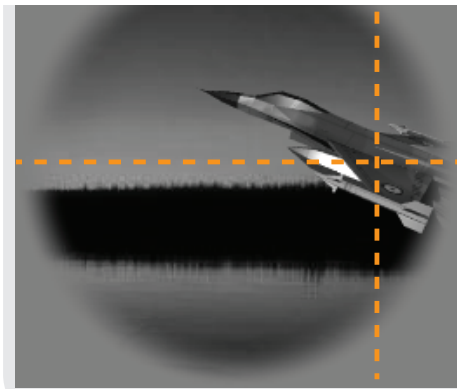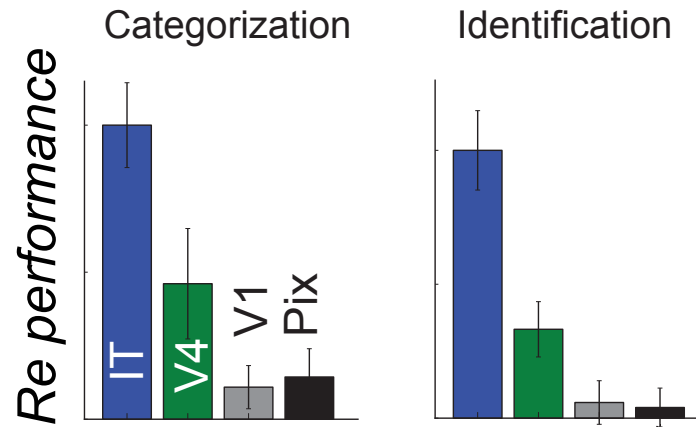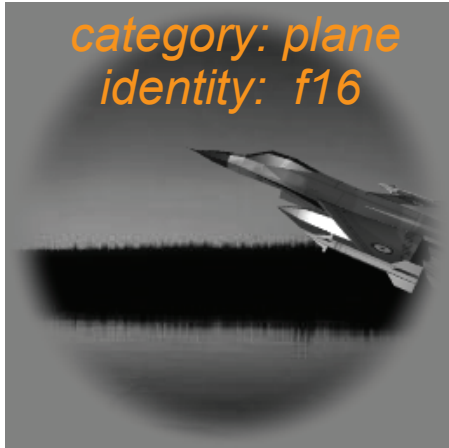Size / Aspect Ratio

Position

Pose

Source: Hong, Ha, Daniel LK Yamins, Najib J. Majaj, and James J. DiCarlo. "Explicit information for category-orthogonal object properties increases along the ventral stream." Nature neuroscience 19, no. 4 (2016): 613-622.

category: plane
identity: f16

Categorization

Identification

Re performance

IT
V4
V1
Pix

0.57

0.66

Site 10

Site 54

Site 43

Site 11

Site 77

Site 102

y

x

Source: Hong, Ha, Daniel LK Yamins, Najib J. Majaj, and James J. DiCarlo. "Explicit information for category-orthogonal object properties increases along the ventral stream." Nature neuroscience 19, no. 4 (2016): 613-622.

84

Source: Hong, Ha, Daniel LK Yamins, Najib J. Majaj, and James J. DiCarlo. "Explicit information for category-orthogonal object properties increases along the ventral stream." Nature neuroscience 19, no. 4 (2016): 613-622.

**But these tasks are not all equally difficult for humans. Does this decoding mechanism predict that pattern of difficulty?**

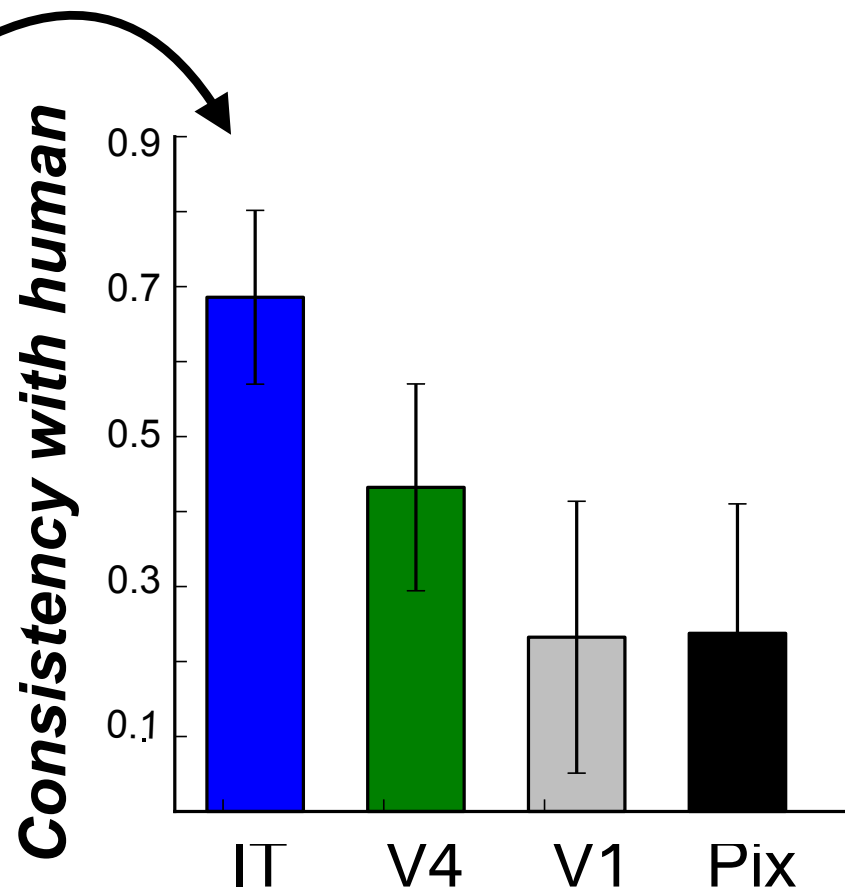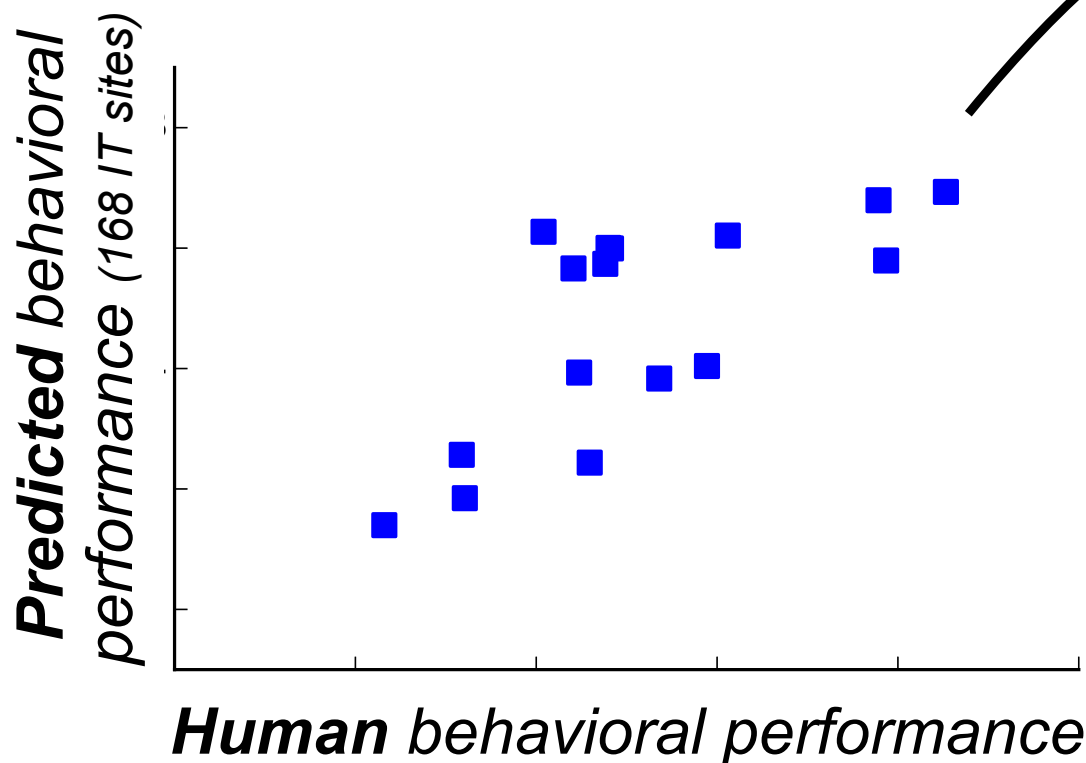**To test this, we collected human performance data on these images/tasks.**

a



Fraction of Human Performance

Number of Neural Sites

Basic Categorization · Subordinate Identification · X-axis Position · Y-axis Position · Bounding Box Size · X-axis Size · Y-axis Size · 3-D Object Scale · Major Axis Length · Aspect Ratio · Z-axis Rotation · X-axis Rotation

IT, V4, V1, Pix

# Number of IT sites needed to match human performance

|  | IT | V4 | V1 | Pix |
|---|---|---|---|---|
| Basic Categorization | 520 +/- 165 | 8.84 x 10^5 | --- | --- |
| Subordinate Identification | 444 +/- 61 | 9.15 x 10^6 | --- | --- |
| X-axis Position | 1624 +/- 44 | 4.5 x 10^6 | 3 x 10^7 | --- |
| Y-axis Position | 647 +/- 215 | 1.1 x 10^5 | 8.7 x 10^6 | --- |
| Bounding Box Size | 234 +/- 91 | 8.4 x 10^3 | --- | --- |
| X-axis Size | 150 +/- 55 | 2.1 x 10^3 | 3.4 x 10^7 | --- |
| Y-axis Size | 182 +/- 62 | 7.8 x 10^3 | 9.5 x 10^6 | --- |

|  | IT | V4 | V1 | Pix |
|---|---|---|---|---|
| 3-D Object Scale | 339 +/- 79 | 1.9 x 10^5 | --- | --- |
| Major Axis Length | 165 +/- 59 | 5.7 x 10^3 | --- | --- |
| Aspect Ratio | 103 +- 37 | 922 +/- 59 | 6.5 x 10^3 | --- |
| Major Axis Angle | 520 +/- 165 | 520 +/- 165 | --- | --- |
| Z-axis Rotation | 1206 +/- 473 | --- | --- | --- |
| Y-axis Rotation | 1317 +/- 459 | 1.1 x 10^5 | --- | --- |
| X-axis Rotation | 775 +/- 248 | --- | --- | --- |



*Predicted behavioral performance (168 IT sites)* vs *Human behavioral performance*

*Consistency with human* — IT, V4, V1, Pix

Source: Hong, Ha, Daniel LK Yamins, Najib J. Majaj, and James J. DiCarlo. "Explicit information for category-orthogonal object properties increases along the ventral stream." Nature neuroscience 19, no. 4 (2016): 613-622.

Category: plane
Identity: f16



Size /
Aspect Ratio

Position

Pose

**Summary: This ventral stream code/decoding mechanism also predicts human patterns of performance for other object latent variables.**

*This suggests that:*

- *the IT population conveys a general purpose object representation*

- *the job of the ventral stream is not to produce category "invariant" representations*

*Edelman (1998), DiCarlo and Cox (2007),
Li et al. (2009), etc.*

*Hong, Yamins, Majaj, and DiCarlo, **Cosyne 2014***

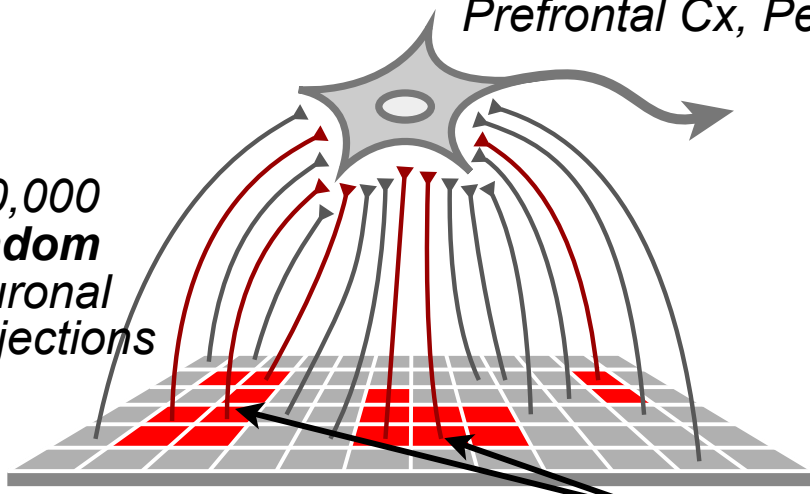*Hong, Yamins, Majaj, and DiCarlo, **(in prep)***

# *Sketch of the inferred anatomy:*

**LaWS of RAD IT** *[70-170ms, 50,000n, 100t]*
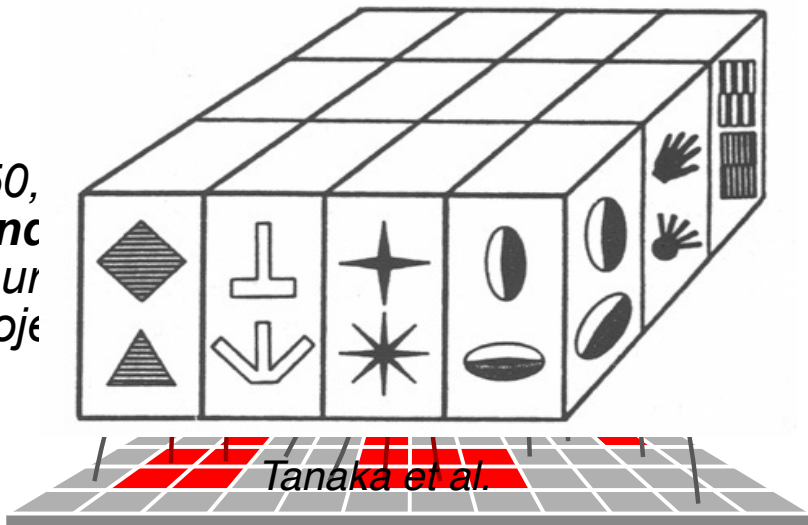
*Prefrontal Cx, Perirhinal Cx, Amygdala*

*~50,000* **random** *neuronal projections*

**IT cortex** *(AIT + CIT)*   *"Face patches"*
*(2-5 mm)*

*~50,*
**ran**
*neu*
*proj*

*Tanaka et al.*

Source: Tanaka, Keiji. "Neuronal mechanisms of object recognition."
Science-New York Then Washington 262 (1993): 685-685.
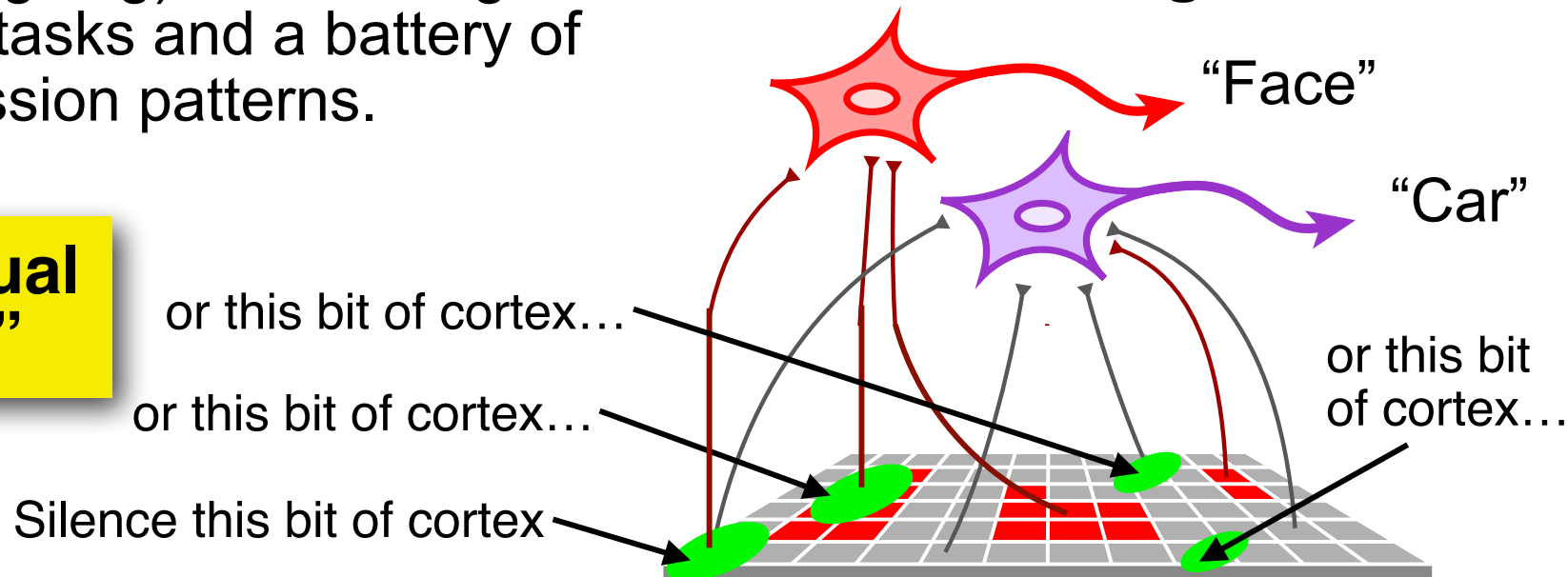
# Causal tests of this model

**The model allows us to predict how much any object recognition task will be disrupted by direct suppression of IT neurons.**

Step 1:  (done) Tool building and testing:   Can we reliably disrupt performance of a recognition task by directly suppressing the activity of ~1mm IT neural sub-populations?

Step 2 (ongoing):  Test a large battery of tasks and a battery of IT suppression patterns.
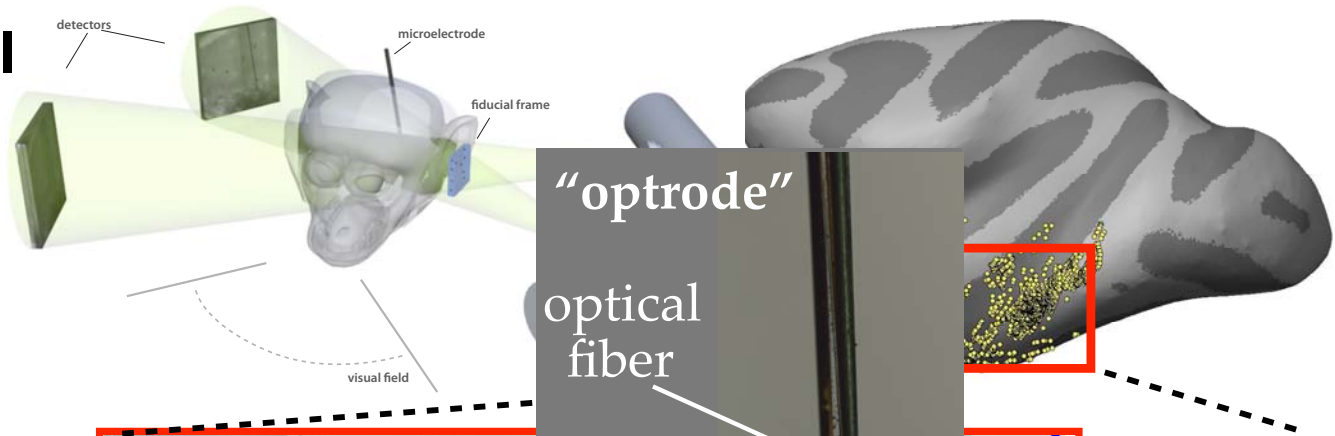
**Post-learning:**

**Towards actual "inception"**

"Face"

"Car"

or this bit of cortex…

or this bit of cortex…

Silence this bit of cortex

or this bit of cortex…

*IT cortex* (AIT + CIT) *~150 IT sub-regions, each ~1 mm in scale*

# Stereo, microfocal x-ray system



detectors
microelectrode
fiducial frame
visual field

# Optogenetic (ArchT, CAG, AAV) suppression of visually-driven IT activity

**"optrode"**

optical fiber

electrode

~1 mm



**Control**

**Laser on**

Neural response (spikes/s)

*Same visual input on interleaved trials*

n = 1384

Time from image onset (msec)

Courtesy of Society for Neuroscience. License CC BY NC SA.
Source: Issa, Elias B., and James J. DiCarlo. "Precedence of the eye region in neural processing of faces." Journal of Neuroscience 32, no. 47 (2012): 16666-16682.
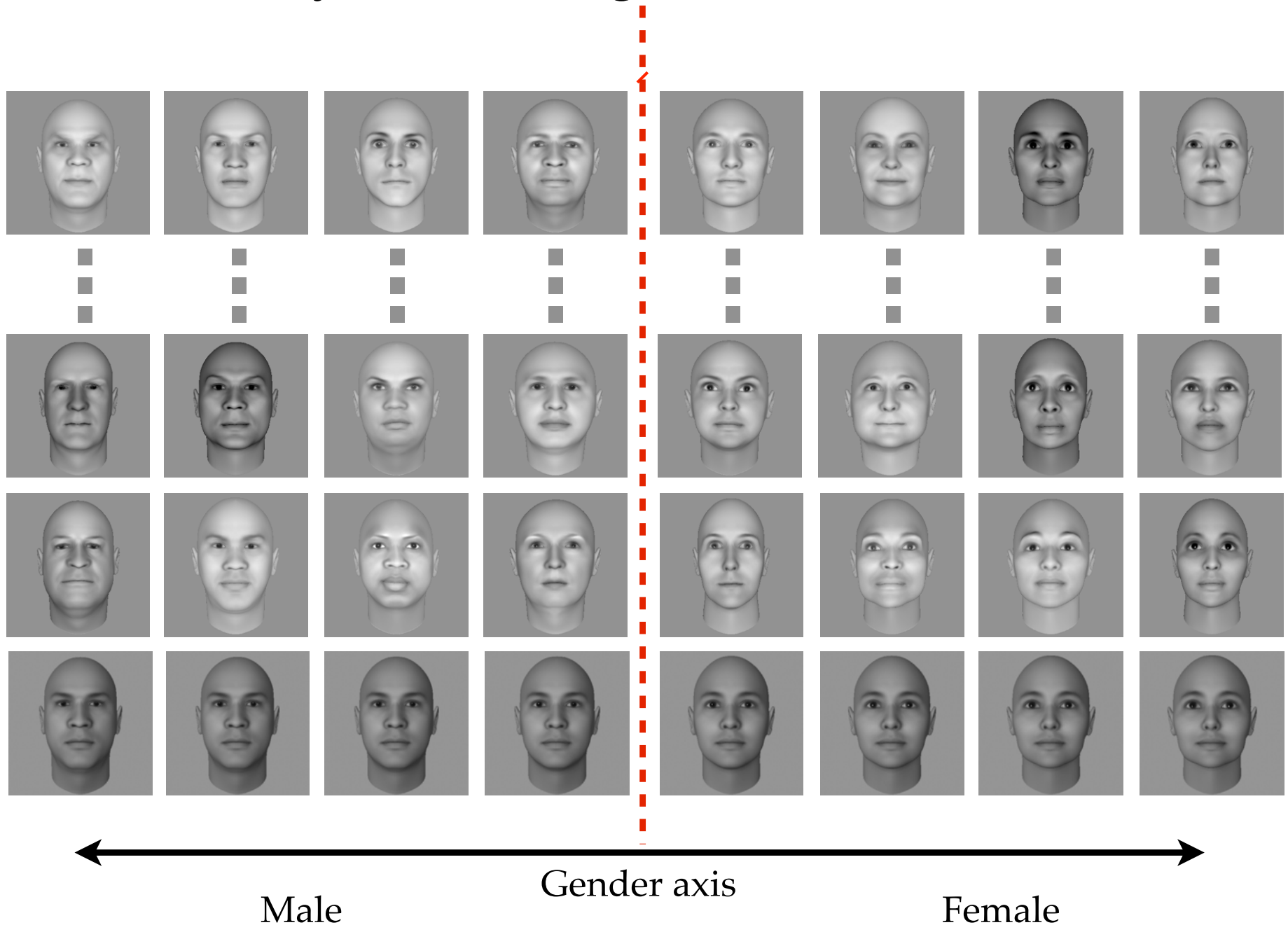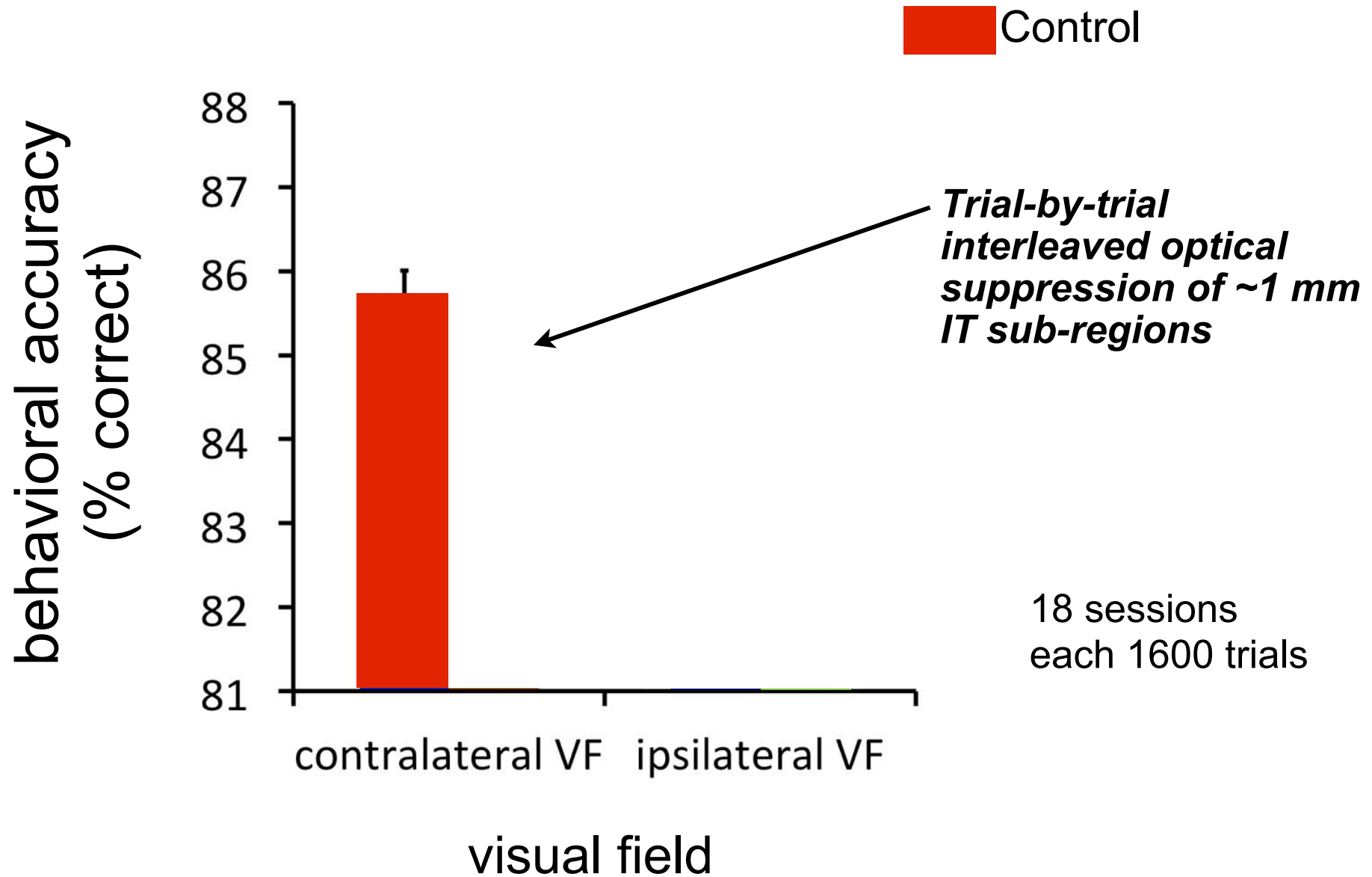
*Afraz, Boyden and DiCarlo, **SFN** (2013)*

*Issa and DiCarlo, **J Neurosci** (2012)*

vs

**face**

**object**

# Monkey task: face gender discrimination



Gender axis

Male  Female

Source: Afraz, Arash, Edward S. Boyden, and James J. DiCarlo."Optogenetic and pharmacological suppression of spatial clusters of face neurons reveal their causal role in face gender discrimination." Proceedings of the National Academy of Sciences 112, no. 21 (2015): 6730-6735.

# We found a spatially-specific behavioral effect on this object discrimination task



**Control**

*Trial-by-trial interleaved optical suppression of ~1 mm IT sub-regions*

18 sessions each 1600 trials

*Afraz, Boyden and DiCarlo,* **SFN** *(2013),* **VSS** *(2014);* **PNAS** *(2015)*
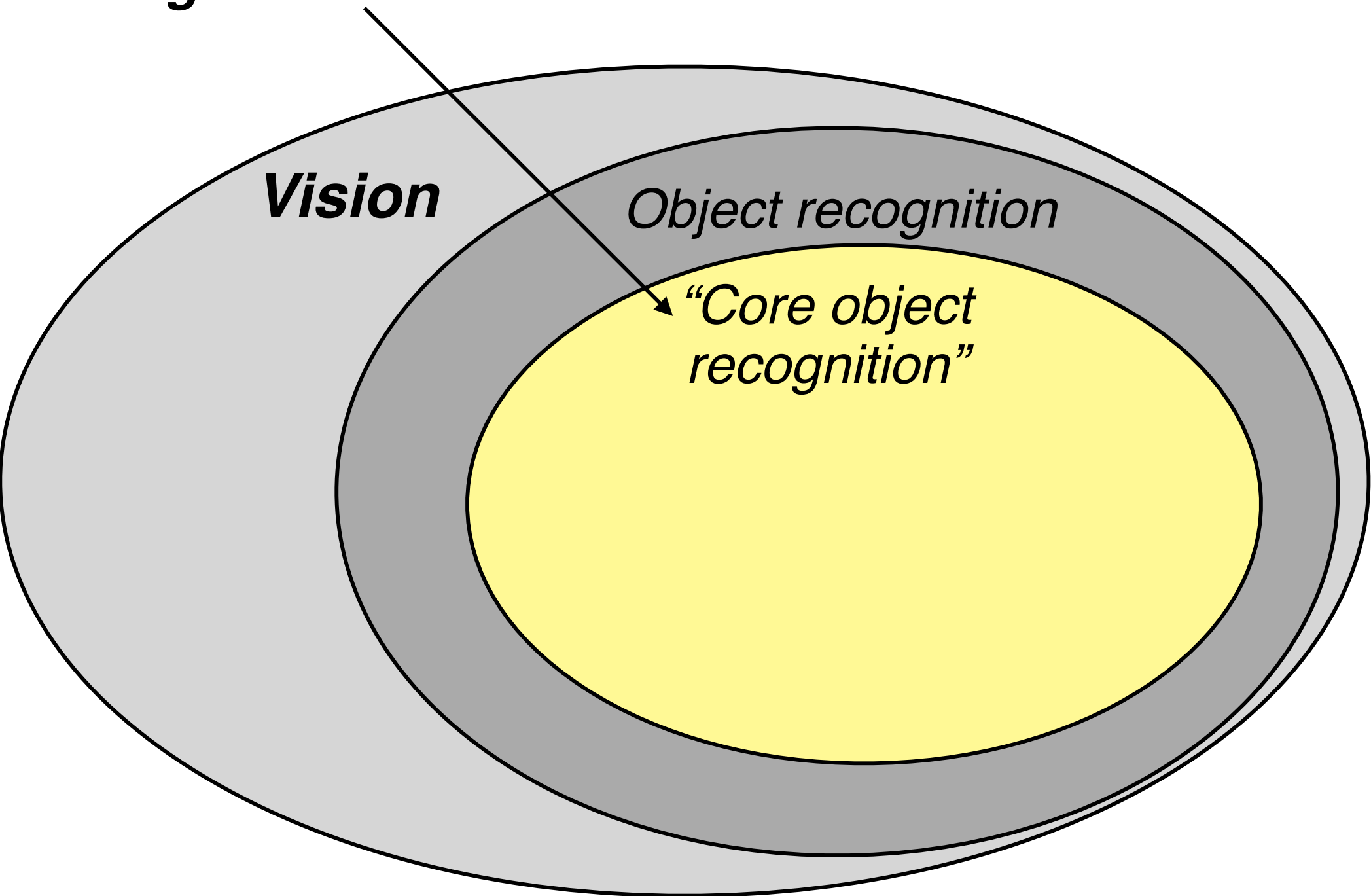
# Pharmacological suppression of different IT sub-regions results in different patterns of deficit in basic level object tasks



Change in behavioral performance

**IT site 1** — Delta d'

**IT site 2** — Delta d'

**IT site 3** — Delta d'

p = 0.10
p = 0.05
p = 0.01

Object discrimination task

*Our current aim is to* **systematically** *measure the specific pattern of behavioral change induced by suppression of each IT sub-region (~100) and compare with model predictions*
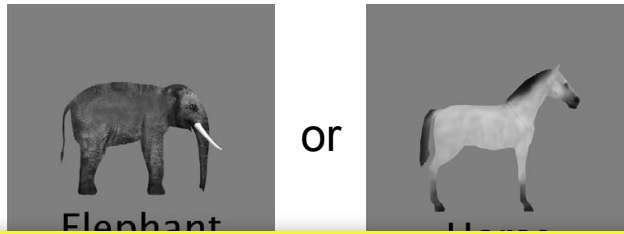
**Can we span the entire domain of core recognition tasks? How?**



Vision

Object recognition

"Core object recognition"

Presentation (100ms)

Choices on this particular trial (post-cue, many possible)

Confusion matrix for an object pair

Elephant or Horse

Stimuli

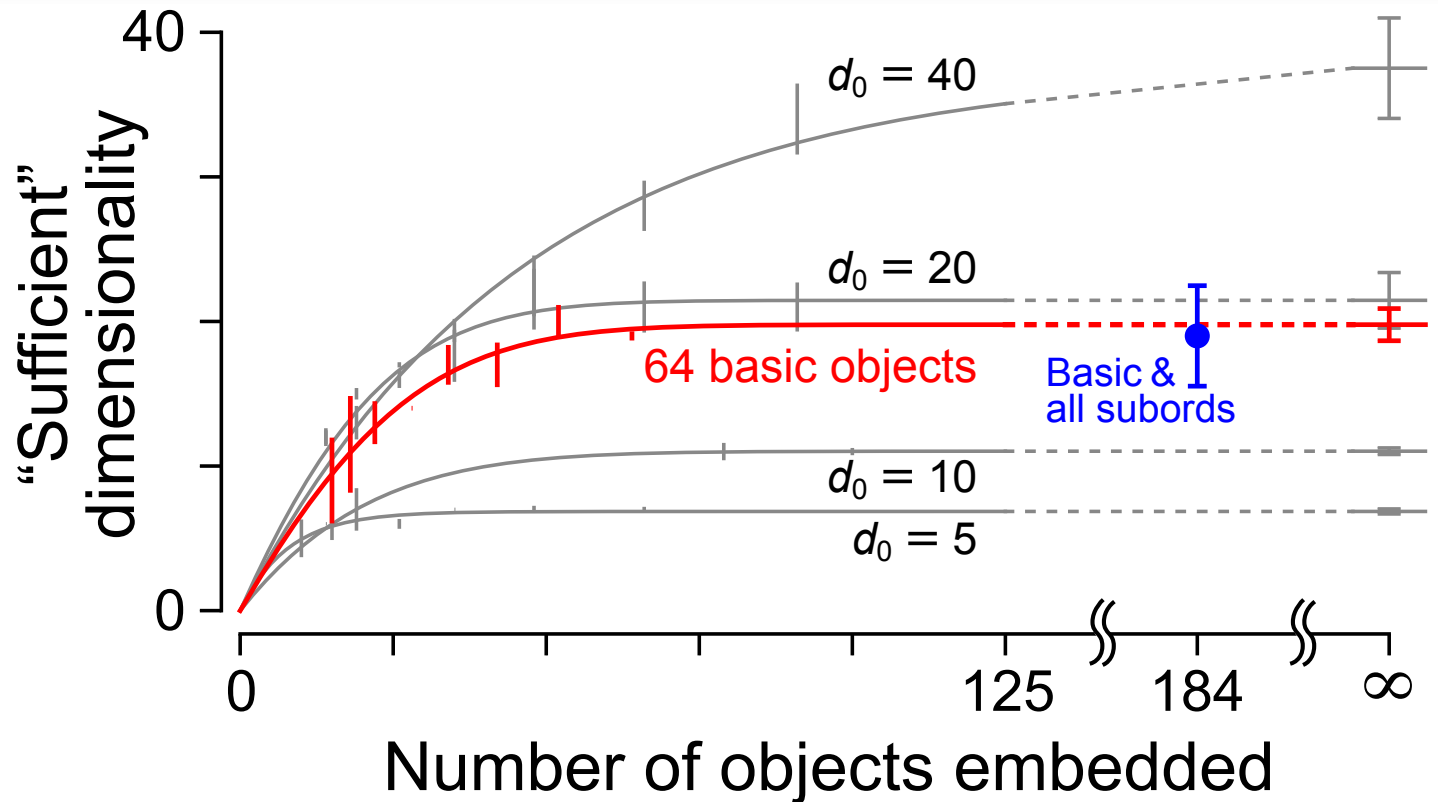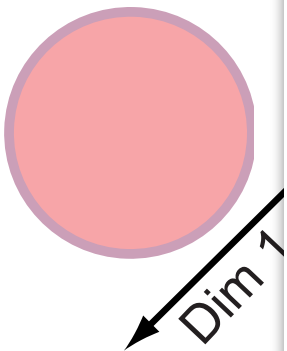| | E | H | |
|---|---|---|---|
| E | 120 | 10 | ... 8,556 matrices |
| H | 5 | 115 | |

Response

**Core recognition: only ~20 dimensions needed to characterize confusions among all basic and subordinate-level objects**

Faces

Cars

Dim 1

"Sufficient" dimensionality

$d_0 = 40$

$d_0 = 20$

64 basic objects

Basic & all subords

$d_0 = 10$

$d_0 = 5$

40

0

0    125    184    ∞

Number of objects embedded

Hong*, Solomon*, Yamins*, and DiCarlo. Large-scale Characterization of a Universal and Compact Visual Perceptual Space. VSS, 2014; in prep

**A**

**C**
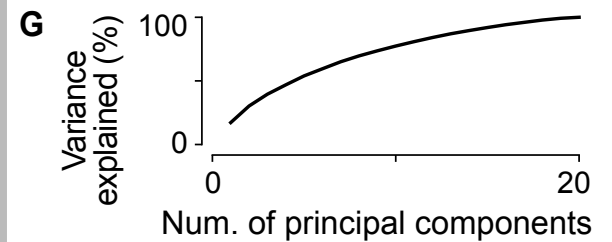
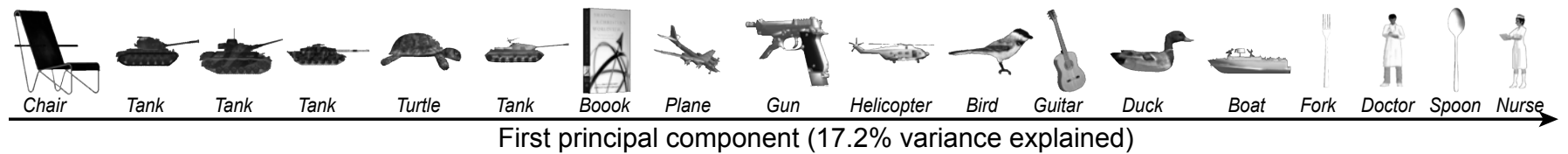Axes in this space correspond to human shape adjectives (subjective magnitude reports)

One important use of this result: for efficient causal testing of the entire domain, we can focus on measuring impacts on object discrimination tasks that span this space
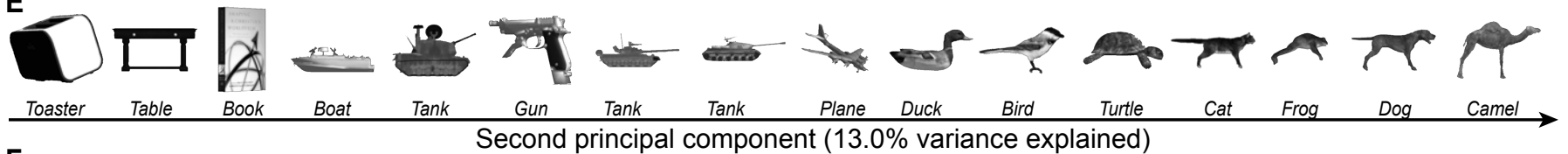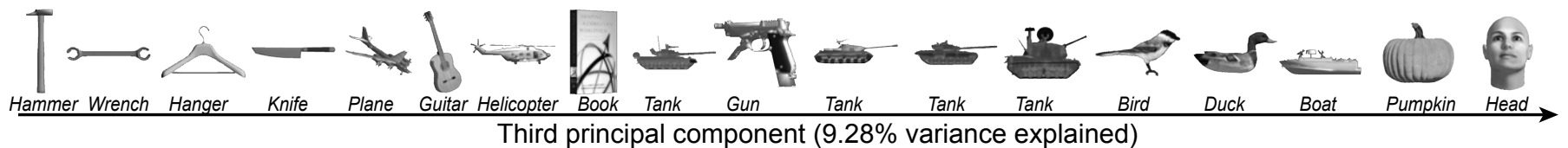
**B**

Ongoing ….

**G**

Variance explained (%)

100

0

0                    20

Num. of principal components

**D**

Chair | Tank | Tank | Tank | Turtle | Tank | Boook | Plane | Gun | Helicopter | Bird | Guitar | Duck | Boat | Fork | Doctor | Spoon | Nurse

First principal component (17.2% variance explained)

**E**

Toaster | Table | Book | Boat | Tank | Gun | Tank | Tank | Plane | Duck | Bird | Turtle | Cat | Frog | Dog | Camel

Second principal component (13.0% variance explained)

**F**

Hammer | Wrench | Hanger | Knife | Plane | Guitar | Helicopter | Book | Tank | Gun | Tank | Tank | Tank | Bird | Duck | Boat | Pumpkin | Head

Third principal component (9.28% variance explained)

97

# Goal:  end-to-end understanding

**1. Can we infer the decoding mechanism that the brain uses to support perceptual reports about visually presented object?**
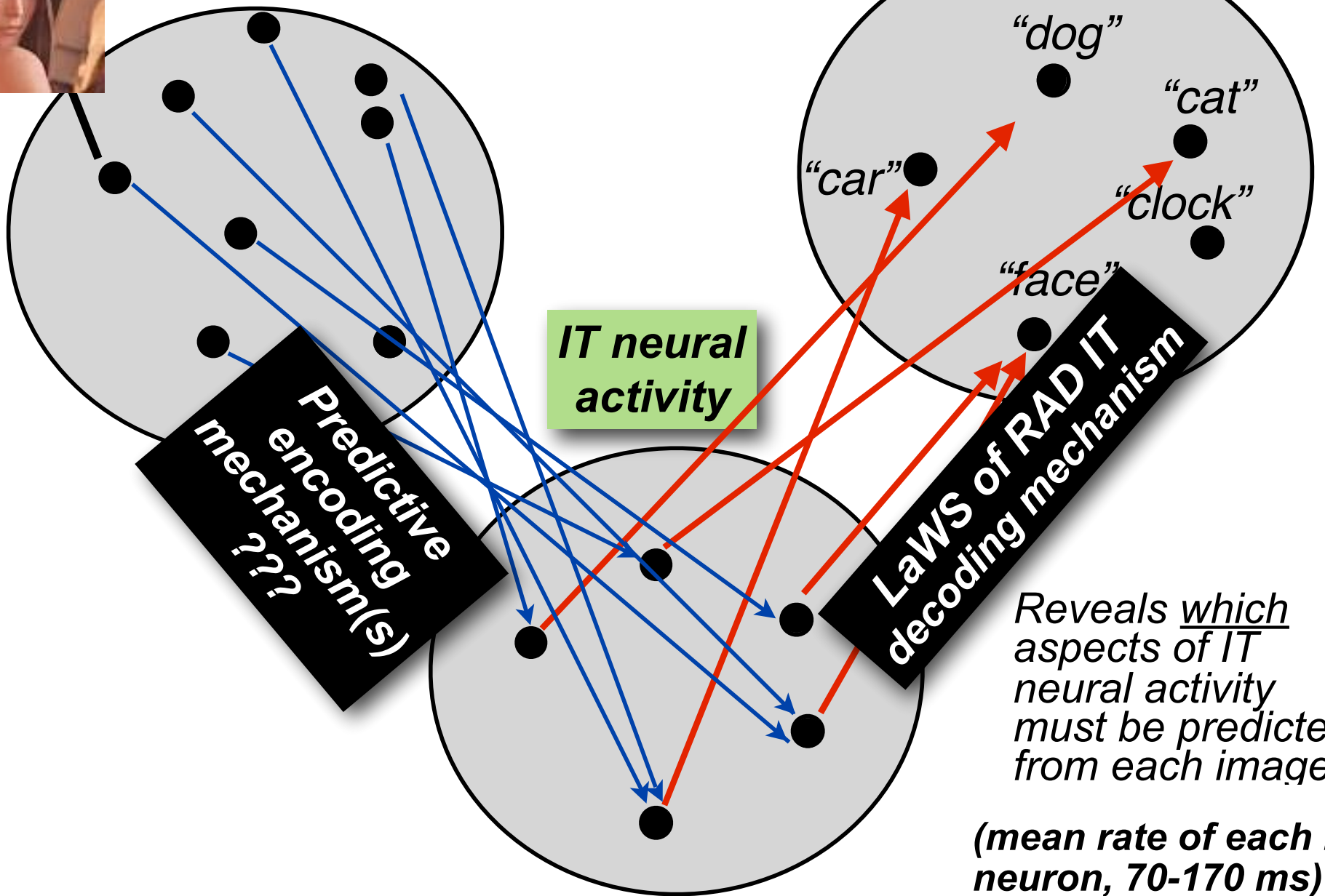
   Note: this must predict behavioral report and it must include a falsifiable statement of the relevant aspects of neural activity (aka "neural code")

**2. Can we infer the encoding mechanism(s) that accurately predict the relevant ventral stream population patterns of neural activity from each image?**
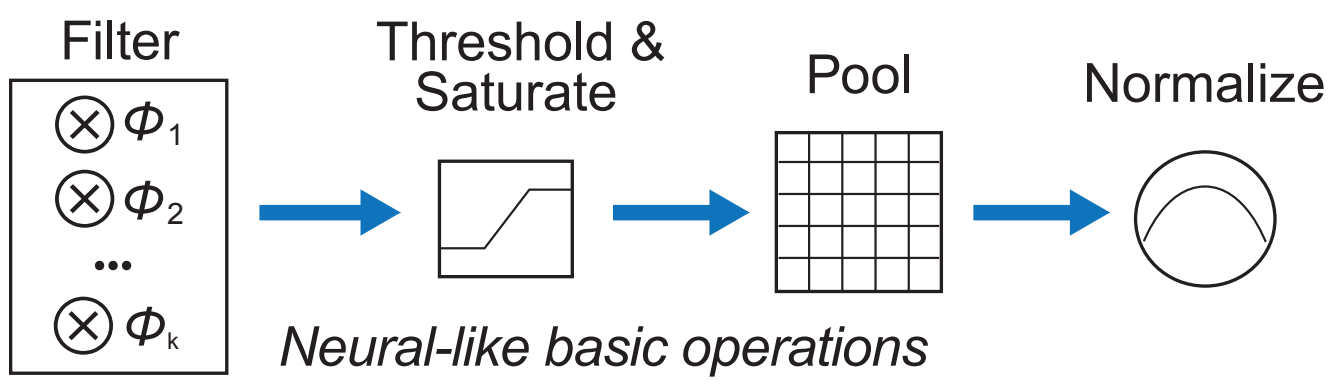
**Images**

**Behavioral reports ("perception")**

"dog"

"cat"

"car"

"clock"

"face"

**IT neural activity**

Predictive encoding mechanism(s) ???

LaWS of RAD IT decoding mechanism

Reveals _which_ aspects of IT neural activity must be predicted from each image

**(mean rate of each IT neuron, 70-170 ms)**

**"Deep convolutional neural networks" (Deep CNN's)**

*Basic operations:* $\Theta = (\theta_{filter}, \theta_{thr}, \theta_{sat}, \theta_{pool}, \theta_{norm})$

Filter $\qquad$ Threshold & Saturate $\qquad$ Pool $\qquad$ Normalize

$\otimes \Phi_1$
$\otimes \Phi_2$
...
$\otimes \Phi_k$

*Neural-like basic operations*

$\Theta^{(1)} \qquad \Theta^{(2)} \qquad \Theta^{(3)}$

Layer 1 $\qquad$ Layer 2 $\qquad$ Layer 3

**Elements** ("neurons") **have large fan-in**

**Simple, bio-known non-linearities**

**Each layer:**
 **is convolutional** (i.e. retinotopy)

 **has many types of tuning functions**

**Deep stack of layers**

Top layer has thousands of visual "neurons"

*Pinto, Doukan, DiCarlo & Cox,* **PLoS Comp Biol** *(2009)*
*Hubel & Wiesel (1962), Fukushima (1980); Perrett & Oram (1993); Wallis & Rolls (1997); LeCun et al. (1998); Riesenhuber & Poggio (1999); Serre, Kouh, et al. (2005), etc....*

100

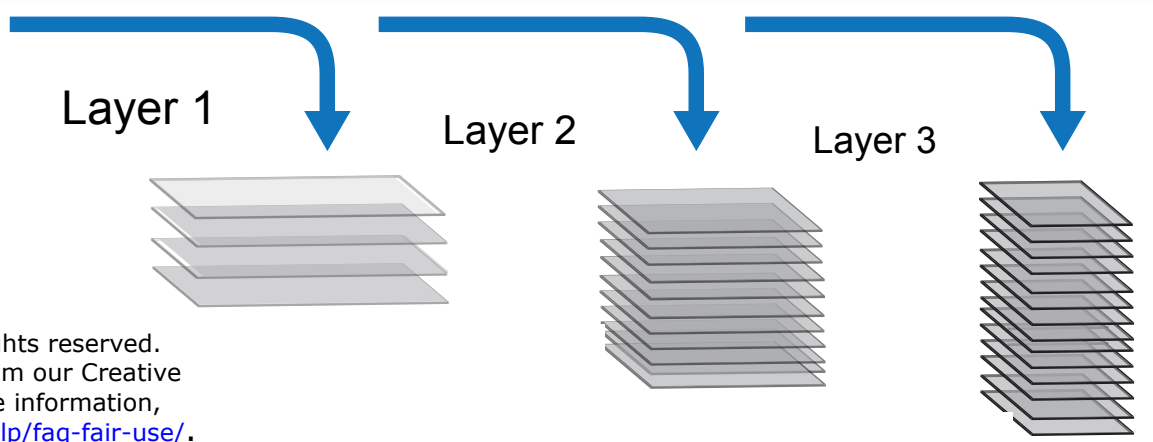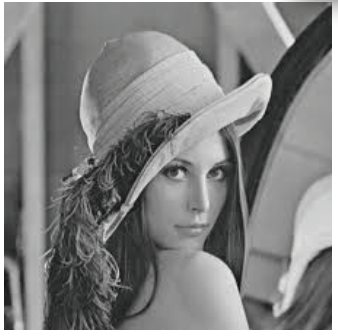**Our goal (2008): explore a family of possible encoding mechanisms**

**"Deep convolutional neural networks" (Deep CNN's)**

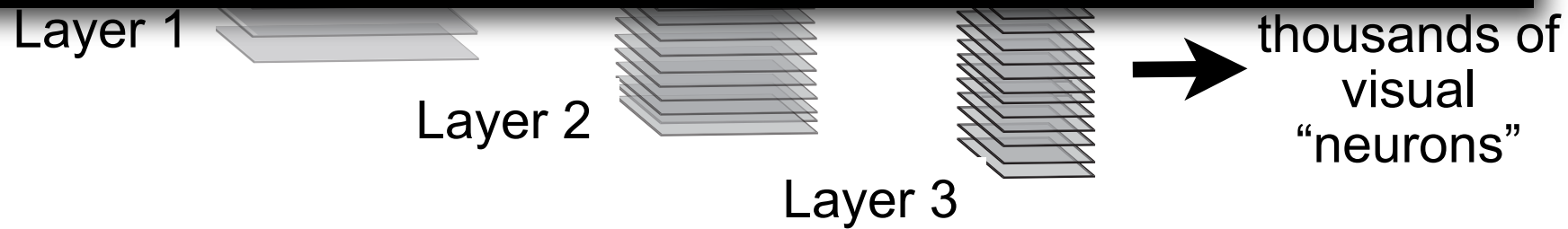*Basic operations:* $\Theta = (\theta_{filter}, \theta_{thr}, \theta_{sat}, \theta_{pool}, \theta_{norm})$

**Thousands of unknown parameters**

(i.e. not directly determined by neurobiology)

Filter

$\otimes \phi_1$
$\otimes \phi_2$
...
$\otimes \phi_k$

Threshold & Saturate

Pool

Normalize

*Neural-like basic operations*

**Set all parameters —> gives a model**

**That model PREDICTS the entire neural population response to ANY image, in each successive visual area**

Layer 1
Layer 2
Layer 3

Top layer has thousands of visual "neurons"

*Pinto, Doukan, DiCarlo & Cox, **PLoS Comp Biol** (2009)*

*Hubel & Wiesel (1962), Fukushima (1980); Perrett & Oram (1993); Wallis & Rolls (1997); LeCun et al. (1998); Riesenhuber & Poggio (1999); Serre, Kouh, et al. (2005), etc....*

**"Deep convolutional neural networks" (Deep CNN's)**

*Basic operations:* $\Theta = (\theta_{filter}, \theta_{thr}, \theta_{sat}, \theta_{pool}, \theta_{norm})$

**Thousands of unknown parameters**

(i.e. not directly determined by neurobiology)

Filter

$\otimes \Phi_1$
$\otimes \Phi_2$
...

Threshold & Saturate

Pool

Normalize

**How do we determine which of these models, if any, is a model of the ventral stream?**

**1. Use optimization methods to find specific models (i.e. parameter settings) in this model family.**

**2. Optimization target = visual tasks that we hypothesize that the ventral stream evolved and/or developed to solve.**

Layer 1

Layer 2

Layer 3

thousands of visual "neurons"

*Hubel & Wiesel (1962), Fukushima (1980); Perrett & Oram (1993); Wallis & Rolls (1997); LeCun et al. (1998); Riesenhuber & Poggio (1999); Serre, Kouh, et al. (2005), etc....*

*Yamins, Hong, Solomon, Seibert and DiCarlo **PNAS (2014)***

▸ **variety of 3D objects (36)** with semantic breadth (e.g. not all faces)

▸ rendered with large amount of **variation**

▸ These are **different objects** that those we will use later in testing

Nine example objects:



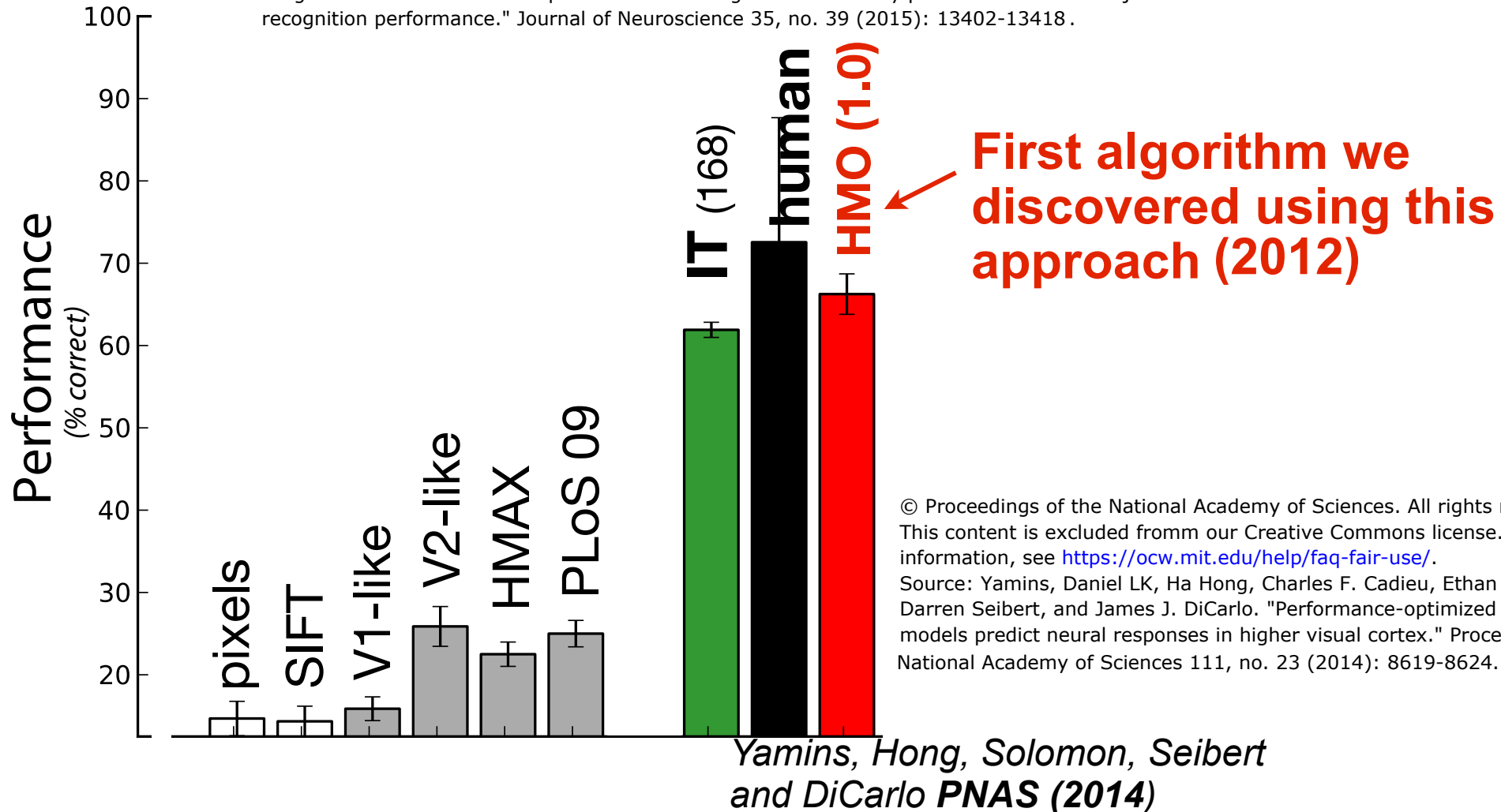Bodies, Buildings, Flowers, Guns, Instruments, Jewelry, Shoes, Tools, Trees

Test on Core Object Recognition 1.0

**First algorithm we discovered using this approach (2012)**

*Yamins, Hong, Solomon, Seibert and DiCarlo* **PNAS (2014)**

**Basic operations:** $\Theta = (\theta_{filter}, \theta_{thr}, \theta_{sat}, \theta_{pool}, \theta_{norm})$

Filter

$\otimes \Phi_1$
$\otimes \Phi_2$
...
$\otimes \Phi_k$

Threshold & Saturate

Pool

Normalize

*Neural-like basic operations*

$\Theta^{(1)}$  $\Theta^{(2)}$  $\Theta^{(3)}$  $\Theta$

**Model layer 1**   **Model layer 2**   **Model layer 3**   **Model layer 4**

**Cross-validated linear regression**   **Predict IT?**

*These are PREDICTIONS: All of these objects and images were never previously seen by the HMO model*

d



Unit 1: $r^2 = 0.48$

**Response\* of IT neural site**

Animals    Boats    Cars    **Chairs**    Faces    Fruits    Planes    Tables

**Prediction of HMO model**

Source: Yamins, Daniel LK, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and James J. DiCarlo. "Performance-optimized hierarchical models predict neural responses in higher visual cortex." Proceedings of the National Academy of Sciences 111, no. 23 (2014): 8619-8624.

**(\* mean rate 70-170 ms after image onset)**

*Yamins, Hong, Solomon, Seibert and DiCarlo PNAS (2014)*

106

# Predictions of single site IT responses from layer 4 of HMO 1.0 model

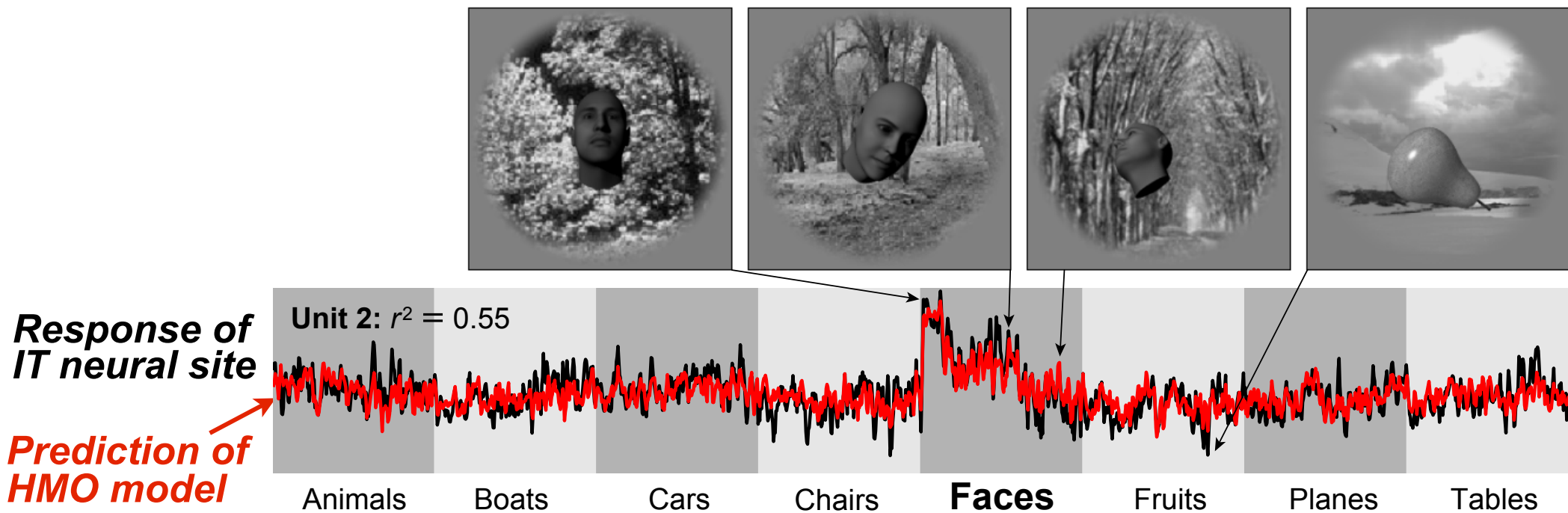*These are PREDICTIONS:  All of these objects and images were never previously seen by the HMO model*



**Response of IT neural site**
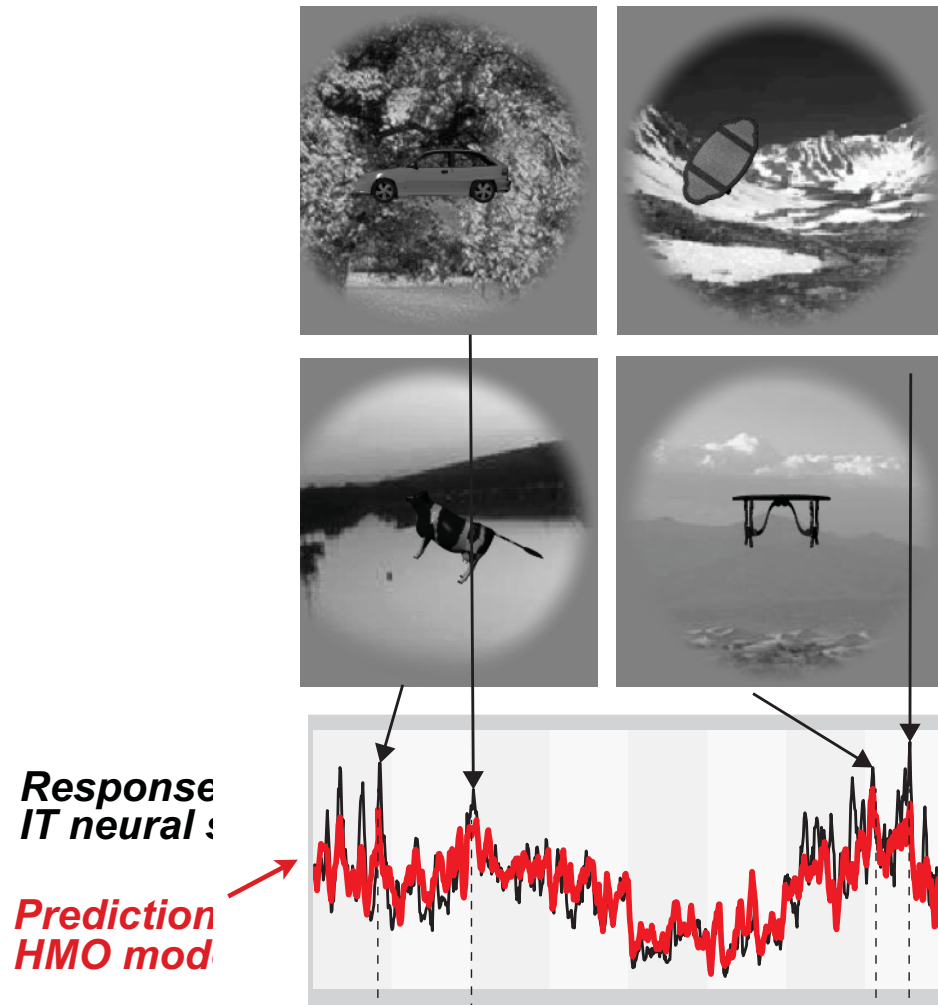
Unit 2: $r^2 = 0.55$

**Prediction of HMO model**

Animals    Boats    Cars    Chairs    **Faces**    Fruits    Planes    Tables

**(* mean rate 70-170 ms after image onset)**

*Yamins, Hong, Solomon, Seibert and DiCarlo* **PNAS (2014)**

## IT Site 42



**Response** ...
**IT neural** ...

**Prediction** ...
**HMO mod** ...

Source: Yamins, Daniel LK, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and James J. DiCarlo. "Performance-optimized hierarchical models predict neural responses in higher visual cortex." Proceedings of the National Academy of Sciences 111, no. 23 (2014): 8619-8624.

*Yamins, Hong, Solomon, Seibert and DiCarlo **PNAS (2014)***

# Ability of various encoding mechanisms (specific models) to predict IT responses to naturalistic images



**~50% of IT single unit response variance predicted. Dramatic improvement over previous models.**

Source: Yamins, Daniel LK, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and James J. DiCarlo. "Performance-optimized hierarchical models predict neural responses in higher visual cortex." Proceedings of the National Academy of Sciences 111, no. 23 (2014): 8619-8624.

*Yamins, Hong, Solomon, Seibert and DiCarlo PNAS (2014)*

Basic operations: $\Theta = (\theta_{filter}, \theta_{thr}, \theta_{sat}, \theta_{pool}, \theta_{norm})$

Filter — Threshold & Saturate — Pool — Normalize

$\otimes \phi_1$
$\otimes \phi_2$
...
$\otimes \phi_k$

Neural-like basic operations

Source: Yamins, Daniel LK, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and James J. DiCarlo. "Peformance-optimized hierarchical models predict neural responses in higher visual cortex." Proceedings of the National Academy of Sciences 111, no. 23 (2014): 8619-8624.
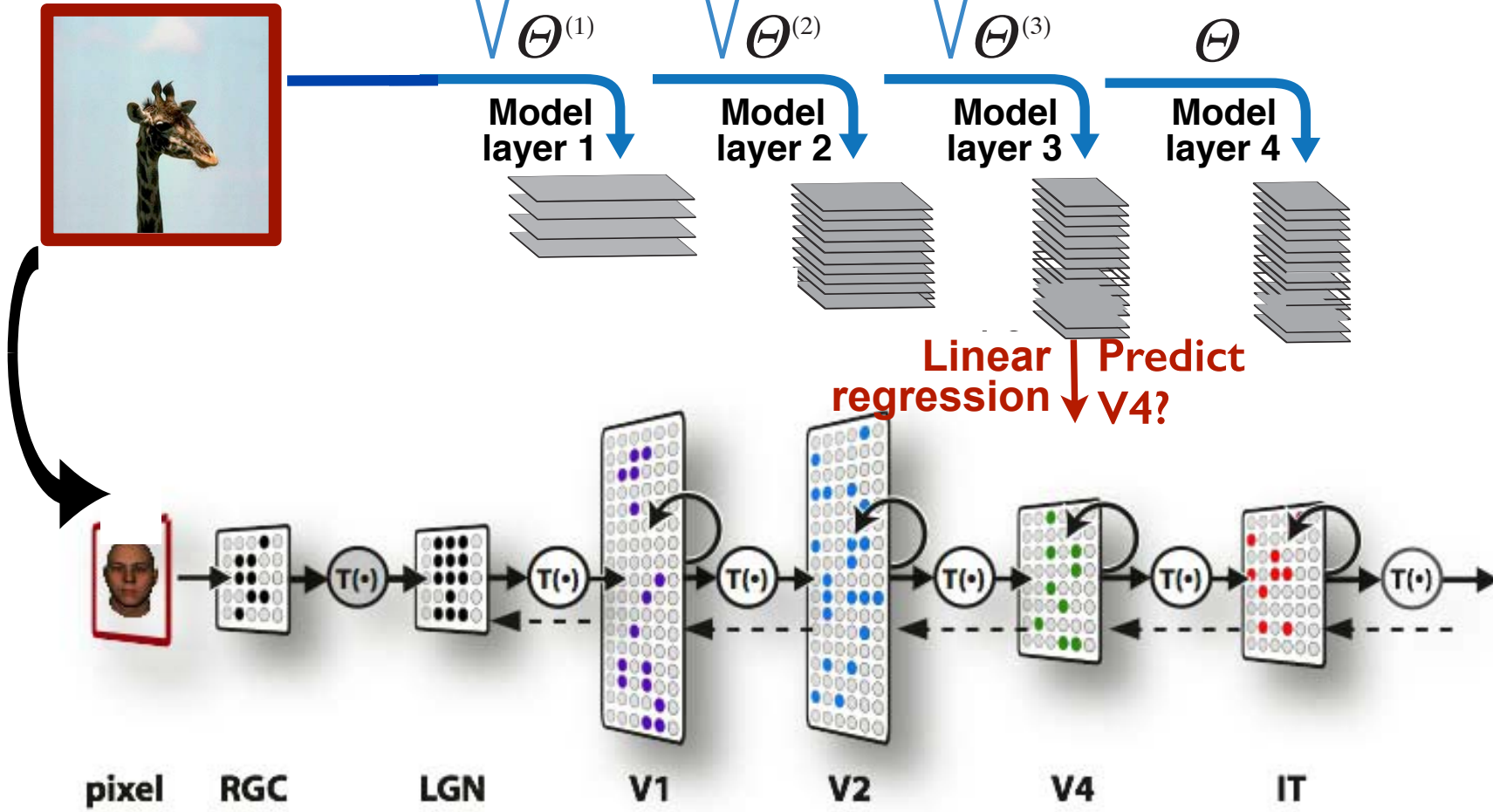
HM0 1.0
(all parameters fixed)

$\Theta^{(1)}$  $\Theta^{(2)}$  $\Theta^{(3)}$  $\Theta$

Model layer 1   Model layer 2   Model layer 3   Model layer 4

Linear regression    Predict V4?

pixel   RGC   LGN   V1   V2   V4   IT

$T(\bullet)$

Source: DiCarlo, James J., and David D. Cox. "Untangling invariant object recognition." Trends in cognitive sciences 11, no. 8 (2007): 333-341.

1F€

**Bio-inspired algorithm class + tasks in domain + optimization ==> neural-like encoding functions!**

**Even in intermediate layers!**

## V4 predictive power
(median over all neurons)

## IT predictive power
(median over all neurons)

*Yamins, Hong, Solomon, Seibert and DiCarlo **PNAS (2014)***

1FF

**Explanatory power of HMO model**

**Current maximum possible\* explanatory power**

HMO Model

Monkey IT

Source: Yamins, Daniel LK, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and James J. DiCarlo. "Performance-optimized hierarchical models predict neural responses in higher visual cortex." Proceedings of the National Academy of Sciences 111, no. 23 (2014): 8619-8624.
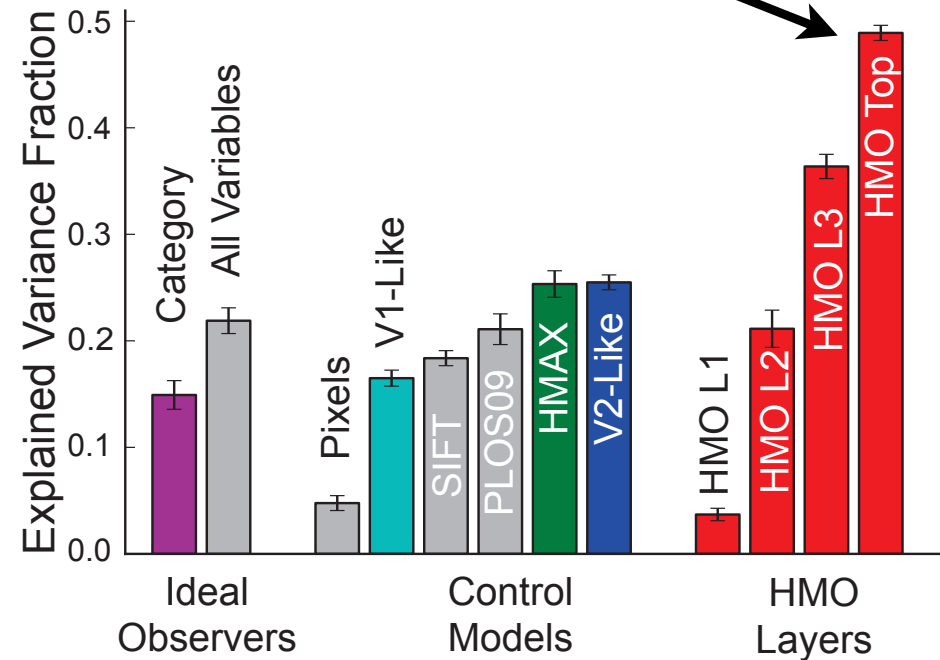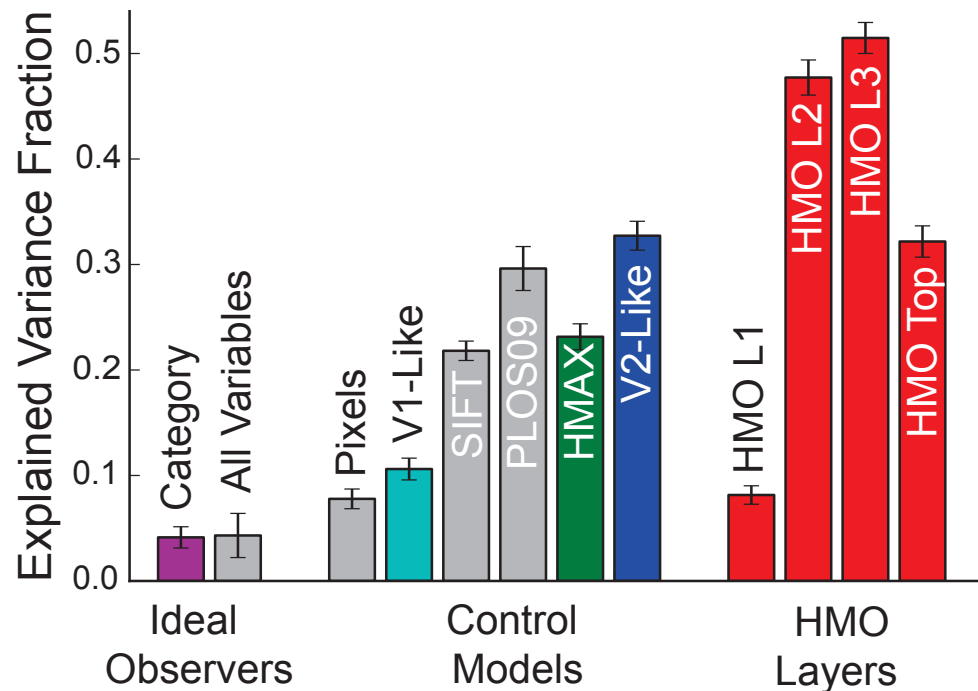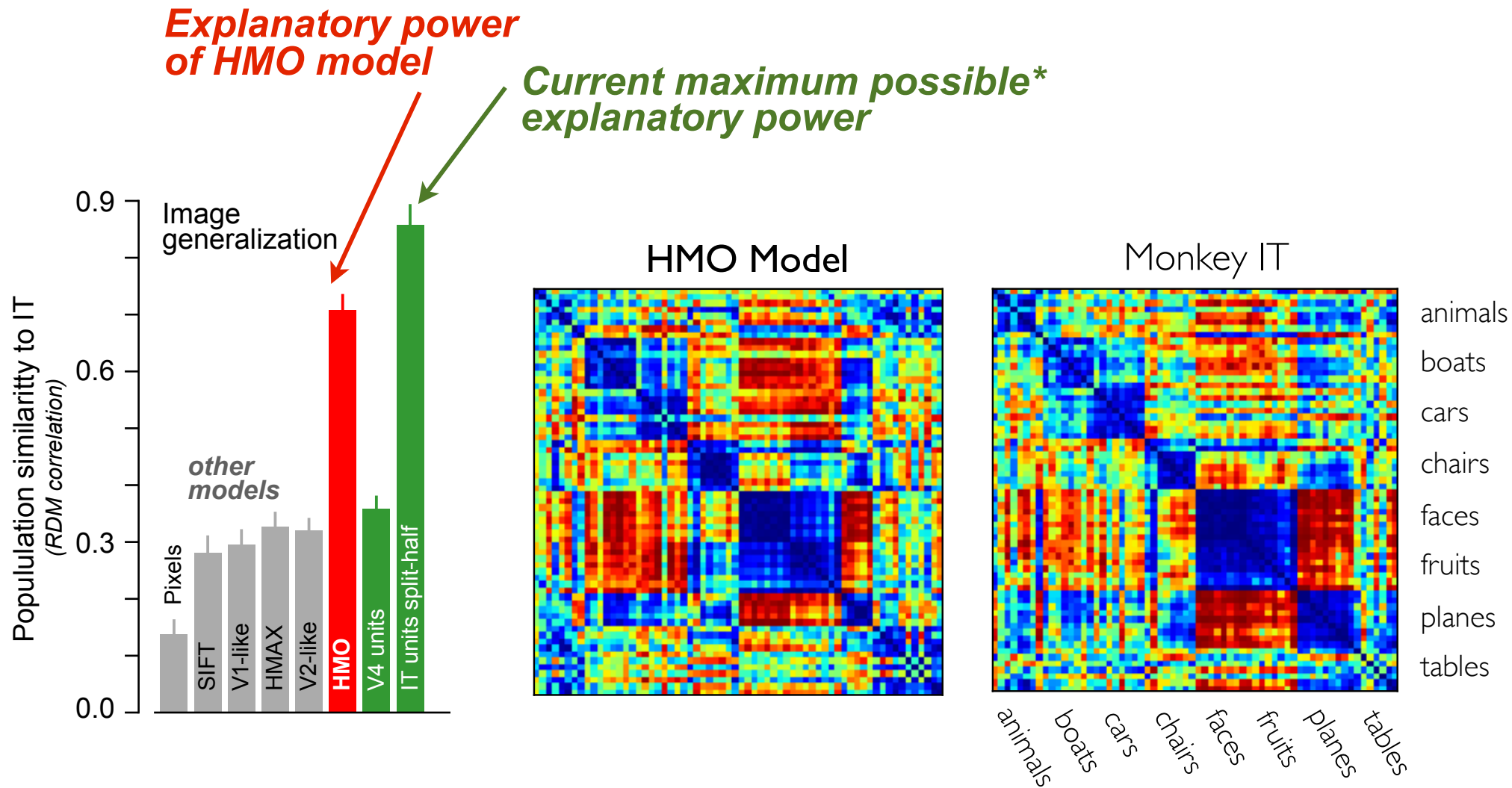
*Yamins, Hong, Solomon, Seibert and DiCarlo* **PNAS (2014)**

112

# Suggests that continued optimization within this family of models would lead to even higher neural predictive power.

Source: Yamins, Daniel LK, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and James J. DiCarlo. "Performance-optimized hierarchical models predict neural responses in higher visual cortex." Proceedings of the National Academy of Sciences 111, no. 23 (2014): 8619-8624.

# Suggests that continued optimization within this family of models would lead to even higher neural predictive power.
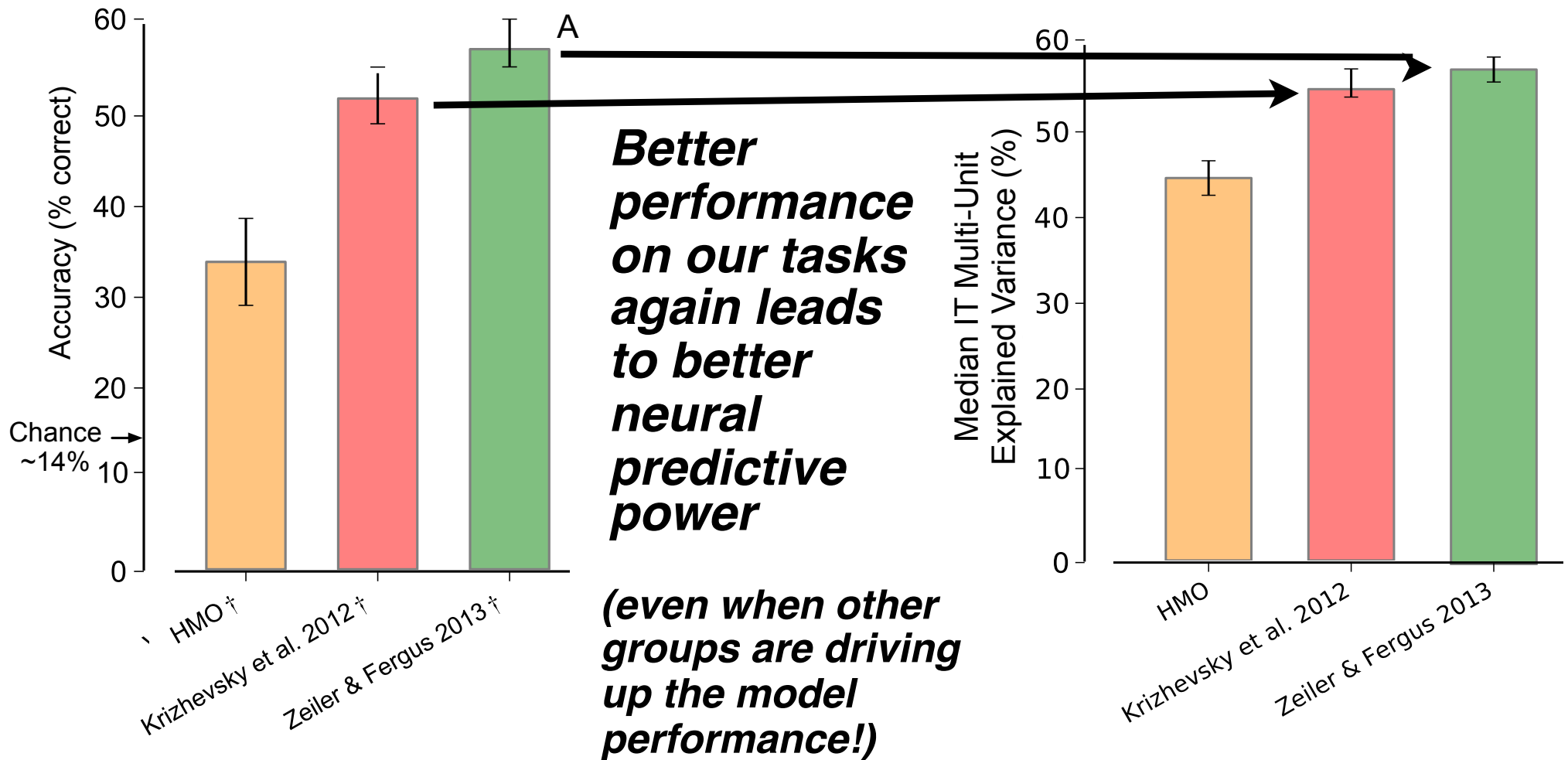


**Better performance on our tasks again leads to better neural predictive power**

*(even when other groups are driving up the model performance!)*

Cadieu, Charles F., Ha Hong, Daniel LK Yamins, Nicolas Pinto, Diego Ardila, Ethan A. Solomon, Najib J. Majaj, and James J. DiCarlo. "Deep neural networks rival the representation of primate IT cortex for core visual object recognition. "PLoS Comput Biol 10, no. 12 (2014): e1003963; https://doi.org/10.1371/journal.pcbi.1003963. License CC BY.

*Cadieu CF, Hong H, Yamins D, Pinto N, Majaj N, and DiCarlo JJ. **ICLR** (2013);*
*Cadieu CF, Hong H, Yamins D, Pinto N, Majaj N, and DiCarlo JJ. **PLoS Comp Bio** (2014)*

# CNN features vs. IT "features"

a)



Cars        Fruits        Animals     Planes Chairs Tables Faces

b)



Retinae Representation    Ventral Stream    IT Cortex Representation

Pixel Representation    Deep Neural Network (DNN) $\phi(x)$    DNN Representation

○ Cars    ● Fruits

Cadieu, Charles F., Ha Hong, Daniel LK Yamins, Nicolas Pinto, Diego Ardila, Ethan A. Solomon, Najib J. Majaj, and James J. DiCarlo. "Deep neural networks rival the representation of primate IT cortex for core visual object recognition. "PLoS Comput Biol 10, no. 12 (2014): e1003963; https://doi.org/10.1371/journal.pcbi.1003963. License CC BY.

Cadieu CF, Hong H, Yamins D, Pinto N, Majaj N, and DiCarlo JJ. **ICLR** (2013);
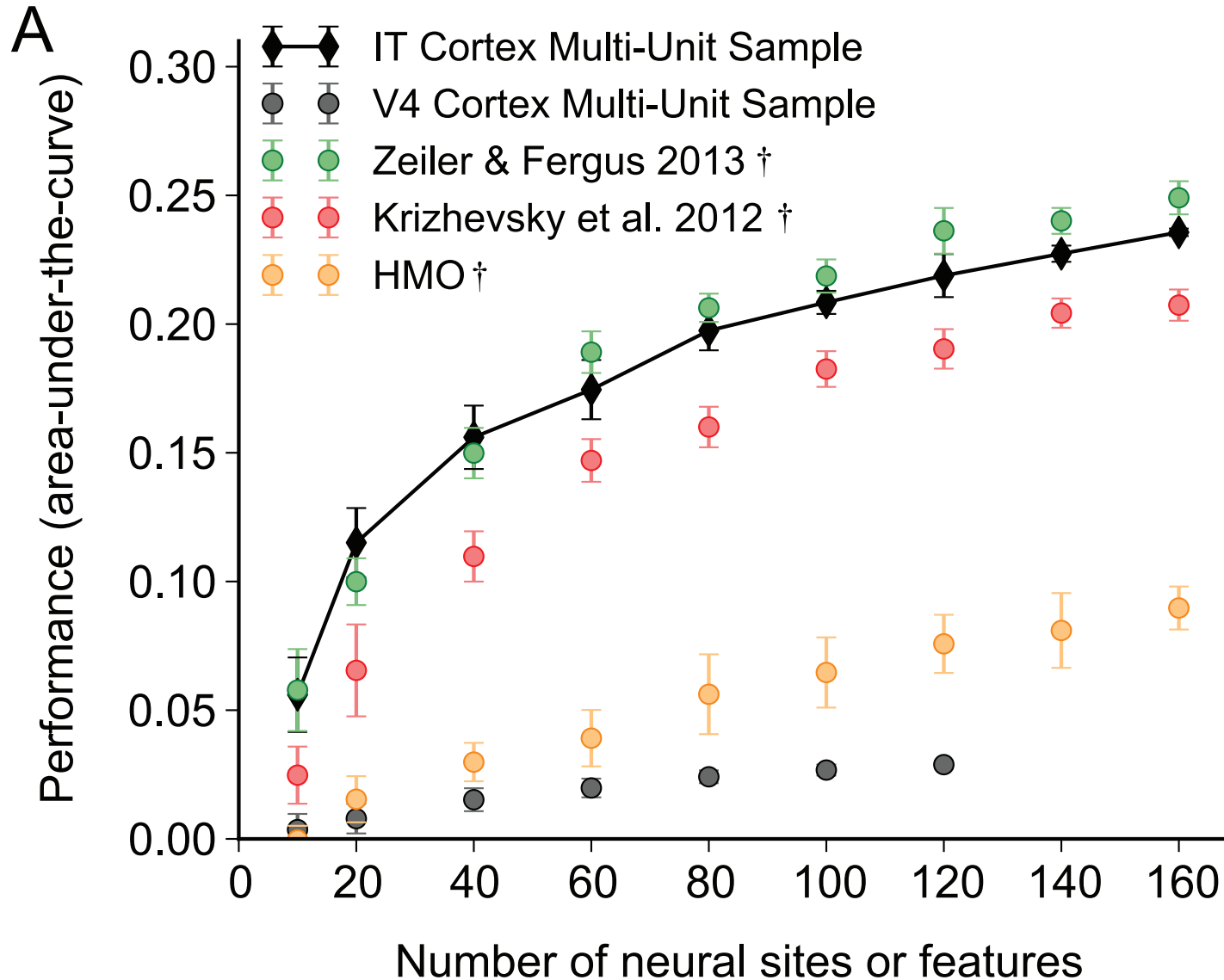Cadieu CF, Hong H, Yamins D, Pinto N, Majaj N, and DiCarlo JJ. **PLoS Comp Bio** (2014)

Cadieu, Charles F., Ha Hong, Daniel LK Yamins, Nicolas Pinto, Diego Ardila, Ethan A. Solomon, Najib J. Majaj, and James J. DiCarlo. "Deep neural networks rival the representation of primate IT cortex for core visual object recognition. "PLoS Comput Biol 10, no. 12 (2014): e1003963; https://doi.org/10.1371/journal.pcbi.1003963. License CC BY.

*Cadieu CF, Hong H, Yamins D, Pinto N, Majaj N, and DiCarlo JJ. **ICLR** (2013);*
*Cadieu CF, Hong H, Yamins D, Pinto N, Majaj N, and DiCarlo JJ. **PLoS Comp Bio** (2014)*

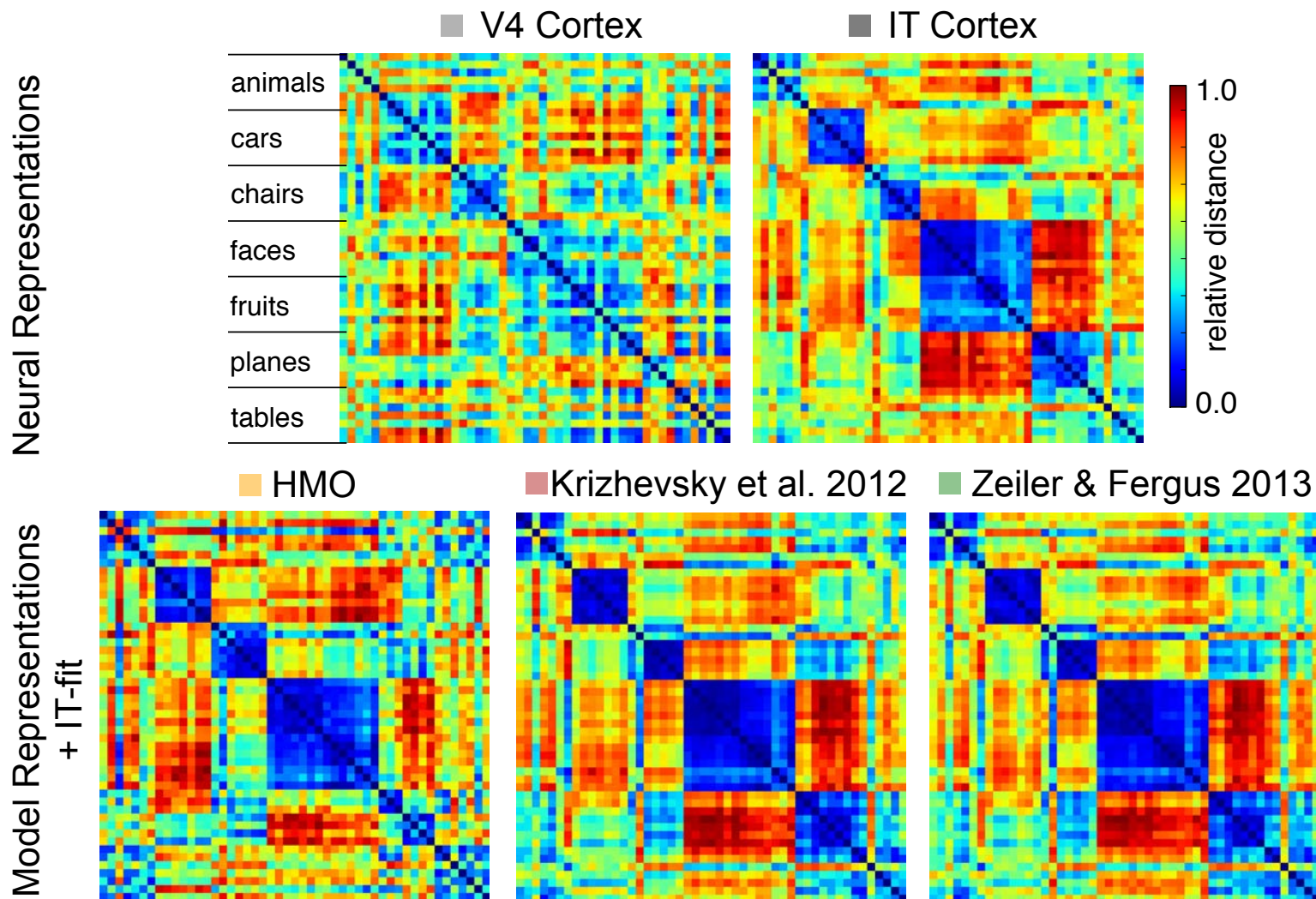# Better performing deep CNN networks also better predict the patterns of IT neural responses



Cadieu, Charles F., Ha Hong, Daniel LK Yamins, Nicolas Pinto, Diego Ardila, Ethan A. Solomon, Najib J. Majaj, and James J. DiCarlo. "Deep neural networks rival the representation of primate IT cortex for core visual object recognition. "PLoS Comput Biol 10, no. 12 (2014): e1003963; https://doi.org/10.1371/journal.pcbi.1003963. License CC BY.

*Cadieu CF, Hong H, Yamins D, Pinto N, Majaj N, and DiCarlo JJ. **ICLR** (2013);*
*Cadieu CF, Hong H, Yamins D, Pinto N, Majaj N, and DiCarlo JJ. **PLoS Comp Bio** (2014)*

# Summary of what I presented today (Domain: Core recognition)

**1. Showed that IT firing rates are a feature basis on which learned object judgements naturally predict human/monkey performance; defined parameters.**

**LaWS of RAD IT**
*[70-170ms, 50,000n, 100t]*

*Inference: this might be the specific neural code and decoding mechanism that the brain uses to support these tasks.*

*Systematic causal tests of this model ongoing, but results thus far are as predicted by the model …*

**2. Showed that optimization of deep CNNs (models) for invariant object recognition tasks led to dramatic improvements in our ability to predict IT and V4 neural responses.** **HMO 1.0, CNN 2.0**

*Inference: the encoding mechanisms in these models are similar to those at work in the ventral stream.*

*This is allowing the field to design experiments to explore what remains unique and powerful about primate object perception.*
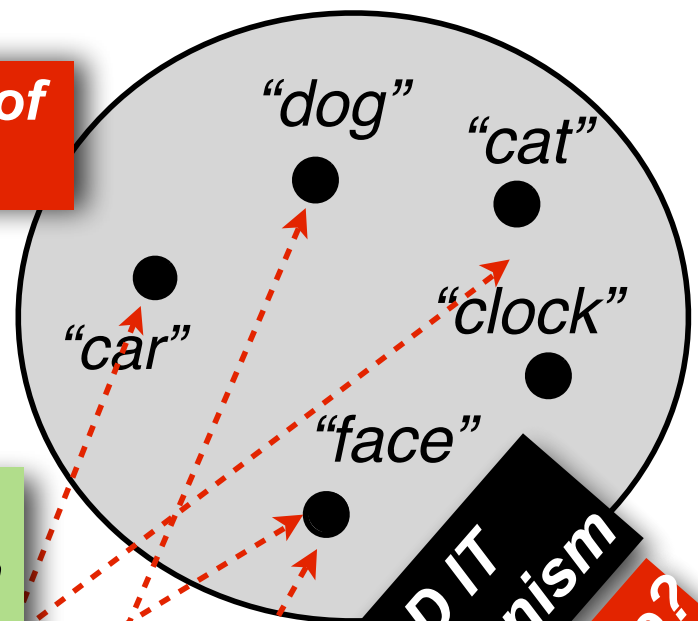
**Behavioral reports ("perception")**

*Images*

*Expand domain of object tasks*

"dog"
"cat"
"car"
"clock"
"face"

*High level ventral stream neural activity (V4, IT)*

HMO encoding mechanism
Other deep CNNS

Learning: Can the next models be less supervised ?

LaWS of RAD IT decoding mechanism

Predict for each image?

Dynamics, feedback?

**Ongoing: Predictable effects of direct neural perturbations of IT?**

# Acknowledgements

MIT OpenCourseWare

Resource: Brains, Minds and Machines Summer Course
Tomaso Poggio and Gabriel Kreiman

The following may not correspond to a particular course on MIT OpenCourseWare, but has been
provided by the author as an individual learning resource.