

A Quest to understand the neural mechanisms of human visual intelligence

James DiCarlo MD, PhD

Peter de Florez Professor of Neuroscience

Director, MIT Quest for Intelligence

Investigator, McGovern Institute for Brain Research

Investigator, Center for Brains, Minds and Machines

Massachusetts Institute of Technology



A Quest to understand the neural mechanisms of human visual intelligence

= “Reverse engineering human visual intelligence”

James DiCarlo MD, PhD

Peter de Florez Professor of Neuroscience

Director, MIT Quest for Intelligence

Investigator, McGovern Institute for Brain Research

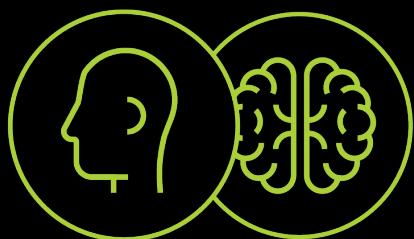
Investigator, Center for Brains, Minds and Machines

Massachusetts Institute of Technology



The neuroscientific goal of reverse engineering:

*Account for human visual intelligence ...
(behavioral capabilities)*

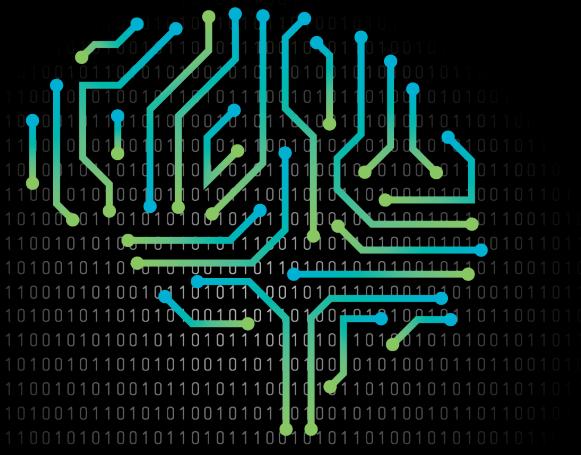


SCIENCE

Mind & Brain

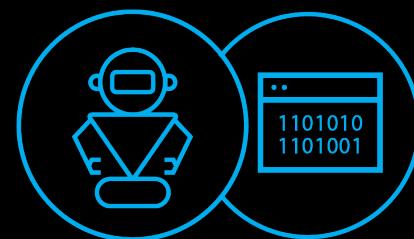
Measurements
& Discoveries

*...using mechanisms
of the brain...
(networks of simulated neurons)*



Specific artificial neural networks as implemented (and leading) scientific hypotheses

...in the language of engineering (predictive, built systems).



ENGINEERING

Software,
Hardware, Robotics

Synthesis & Creation

My talk today: Ongoing progress on a foundational piece of visual intelligence

Local reverse engineering team:



brain+cognitive
sciences



Current group:

Robert Ajemian
Yoon Bai
Joel Dapello
Kohitij Kar
Micheal Lee
Tiago Marques
Ratan Murty
Alina Peter
Jon Prescott-Roy
Sachi Sanghavi
Martin Schrimpf
Christopher Shay
Chris Stawarz



Alumni scientists:

Arash Afraz (=> Prof., NIH)
Paul Aparicio (=> NIH)
Pouya Bashivan (=> Prof., McGill)
Charles Cadieu (=> Caption Health)
David Cox (=> IBM, VP AI research)
Ha Hong (=> Caption Health)
Chou Hung (=> Army research)

Elias Issa (=> Prof., Columbia)
Xiaoaxuan Jia (=> Allen Institute)
Kamila Jozwik (=> U. Cambridge)
Gabriel Kreiman (=> Prof., Harvard)
Hyodong Lee (=> Google)
Nuo Li (=> Prof., Baylor College Med)
Najib Majaj (=> NYU)

Shay Ohayon (=> Google)
Nicolas Pinto (=> Apple => Cygni)
Rishi Rajalingham (=> MIT)
Nicole Rust (=> Prof., U Penn)
Dan Yamins (=> Prof., Stanford)
Davide Zoccolan (=> Prof., SISSA)

- Office of Naval Research
- National Science Foundation (CBMM)
- Simons Global Brain
- IBM/MIT Watson Lab
- Semiconductor Research Corporation (SRC)/DARPA

Key collaborators:

Ed Boyden (MIT)
SueYeon Chung (Columbia)
David Cox (IBM)
Danny Gutfreund (IBM)
Nancy Kanwisher (MIT)
Lynne Kiorpes (NYU)
Fei-Fei Li (Stanford)
Jitendra Malik (UC Berkeley)
J. Anthony Movshon (NYU)
Tomaso Poggio (MIT)
Kaushik Roy (Purdue)
Josh Tenenbaum (MIT)
Andreas Tolias (Baylor CM)
Dan Yamins (Stanford)

Human visual intelligence...



Guidance from brain and cognitive sciences: ~10 deg at center of gaze, ~200 msec snapshots



Guidance from brain and cognitive sciences: ~10 deg at center of gaze, ~200 msec snapshots



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

Guidance from brain and cognitive sciences: ~10 deg at center of gaze, ~200 msec snapshots

***Foundational component of visual intelligence:
Core object recognition***

Is a car here?

Is a person here?

What is the pose of the car?

...



~200 msec

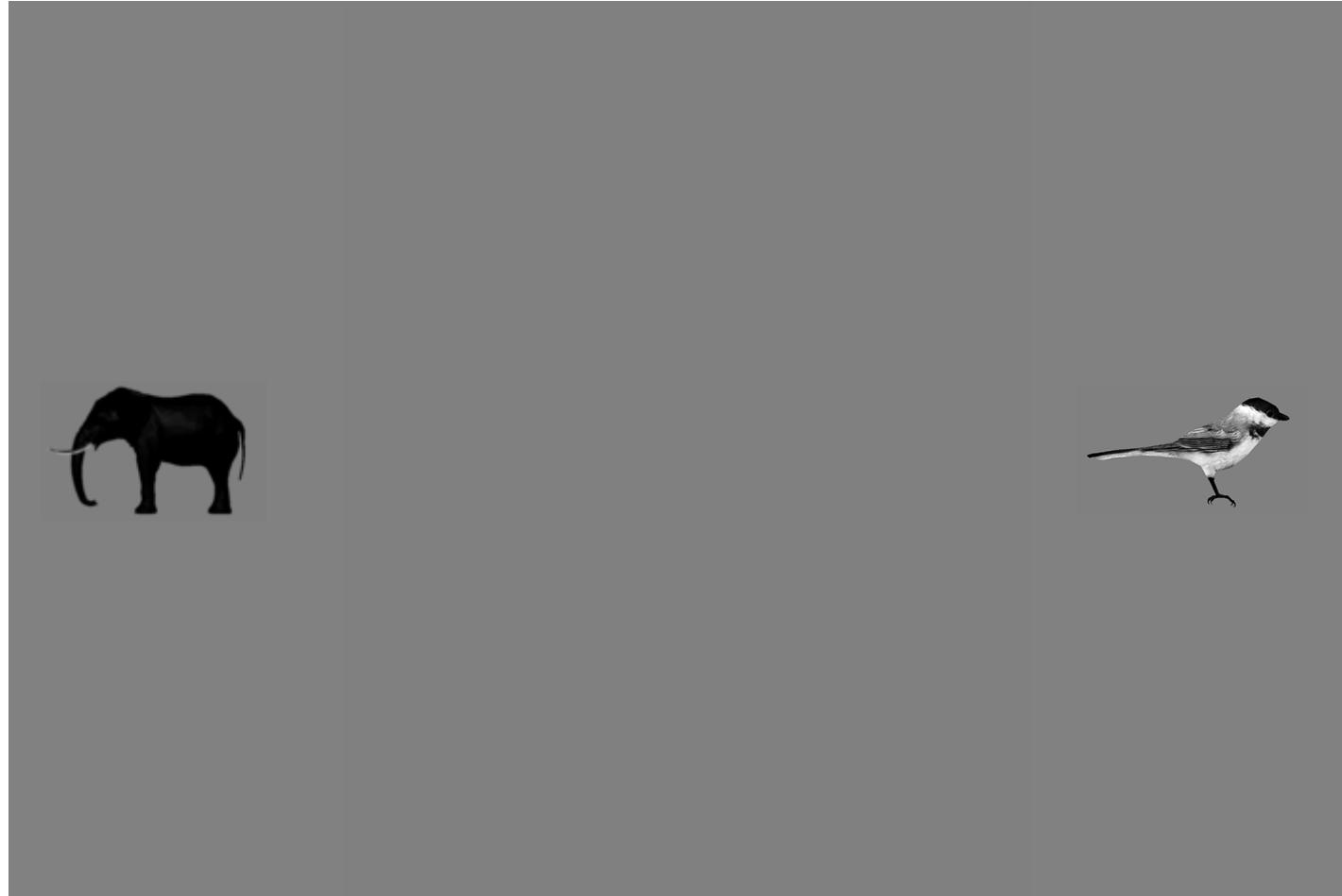


~200 msec



~200 msec

Example behavioral test trials



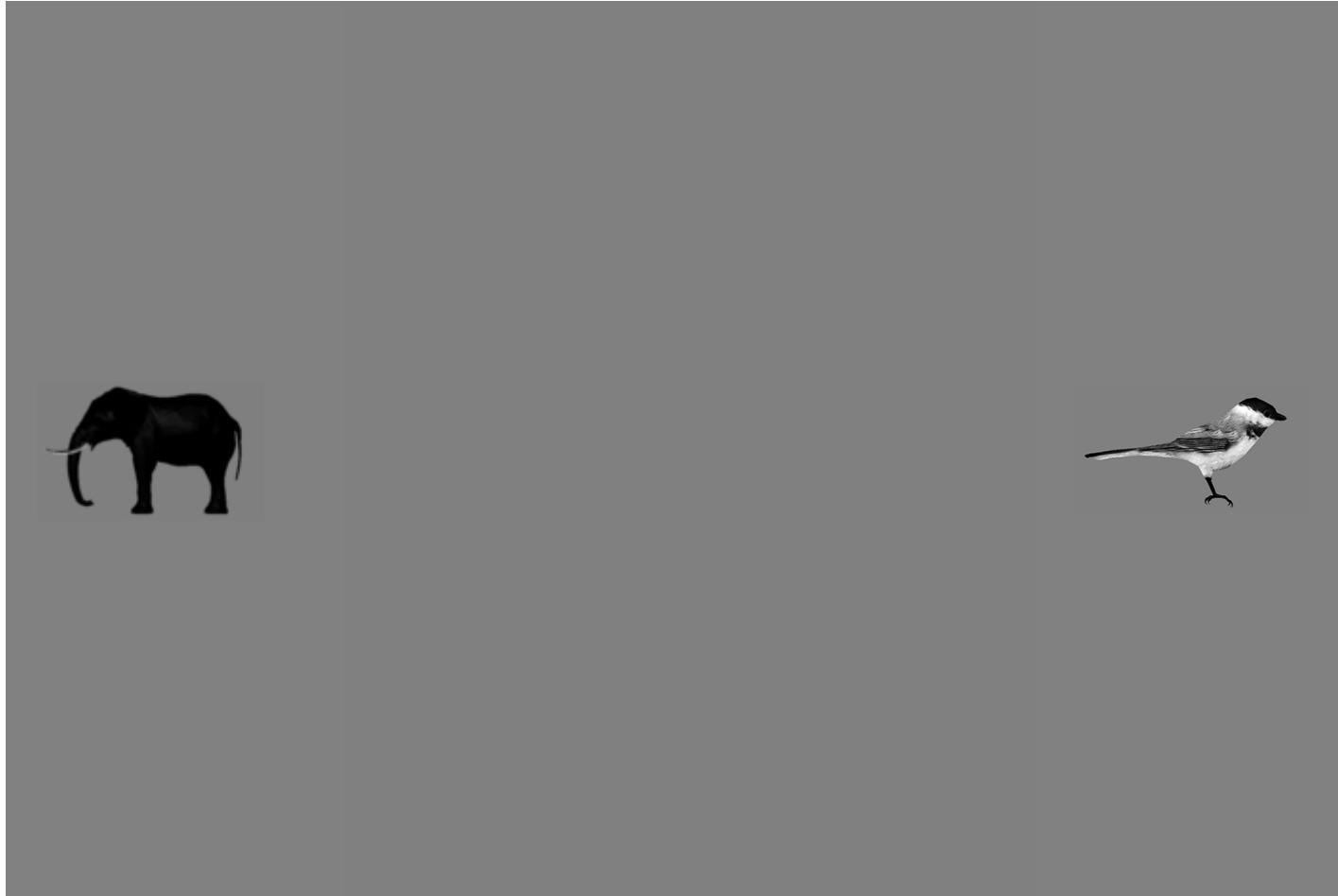
8 deg image at center of gaze, ~100 msec viewing time

Example behavioral test trials



8 deg image at center of gaze, ~100 msec viewing time

Example behavioral test trials



8 deg image at center of gaze, ~100 msec viewing time

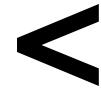
Species B



Species A

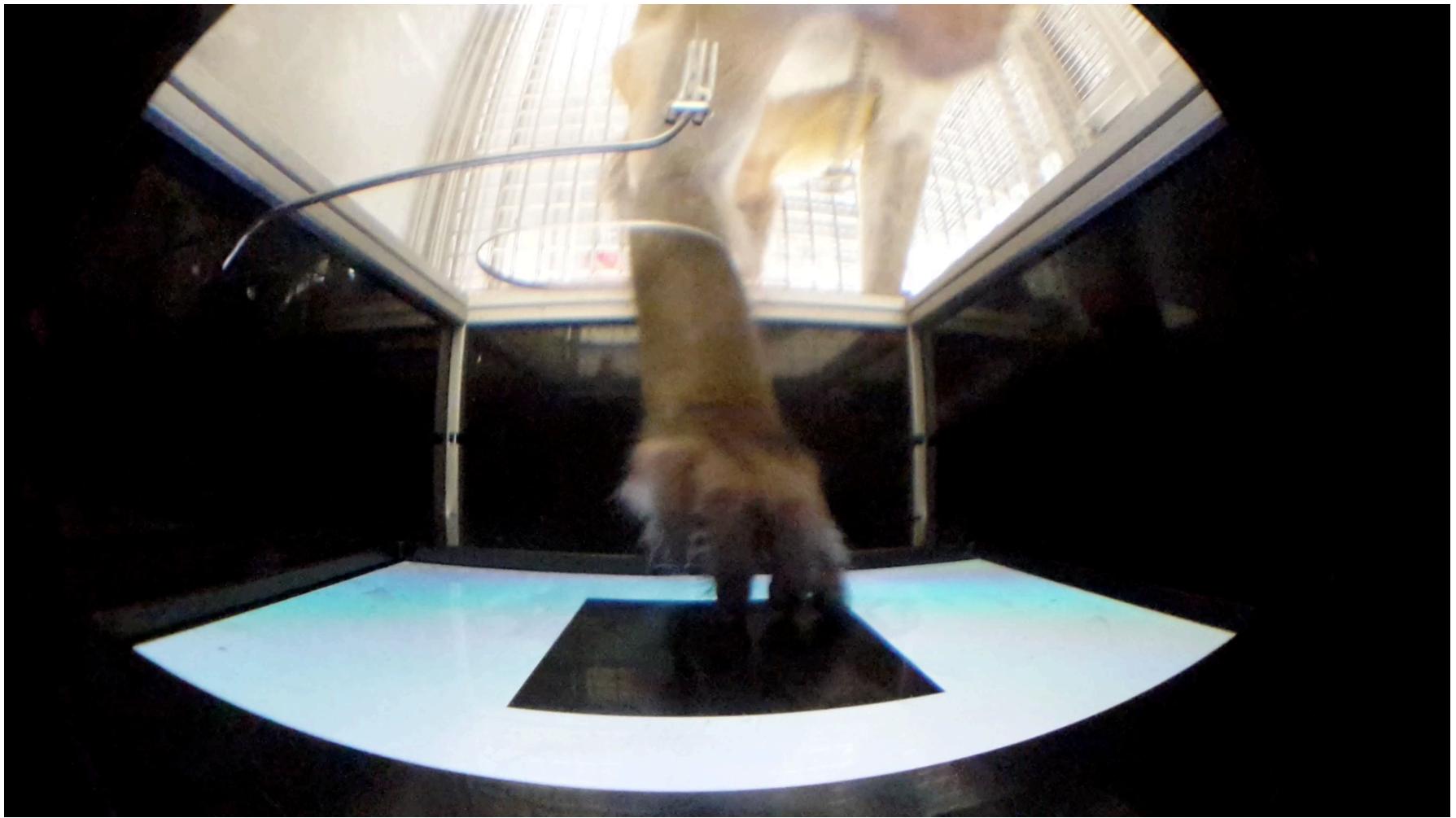


**Computer
vision systems**



**primate,
*Homo sapien***

Intelligence test domain: Core Visual Object Perception

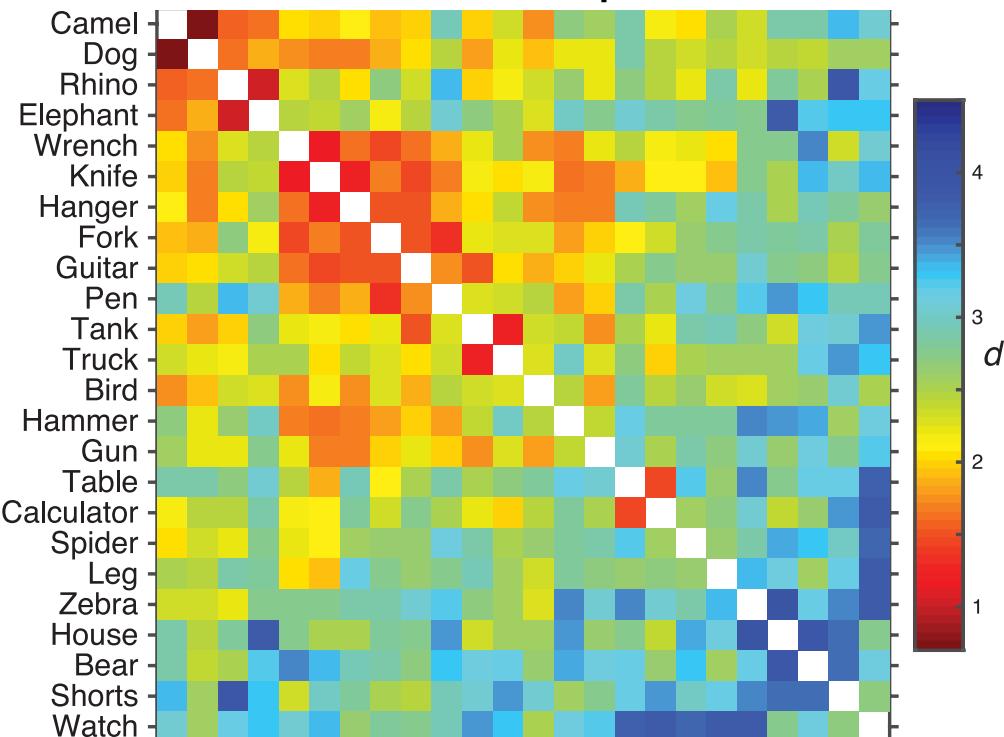


Intelligence test domain: Core Visual Object Perception

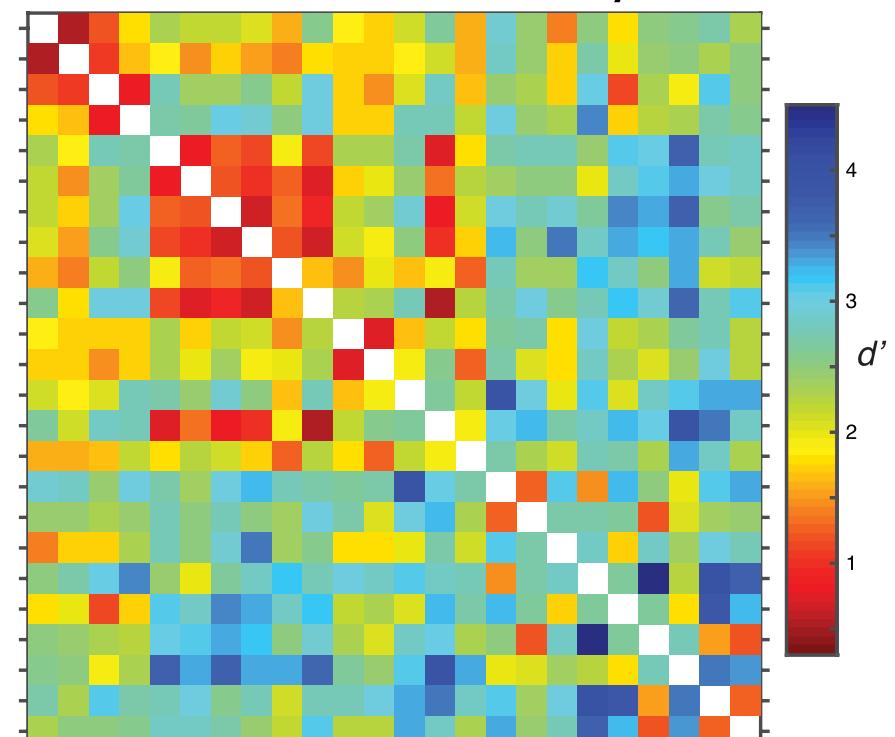
Behavior test performance for 276 core object recognition tasks

(note: many images in each task)

primate,
Homo sapien



primate,
rhesus monkey



Rajalingham, Schmidt, & DiCarlo, *Vision Sciences Society* (2014)

Rajalingham, Schmidt, & DiCarlo, *J. Neuroscience* (2015)

Rajalingham, Issa, Kar, Schmidt, & DiCarlo, *CCN* (2017)

Intelligence test domain: Core Visual Object Perception

Want-to-be
primates



**Computer
vision systems**



**primate,
Homo sapien**

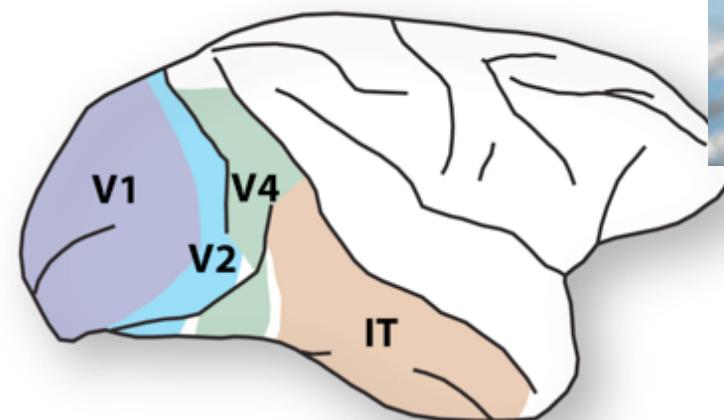
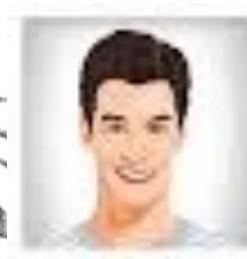
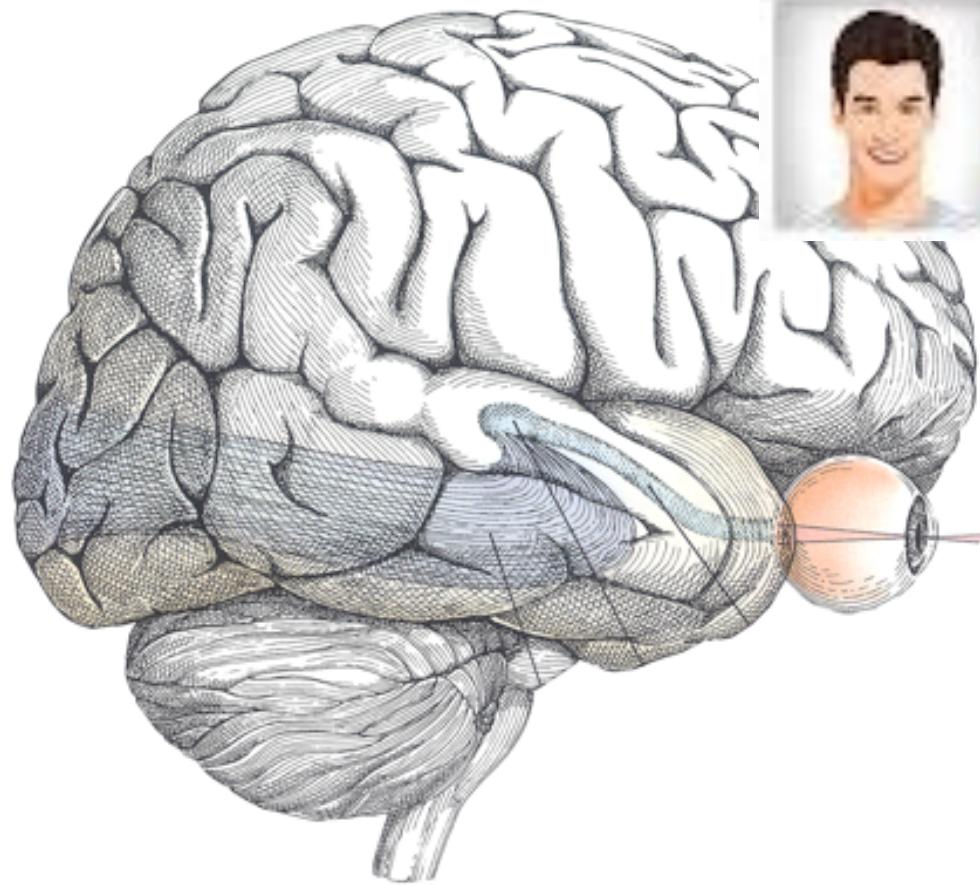


**primate,
rhesus monkey**

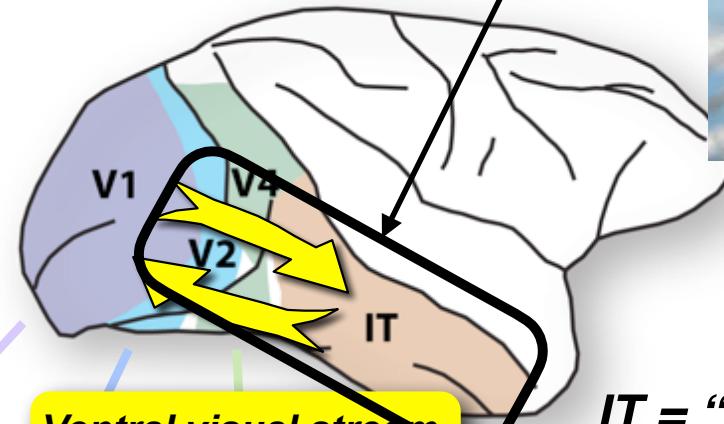
Primates



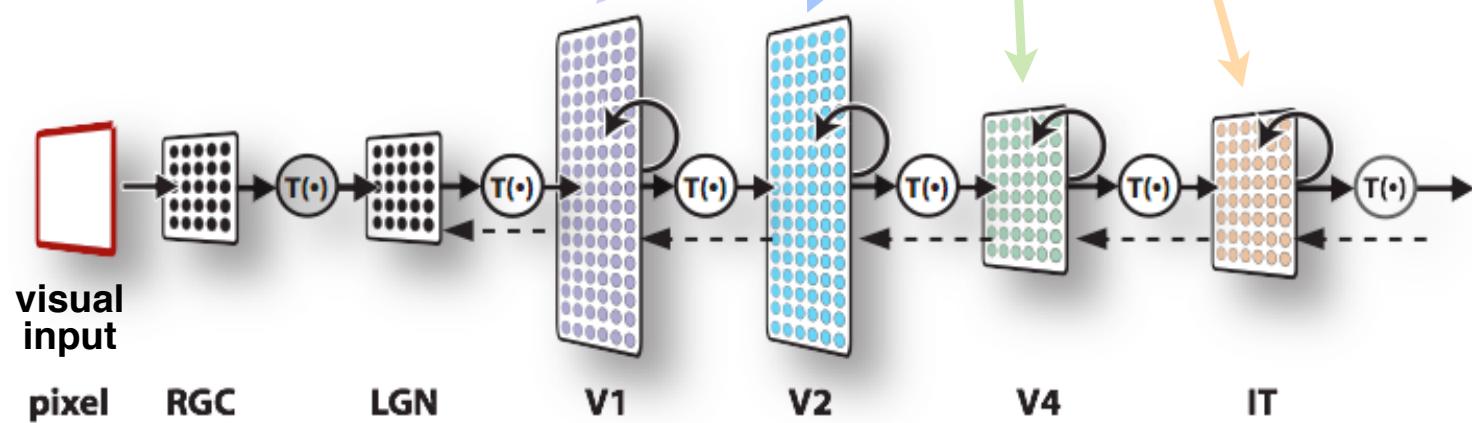
Intelligence test domain: Core Visual Object Perception

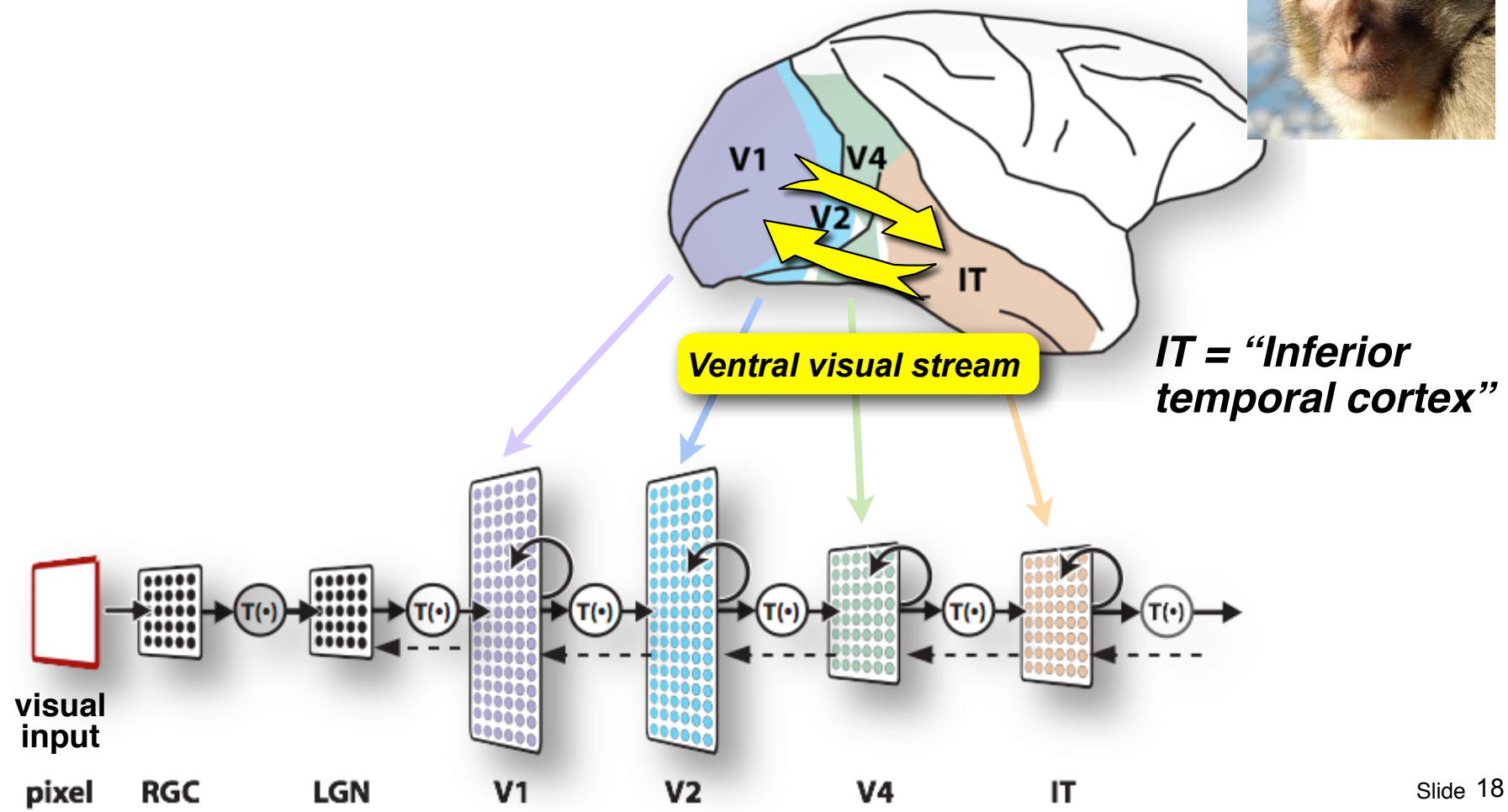


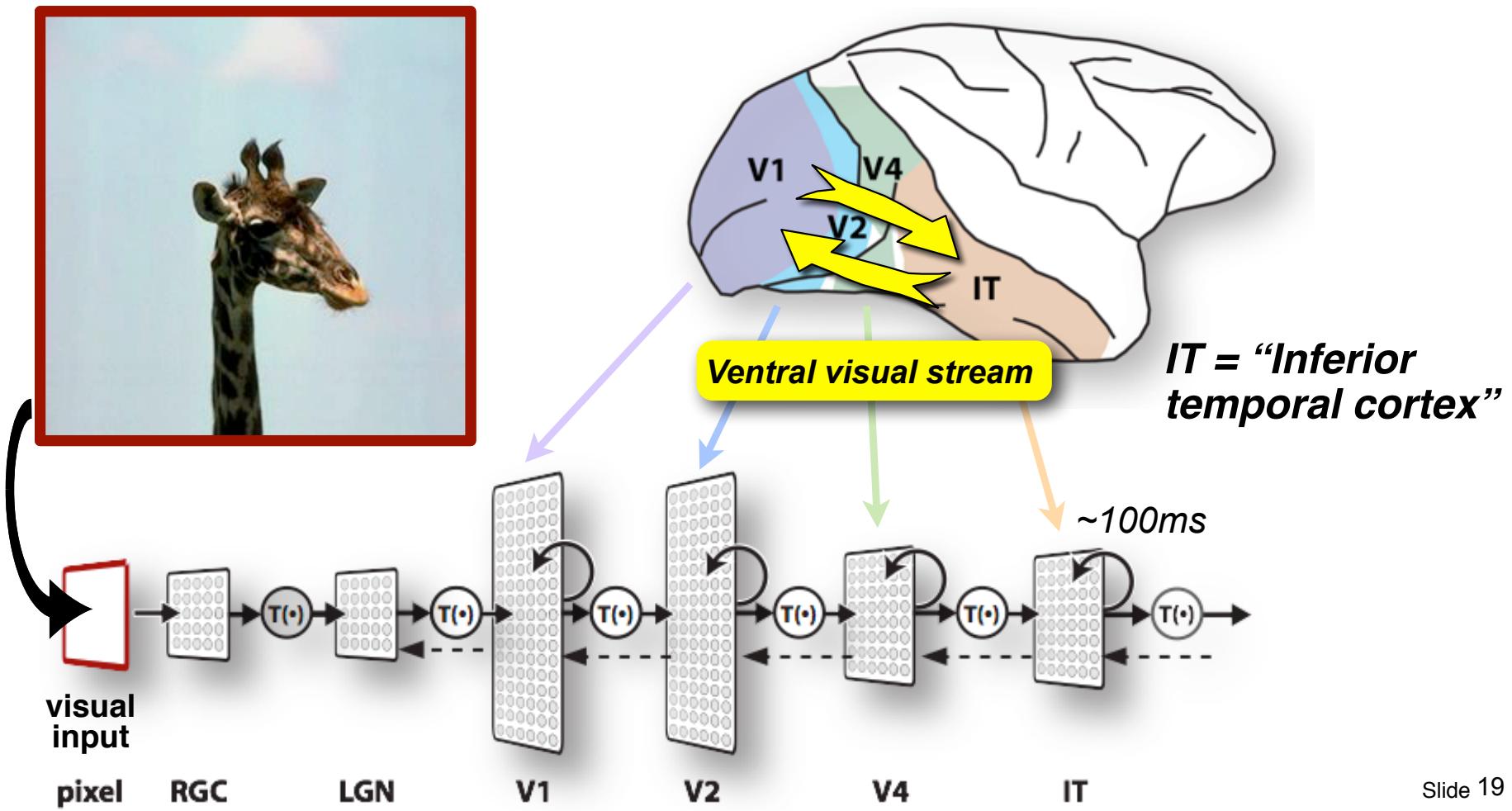
Lesions here result in deficits in object recognition.

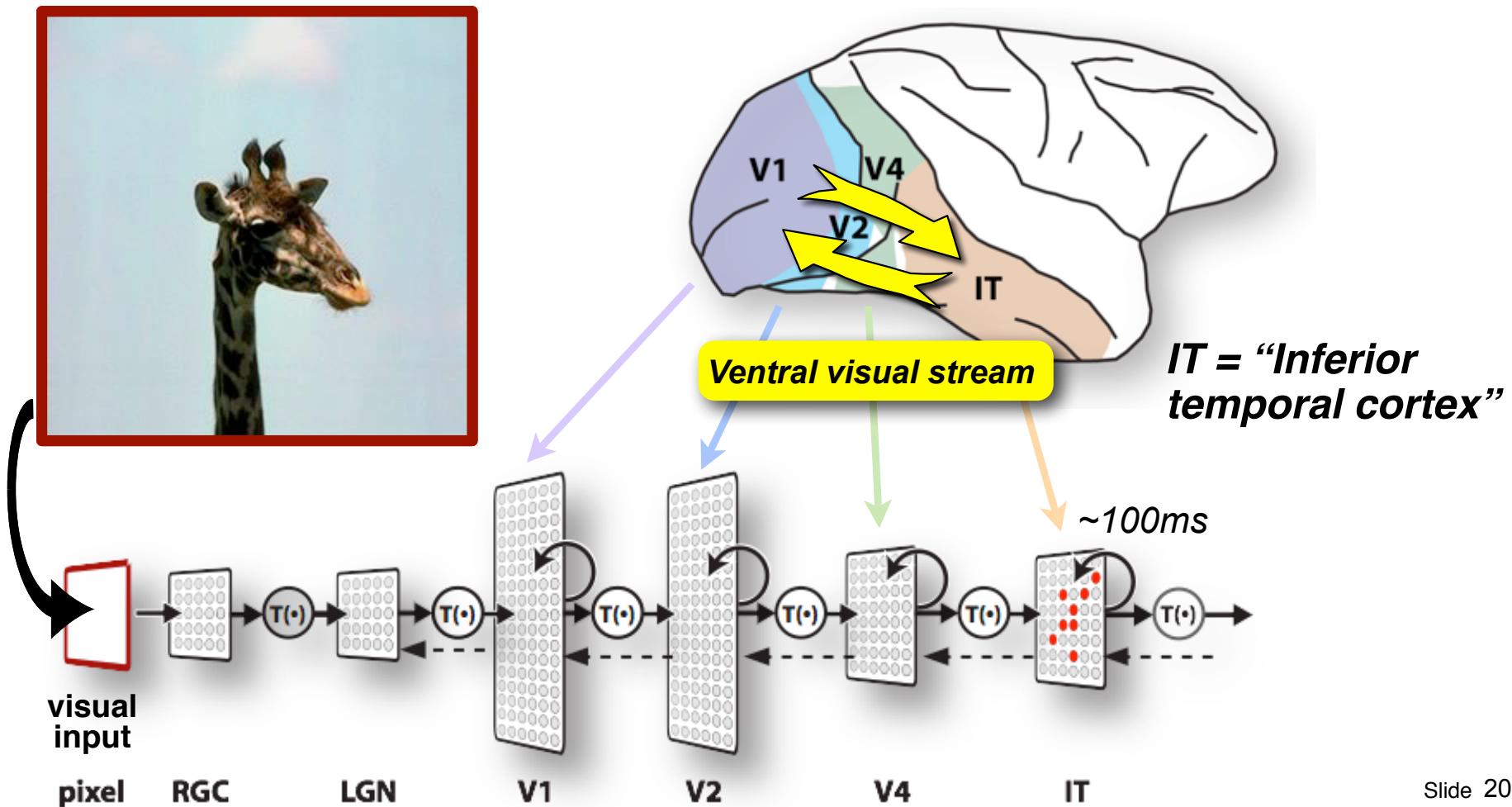


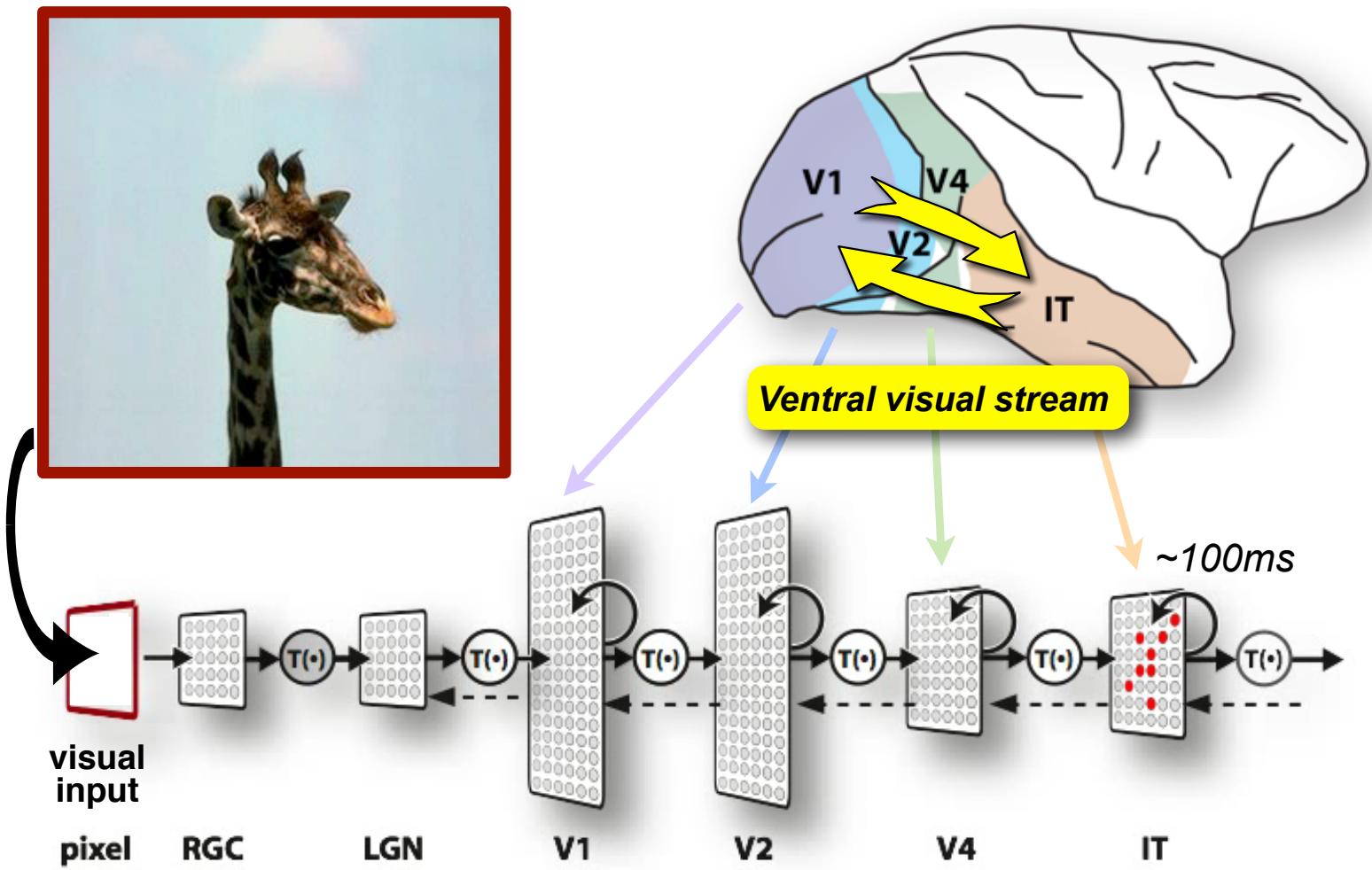
IT = “Inferior temporal cortex”



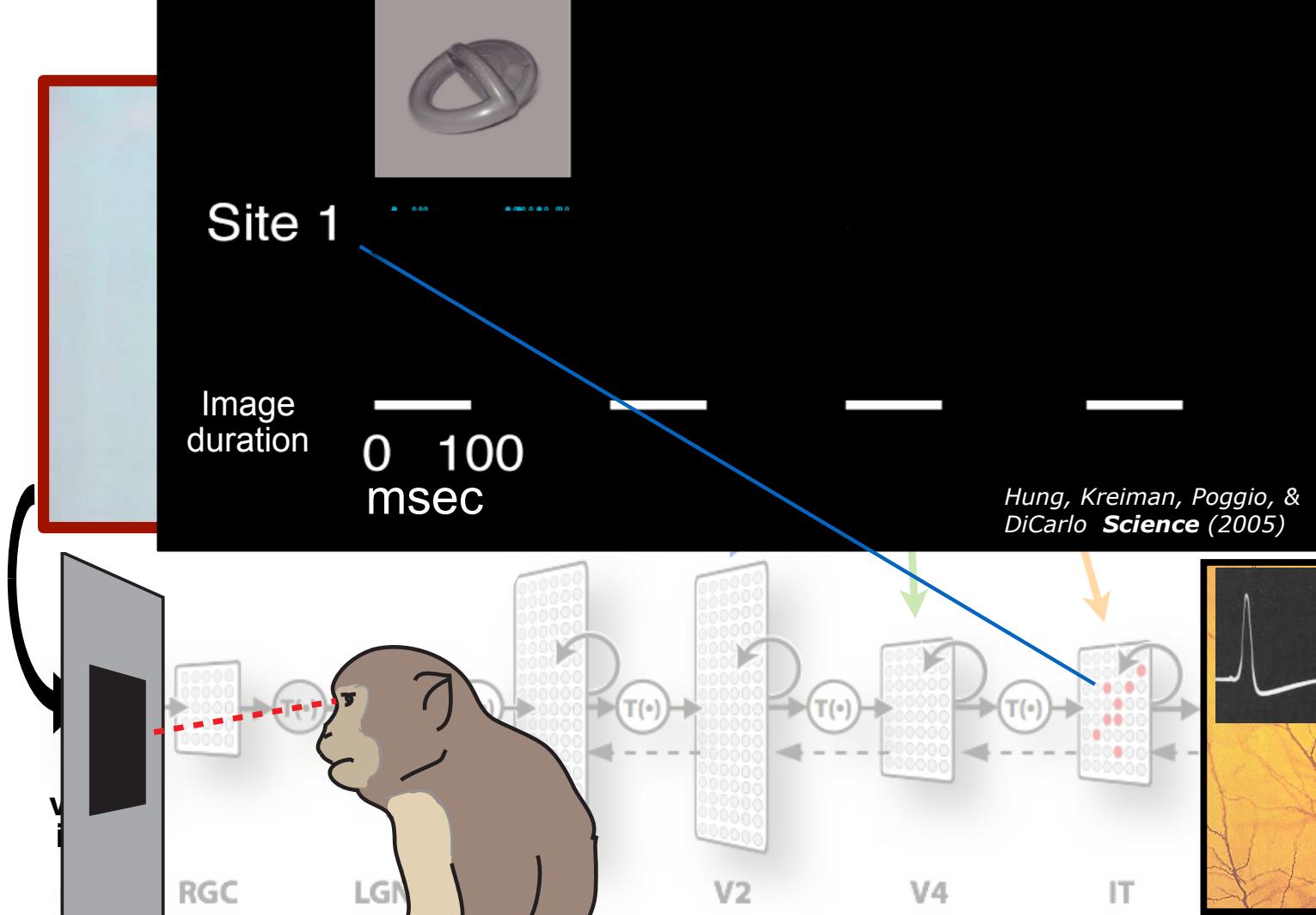




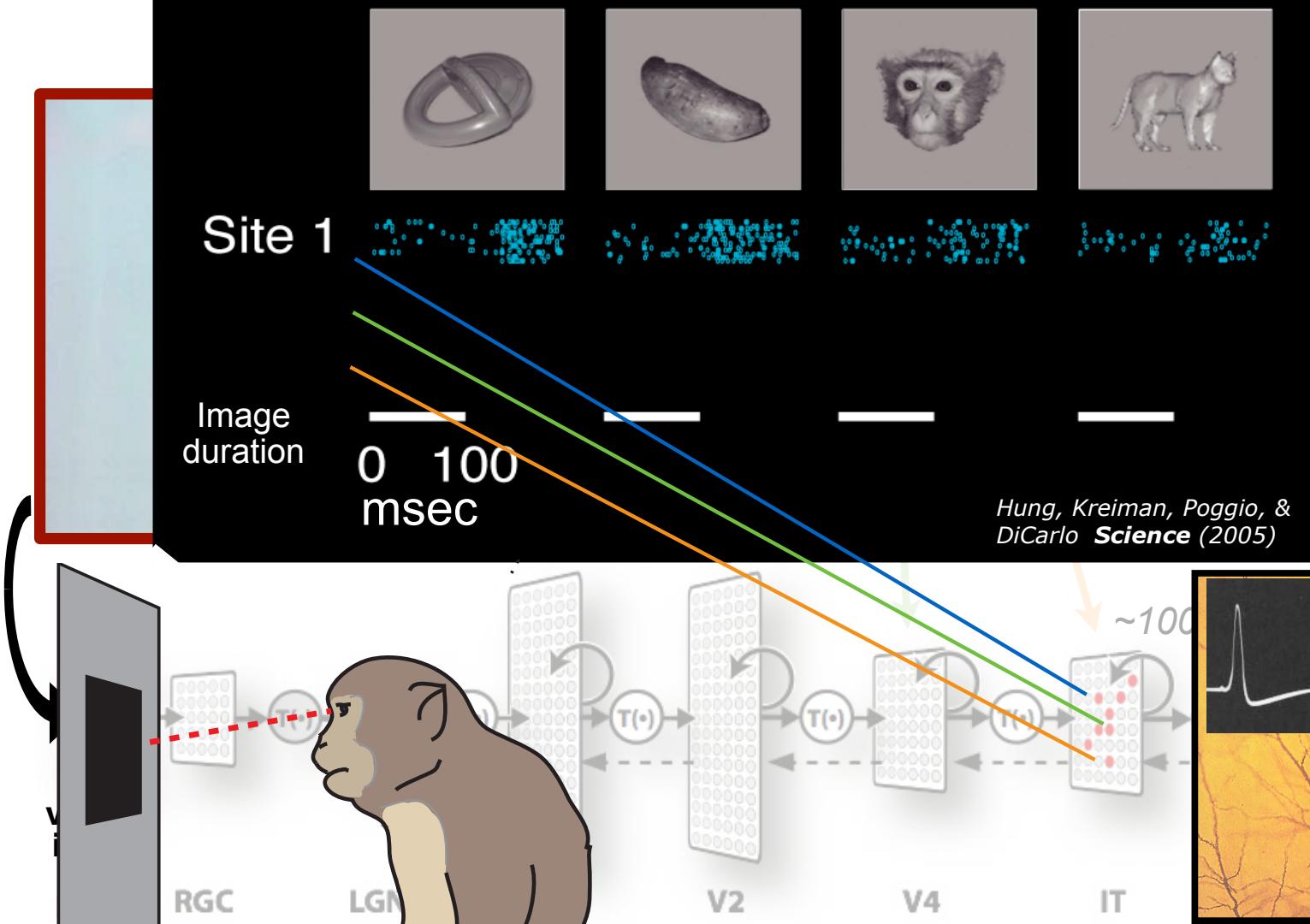




Examples of IT neuronal spiking responses



Examples of IT neuronal spiking responses



Examples of IT neuronal spiking responses

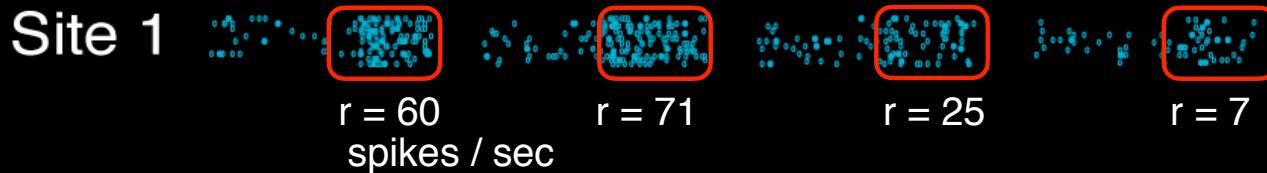
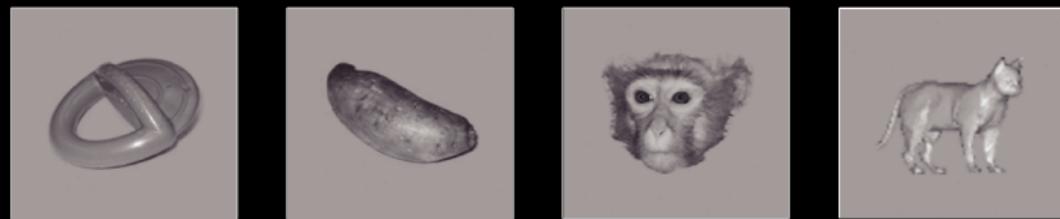
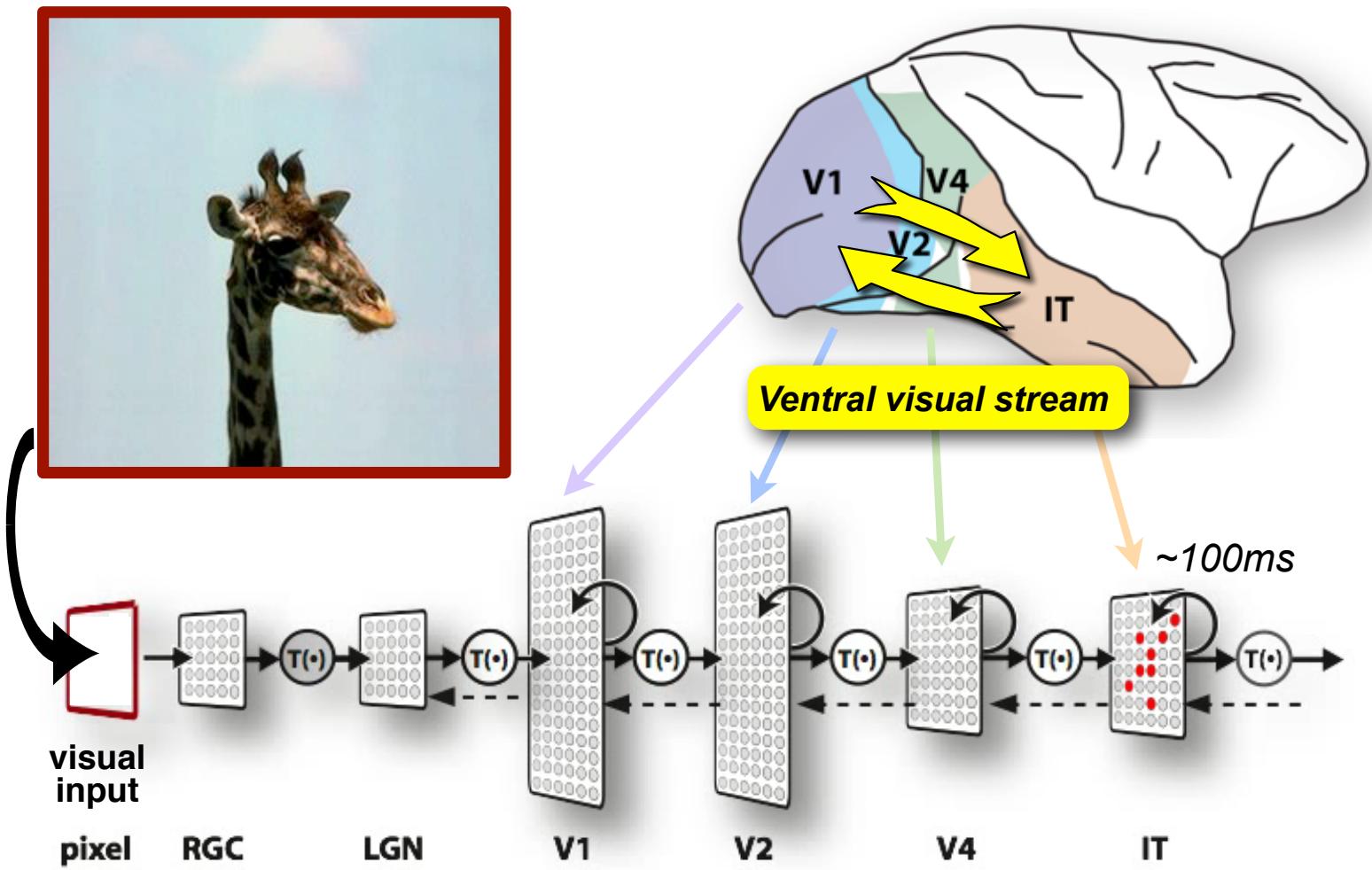
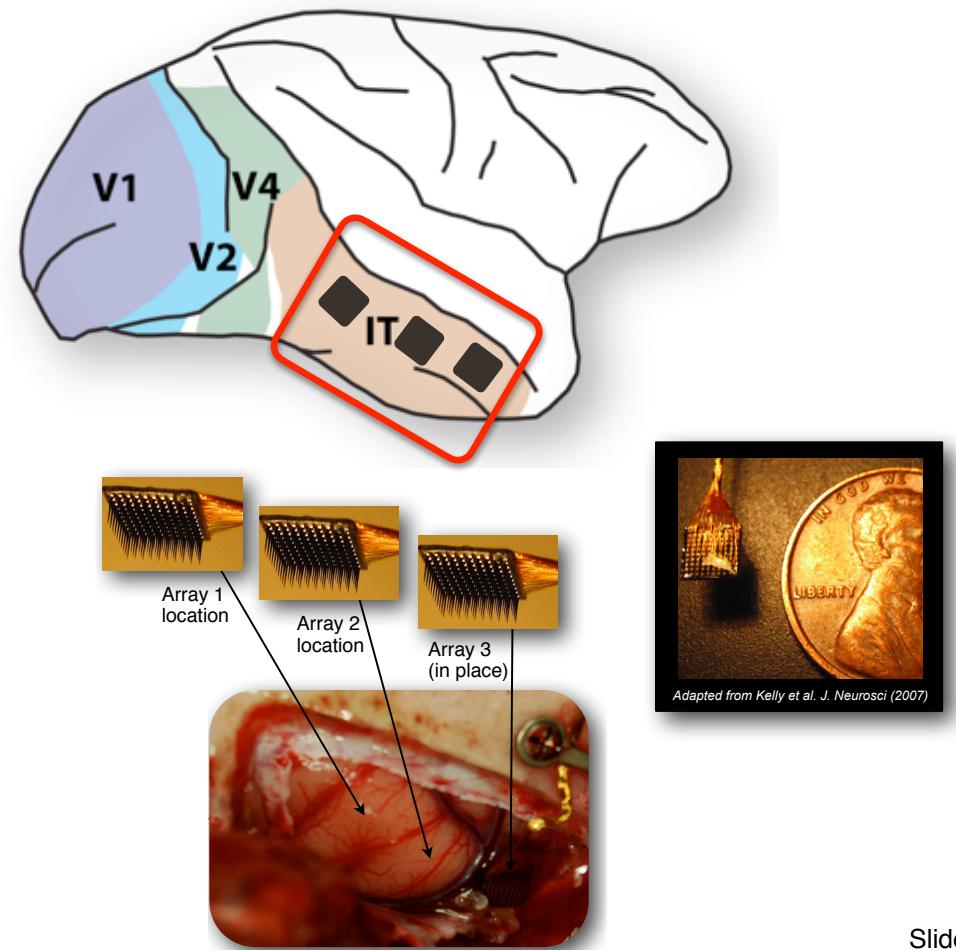
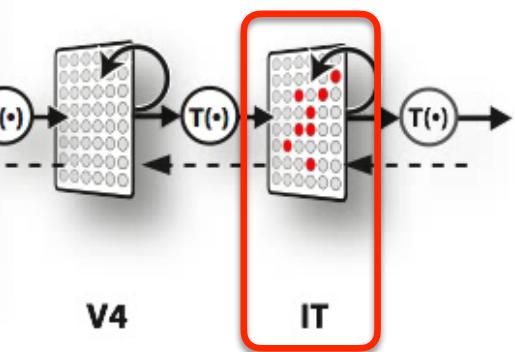


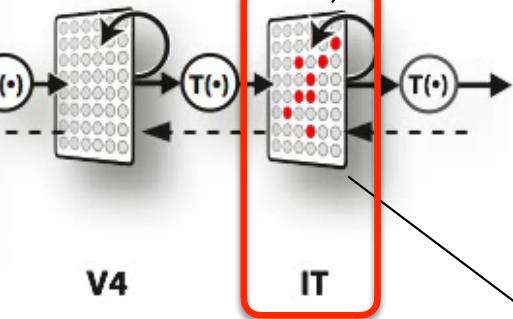
Image duration
0 100 msec

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





Neural response



100-1000

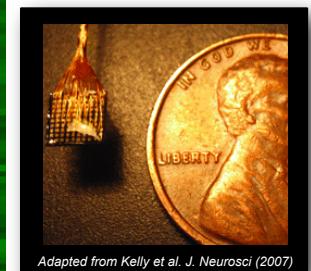
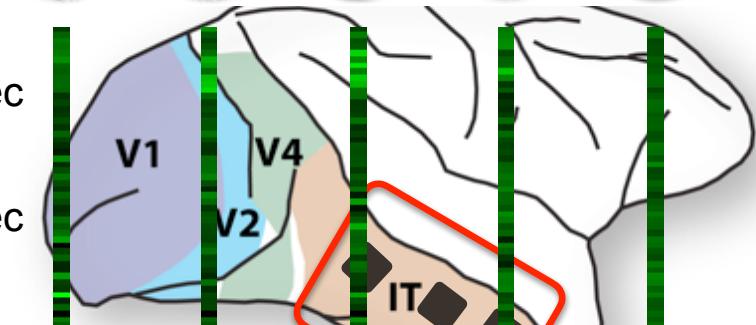
IT neuron sample number

1 *Image #*

8



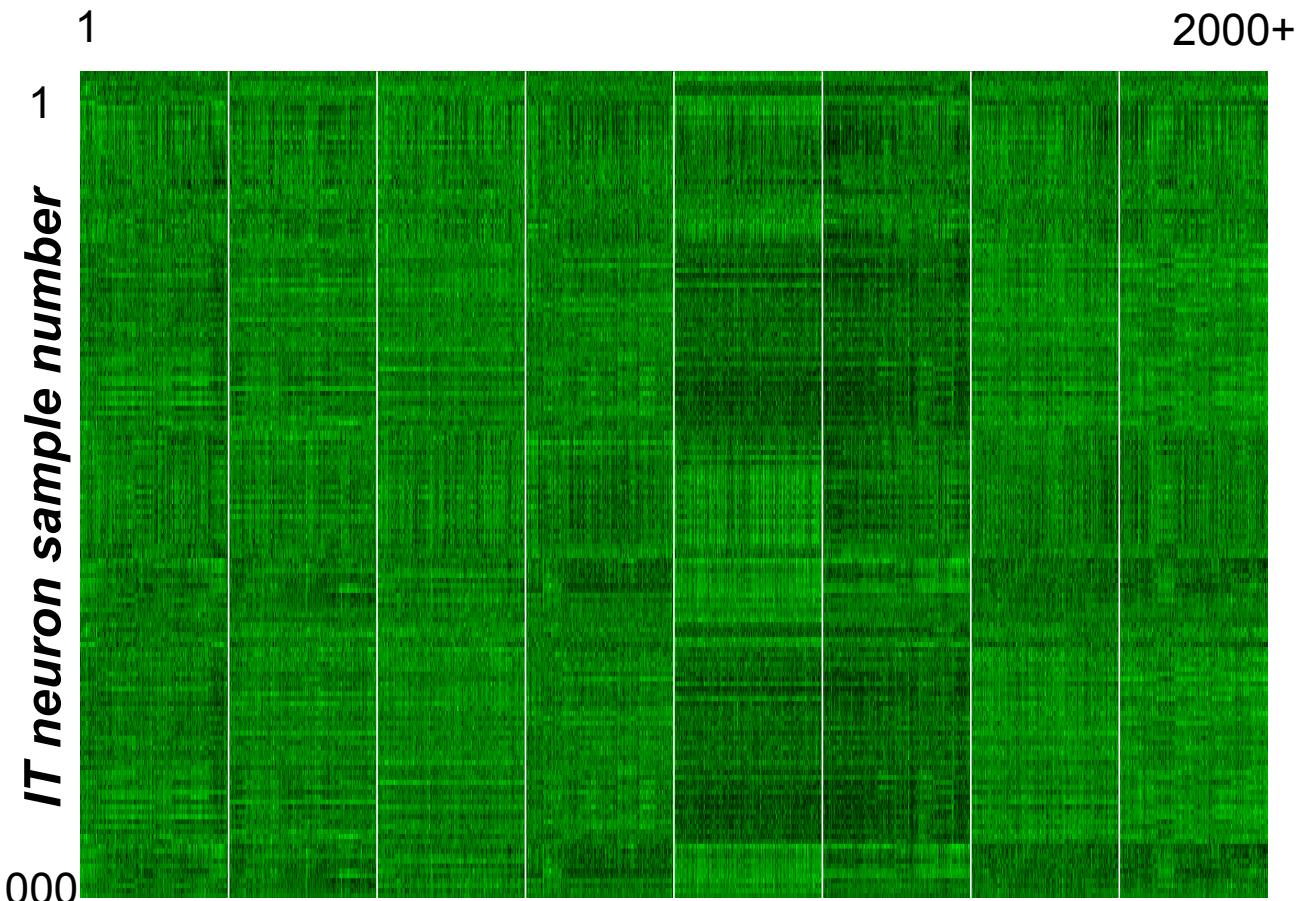
- 1 $r = 3$ spikes/sec
- 2 $r = 12$ spikes/sec
- 3 $r = 4$ spikes/sec
- 4 $r = 35$ spikes/sec
- ⋮



Neural response



Image #



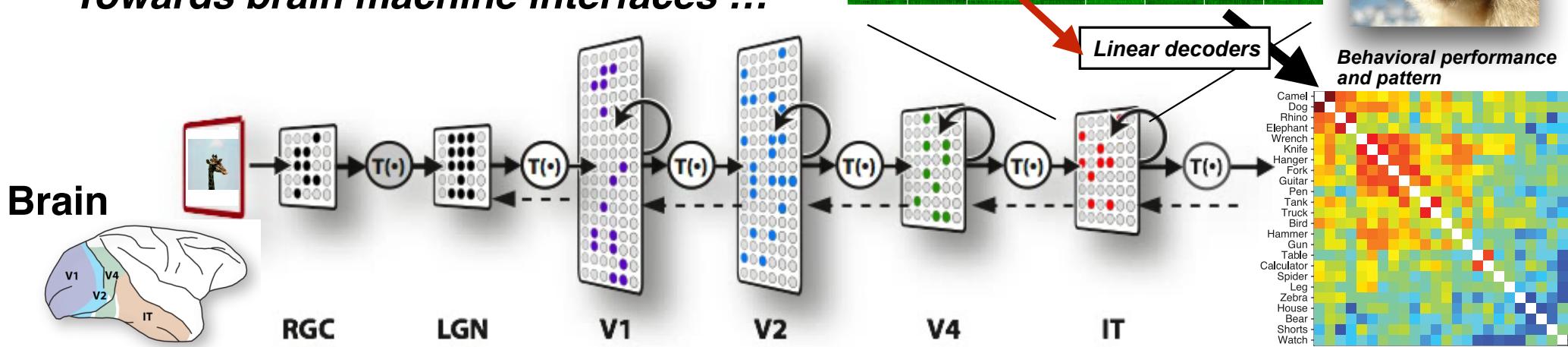
Hung*, Kreiman*, Poggio and DiCarlo, **Science** (2005);
Rust & DiCarlo, **J Neuroscience** (2010);
Majaj et al. **J Neuroscience** (2015)

The IT neural population representation explains & predicts object recognition behavior !

The parameters of this model of object perception tell us how we should manipulate IT neural responses to predictably modify object percepts.

(Afraz et al. **PNAS** 2015 ; Rajalingham, **Neuron** 2018)

Towards brain machine interfaces ...

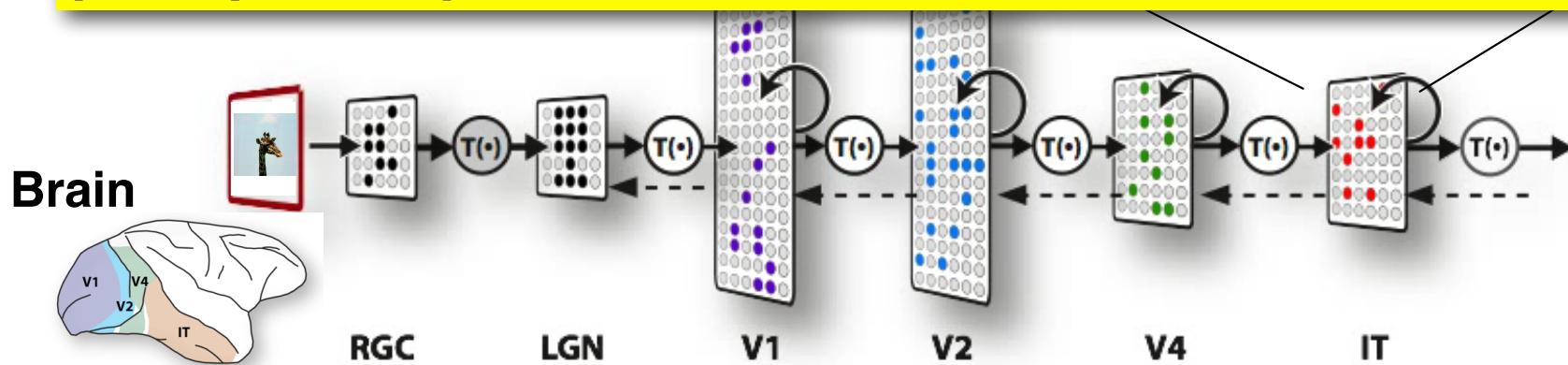


The IT neural population representation explains & predicts object recognition behavior

AI relevance: Primates are behaviorally higher performing than computer vision systems because their brain can compute this IT neural representation !

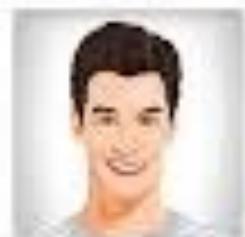
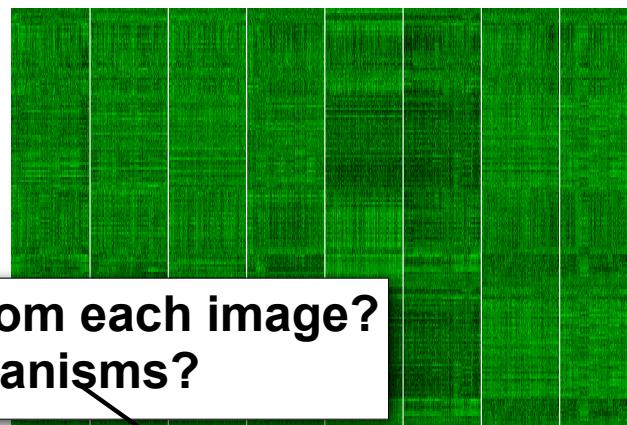


One key take-away: explaining the mean IT firing rates is ~sufficient to (computationally) explain behavior & perceptual report



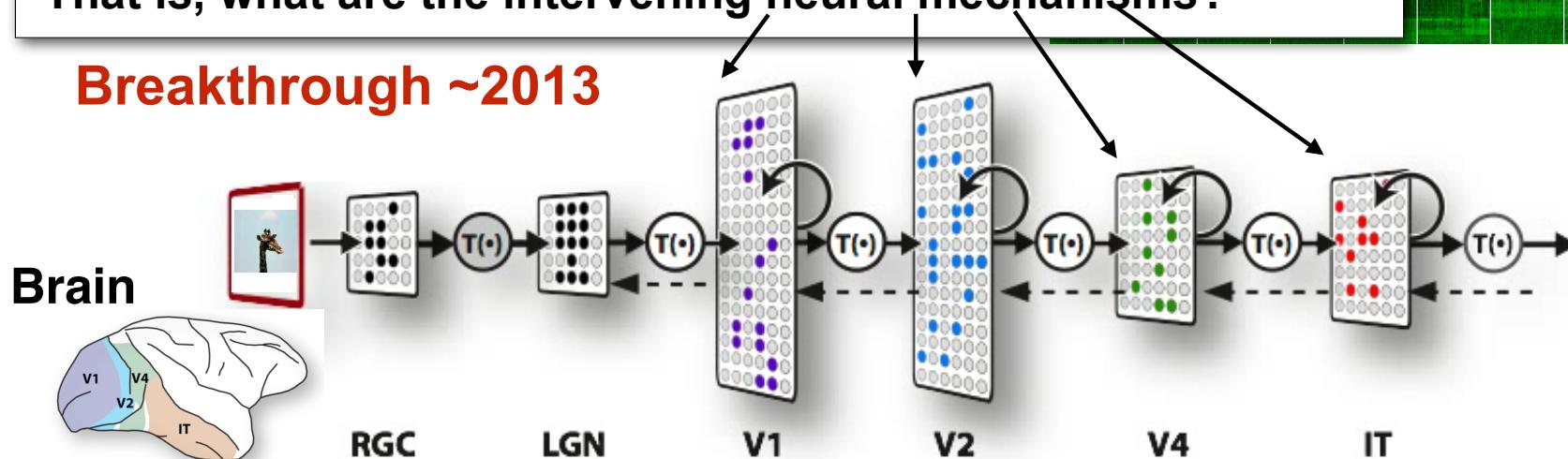
The IT neural population representation explains & predicts object recognition behavior

AI relevance: Primates are behaviorally higher performing than computer vision systems because their brain can compute this IT neural representation !

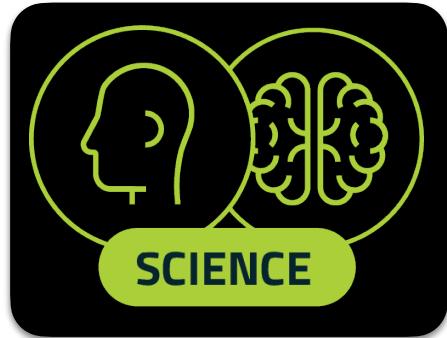


But how is the IT representation computed from each image?
That is, what are the intervening neural mechanisms?

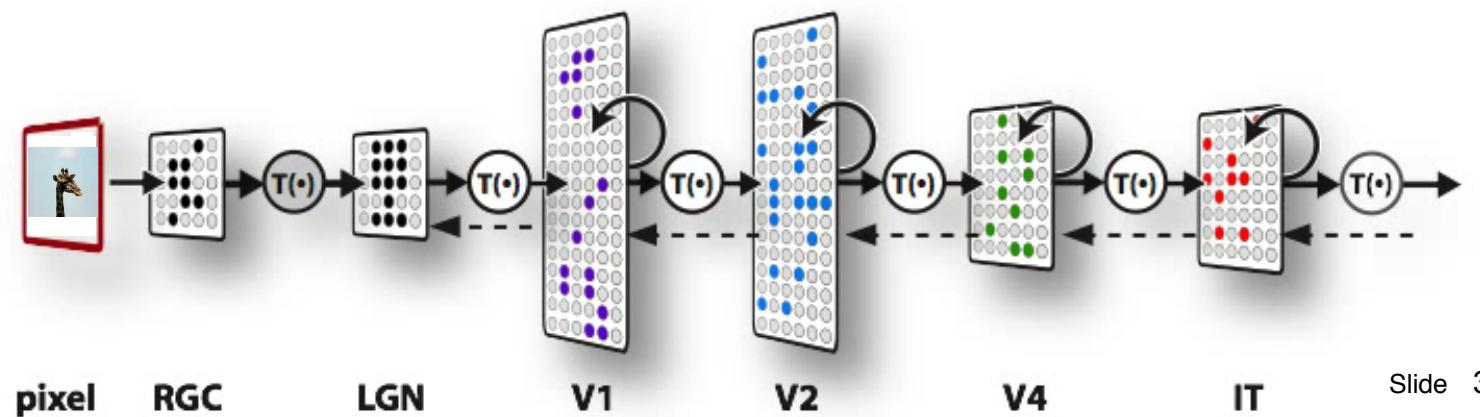
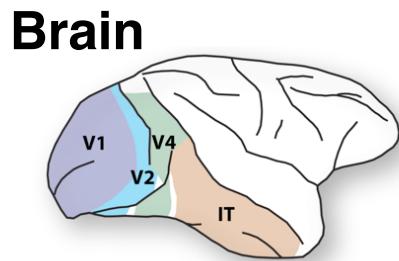
Breakthrough ~2013



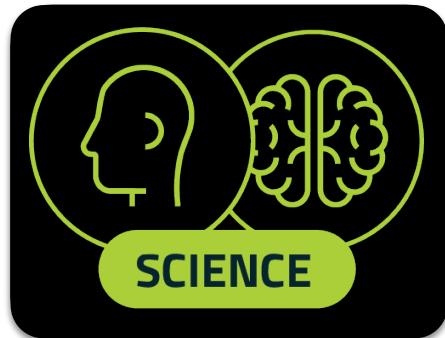
Background



GUIDANCE FROM
NEUROSCIENCE (many labs):



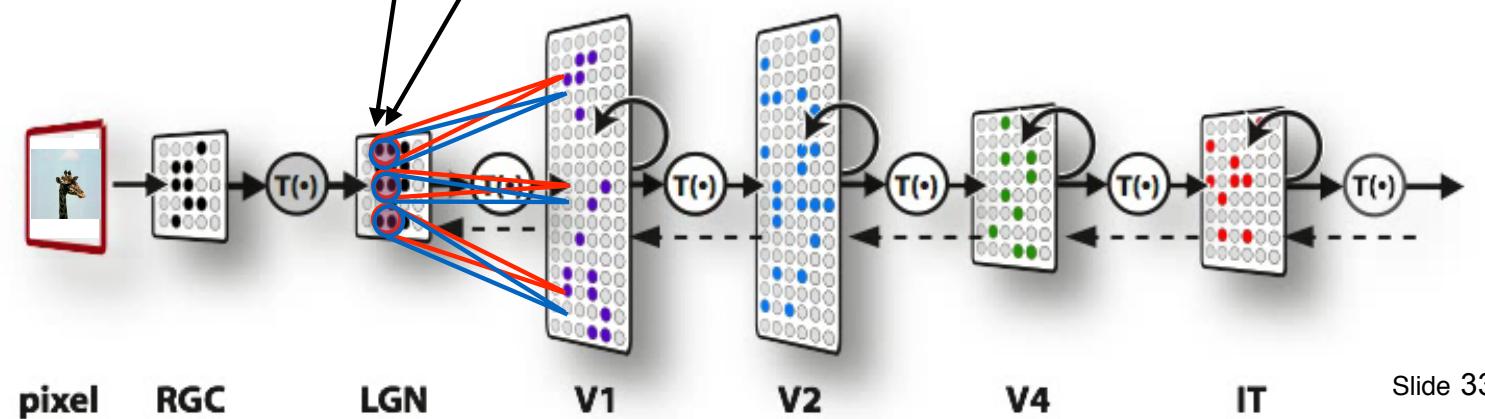
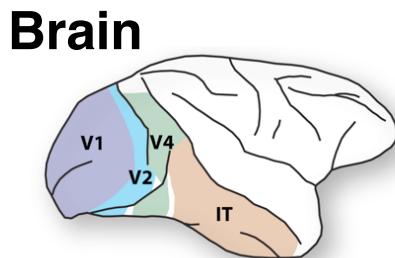
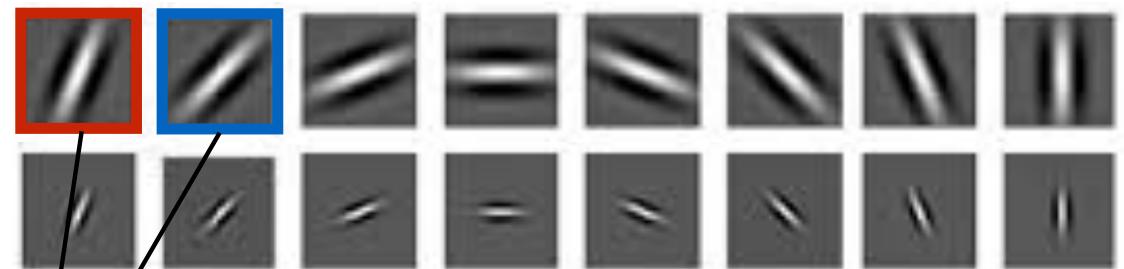
Background



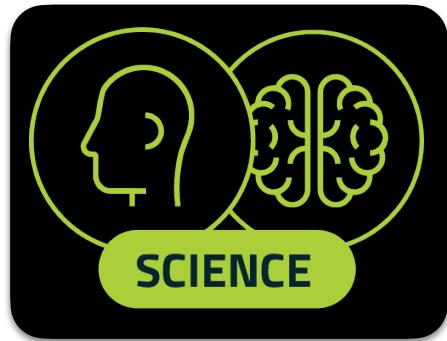
GUIDANCE FROM NEUROSCIENCE (many labs):

Each “area” of processing:

- *Spatially local, ~linear filters*
- *Different types of such filters*
- *Each repeated spatially over the input (~convolution)*



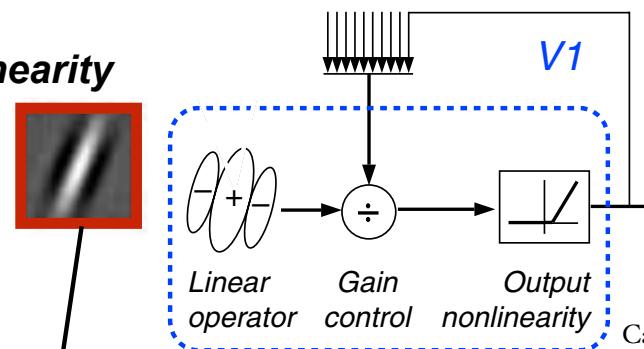
Background



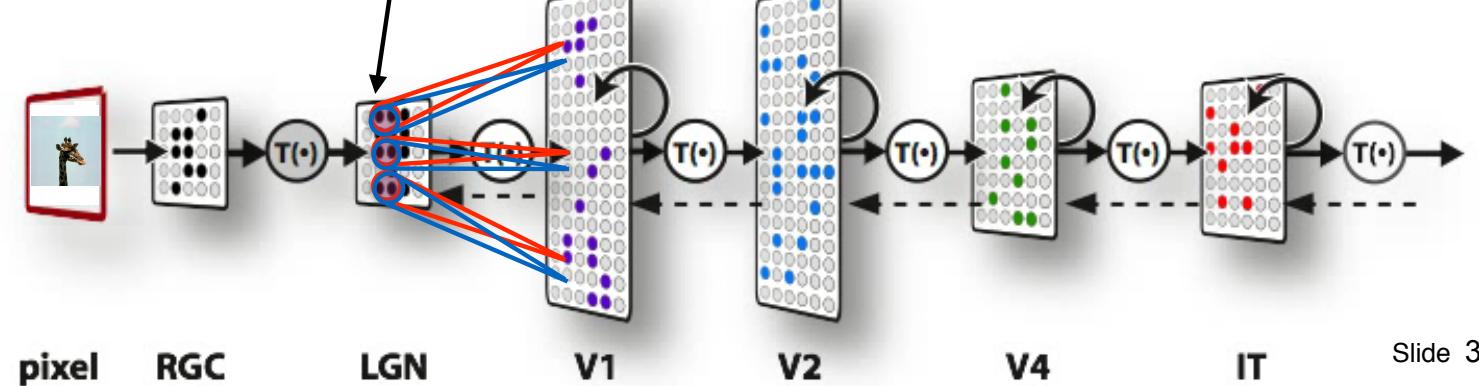
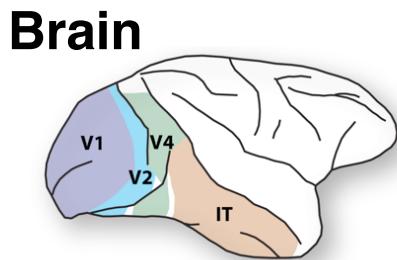
GUIDANCE FROM NEUROSCIENCE (many labs):

Each “area” of processing:

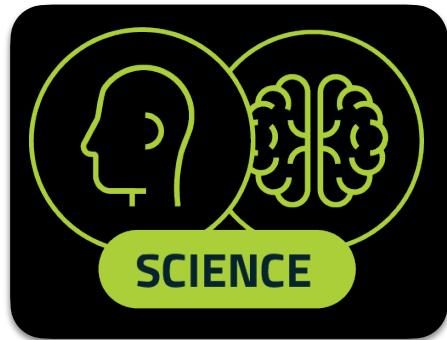
- *Spatially local, ~linear filters*
- *Different types of such filters*
- *Each repeated spatially over the input (~convolution)*
- *Rectifying non-linearity*
- *Normalization*



Carandini & Heeger, 1994



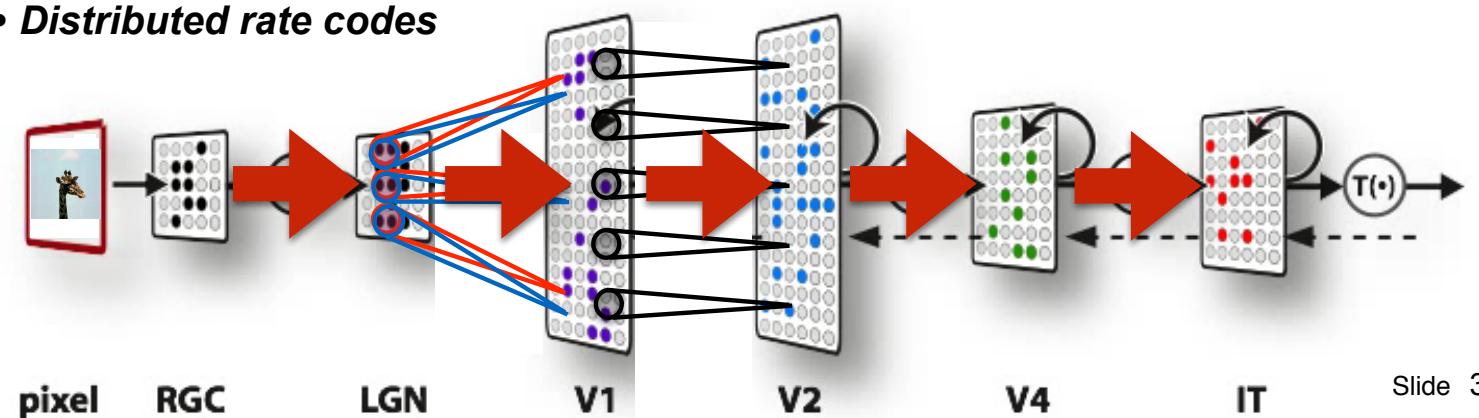
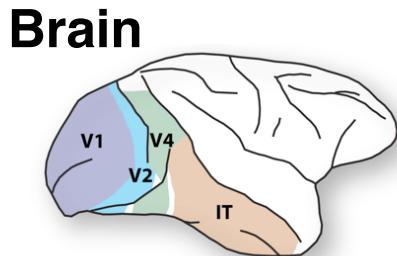
Background



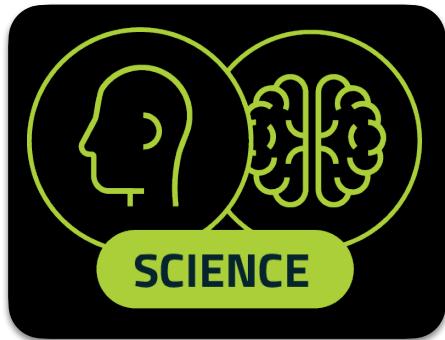
GUIDANCE FROM NEUROSCIENCE (many labs):

Each “area” of processing:

- *Spatially local, ~linear filters*
- *Different types of such filters*
- *Each repeated spatially over the input (~convolution)*
- *Rectifying non-linearity*
- *Normalization*
- **“Deep” series of areas**
- **Similar “style” operations at each successive area**
- **Fast ~feedforward does a lot!**
- **Distributed rate codes**



Background



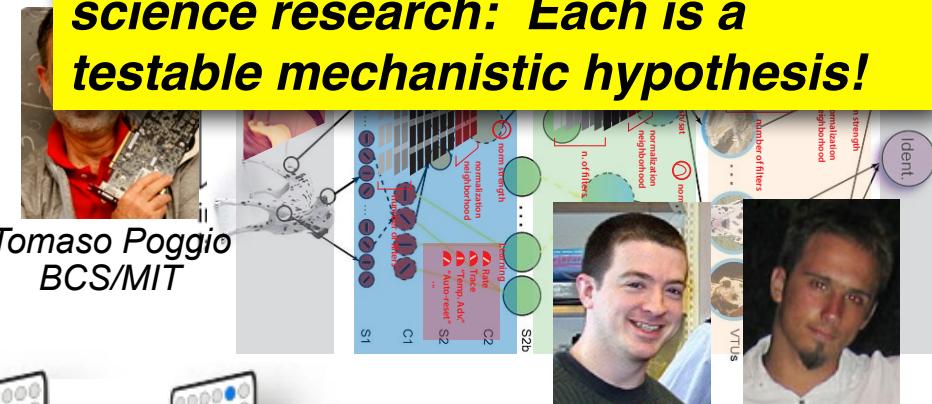
GUIDANCE FROM NEUROSCIENCE (many labs):

- **Each “area” of processing:**
- *Spatially local, ~linear filters*
- *Different types of such filters*
- *Each repeated spatially over the image (~convolution)*
- *Rectifying non-linearity*
- *Normalization*
- **“Deep” series of areas**
- *Similar “style” operations at each successive area*
- *Fast ~feedforward does a lot!*
- *Distributed rate codes*

Resulted in proposed feedforward artificial neural networks (ANNs):

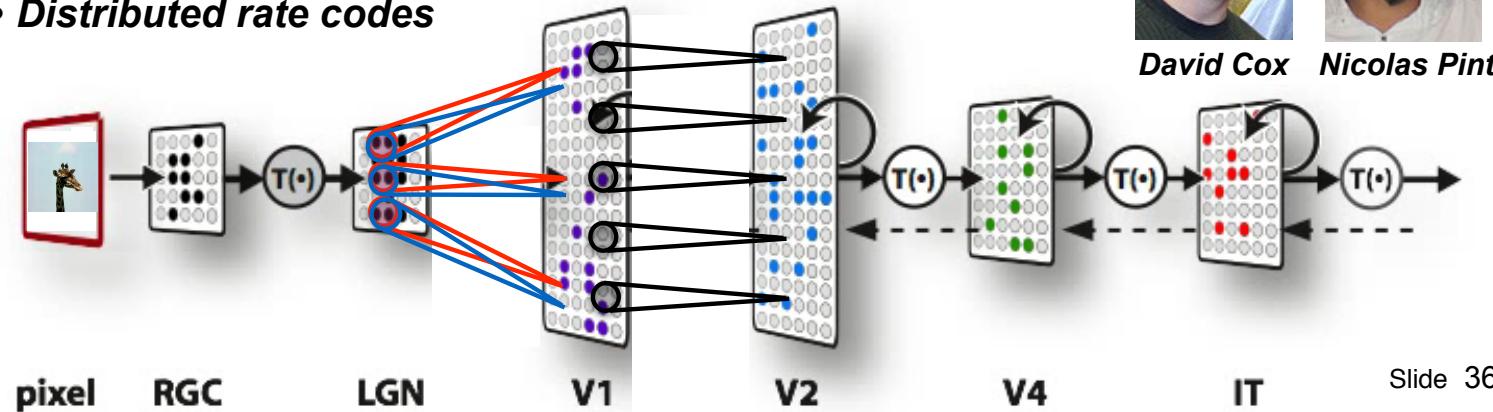
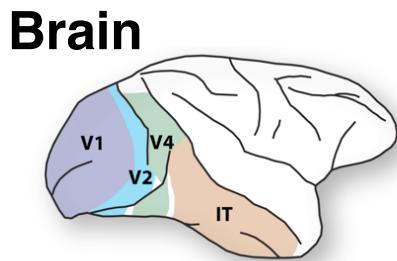
Pinto and Cox 2008-2010

The building of such models is critically important to basic science research: Each is a testable mechanistic hypothesis!



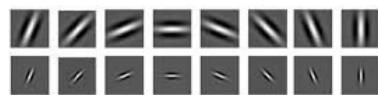
Tomaso Poggio
BCS/MIT

David Cox *Nicolas Pinto*



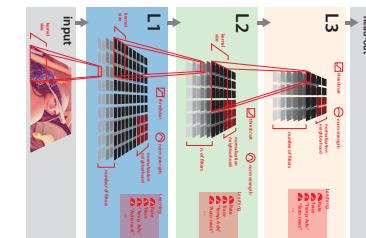
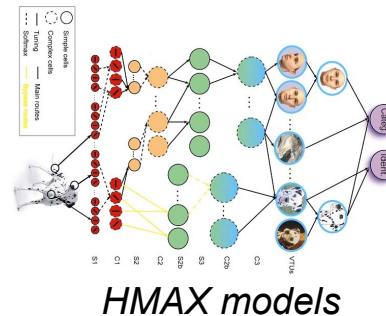
Hypotheses (specific ANN models):

~1980-2010



“V1-like” models

“V2-like” models



PLoS09 models

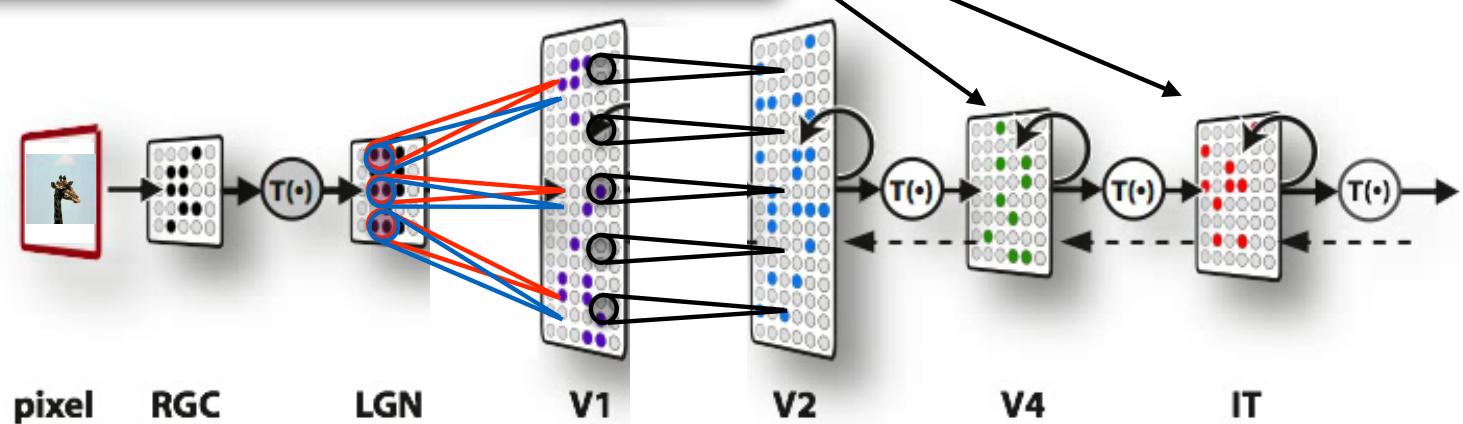
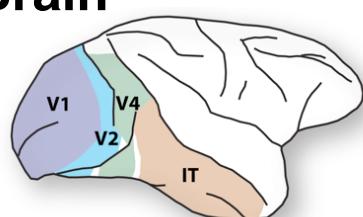
Unfortunately, all of these specific mechanistic hypotheses were inadequate.

E.g. they each failed to accurately predict the internal neural responses to new sets of test images.

GUIDANCE FROM NEUROSCIENCE:

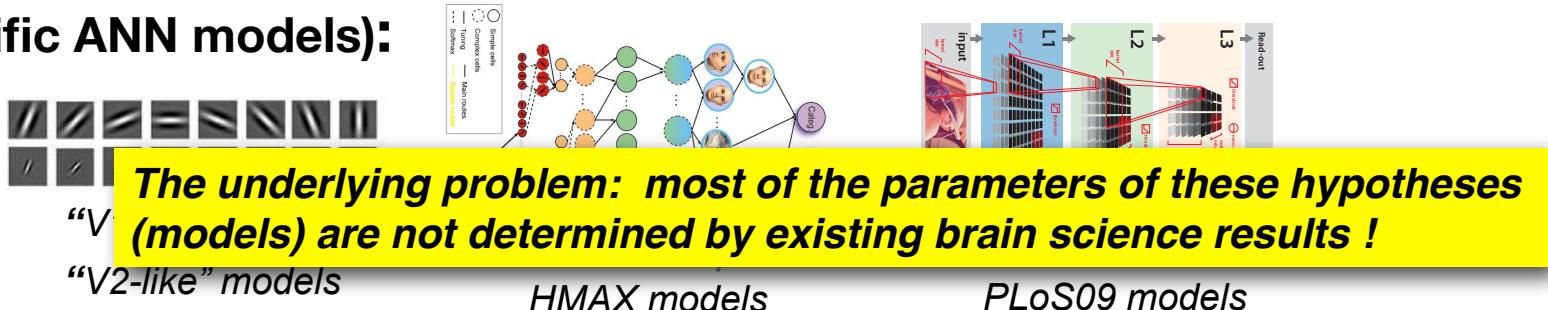


Brain



Hypotheses (specific ANN models):

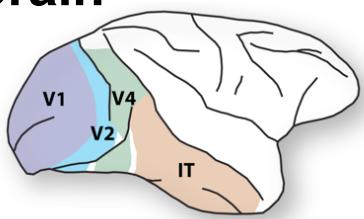
~1980-2010



GUIDANCE FROM NEUROSCIENCE:

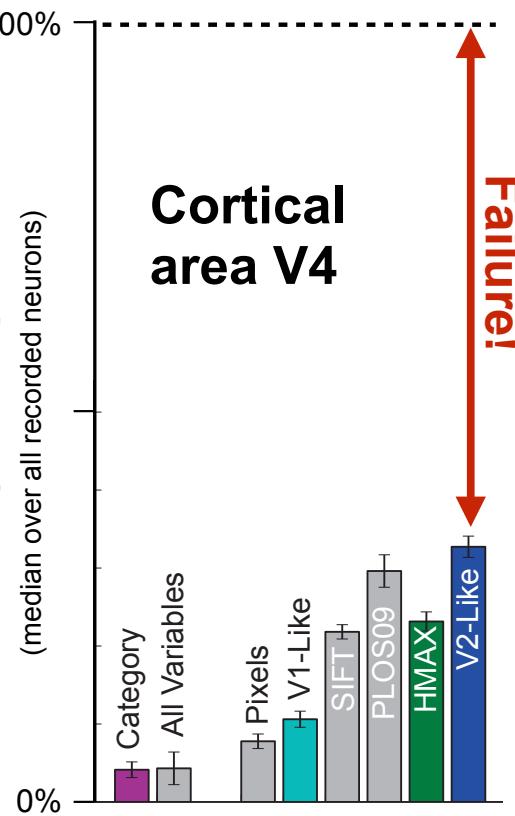


Brain

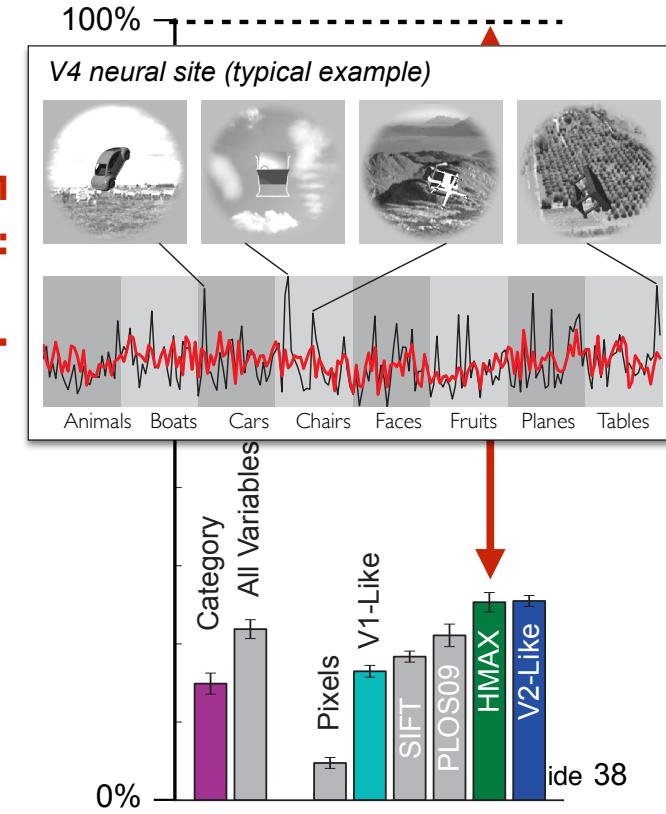


"explained" means predicted for new images

Fraction of visually driven neural response explained



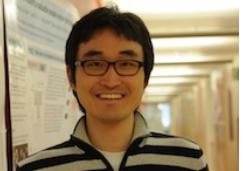
Cortical area V4



~2013: Collaborative breakthrough



Dan Yamins



Ha Hong

Yamins, Hong,
Solomon, Seibert and
DiCarlo **NIPS** (2013),
PNAS (2014)

The underlying problem: most of the parameters of these hypotheses (models) are not determined by existing brain science results !

~2013: Collaborative breakthrough



Dan Yamins

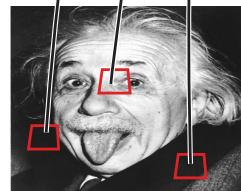
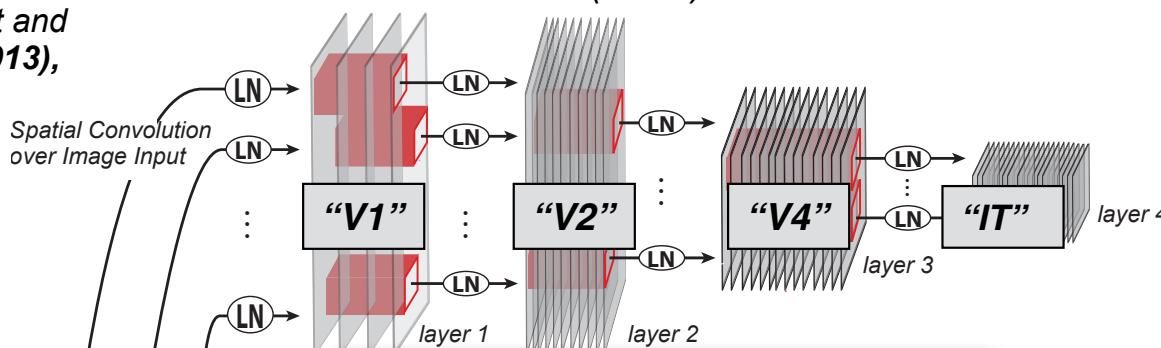
Ha Hong

Cognitive science guided the task

The underlying problem: most of the parameters of these hypotheses (models) are not determined by existing brain science results !

Yamins, Hong,
Solomon, Seibert and
DiCarlo **NIPS (2013)**,
PNAS (2014)

Artificial neural network (ANN)



Example test image
(one of many)

Neuroscience guided the
parameters of system macro-
and meso- architecture

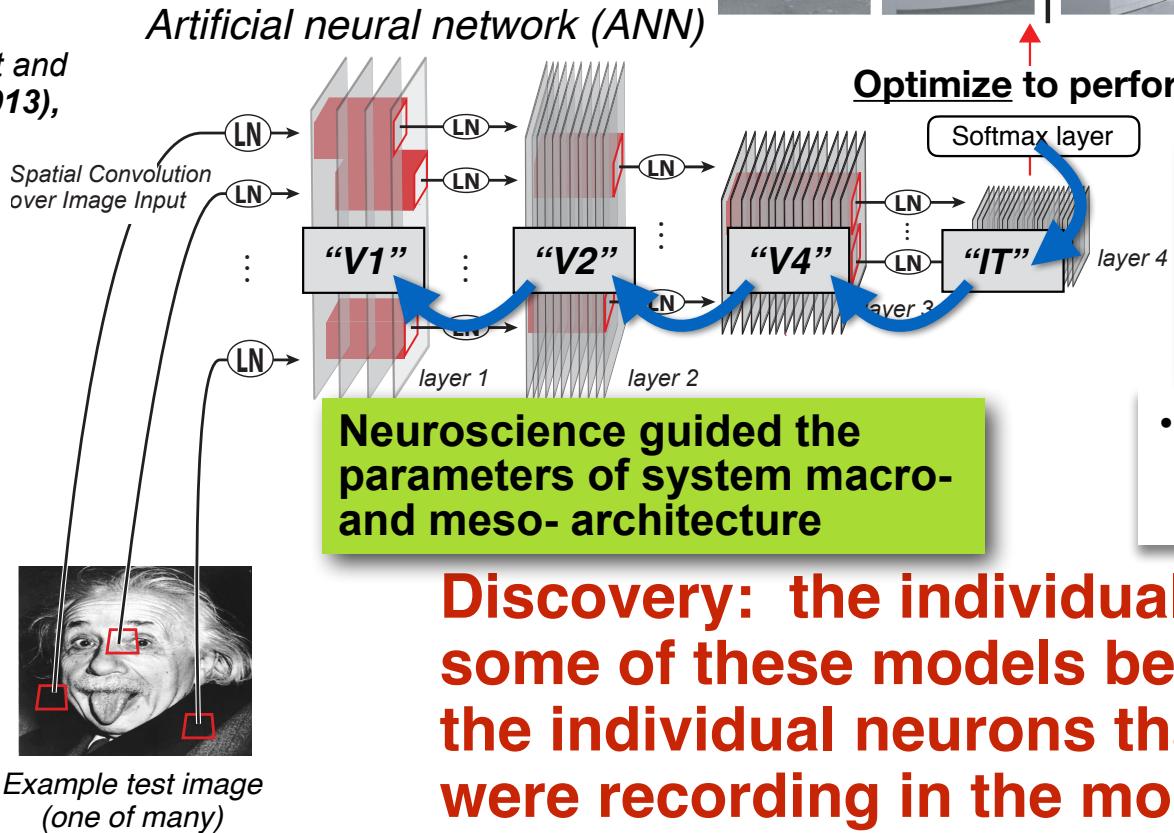
~2013: Collaborative breakthrough



Dan Yamins

Ha Hong

Yamins, Hong,
Solomon, Seibert and
DiCarlo **NIPS** (2013),
PNAS (2014)



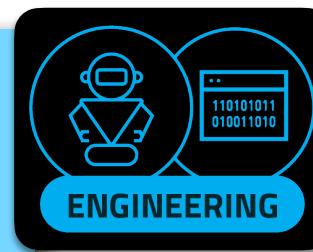
Cognitive science guided the task

Core object recognition



Optimize to perform this task!

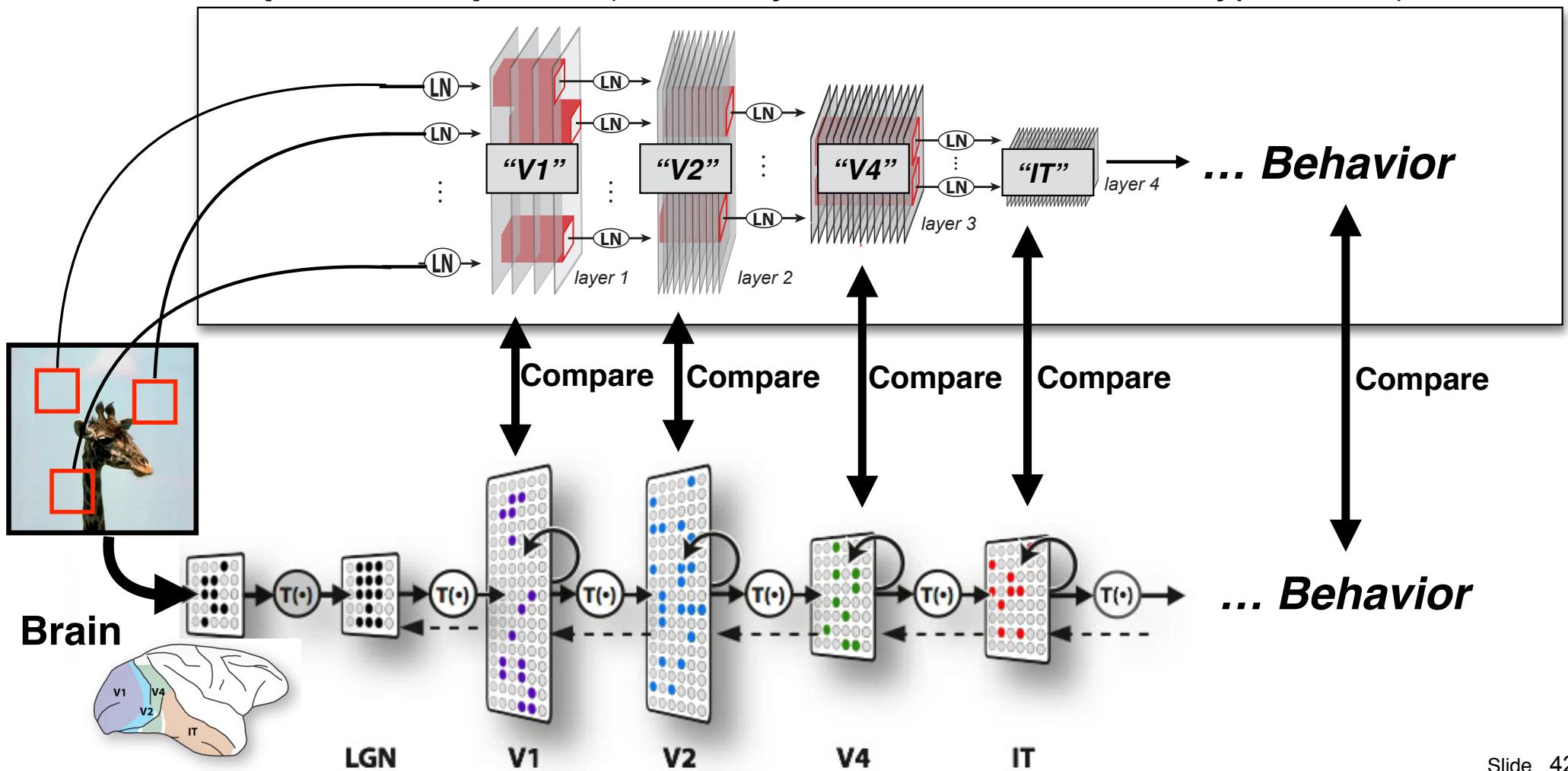
**Engineering tools
to tune the
microarchitecture
parameters to
perform well on
this task !**

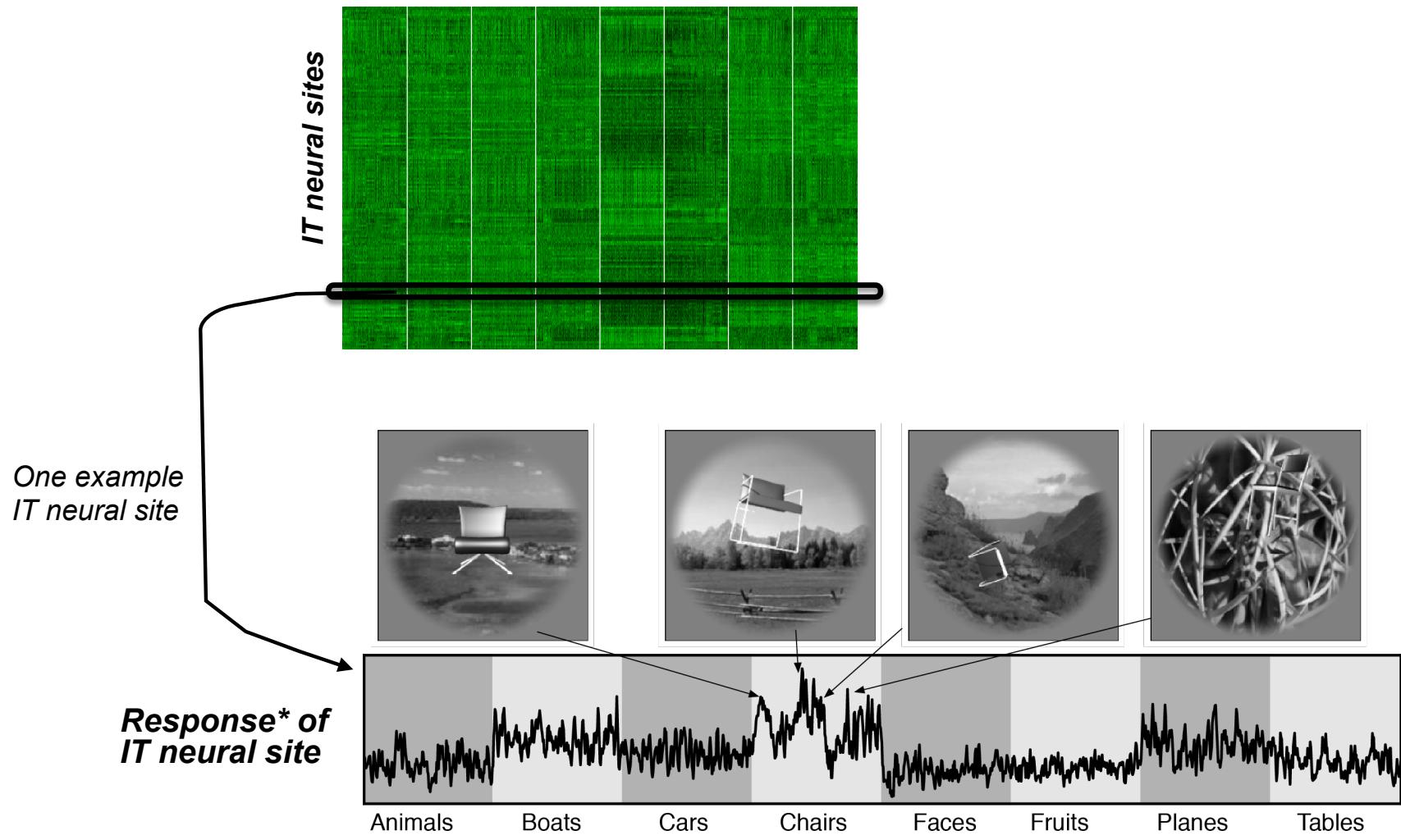


- Each choice of all system parameters is an entire new artificial ventral stream!

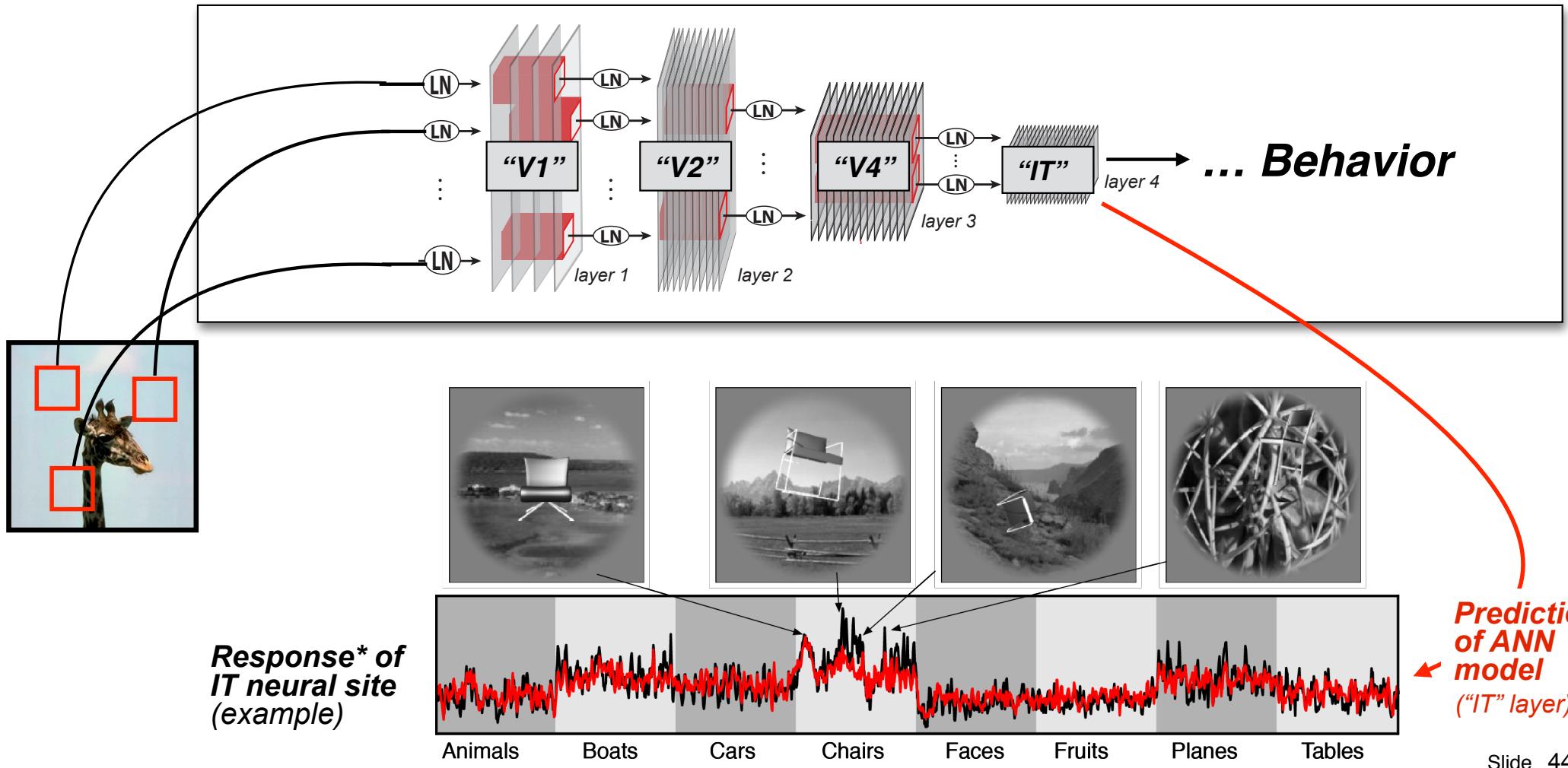
Discovery: the individual “neurons” inside some of these models behave very much like the individual neurons that we and others were recording in the monkey brain !

A specific deep ANN (a neurally-mechanistic scientific hypothesis!)



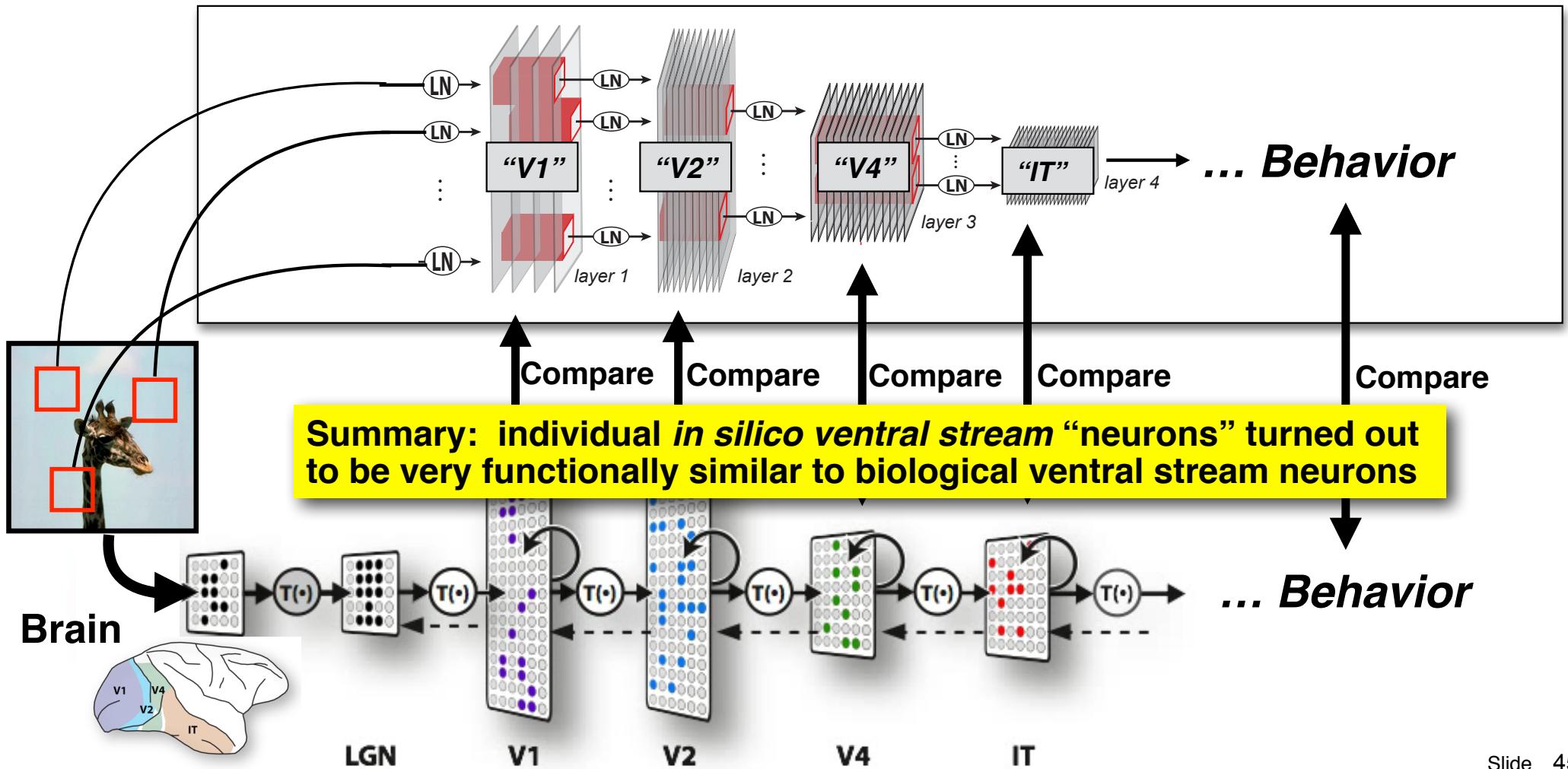


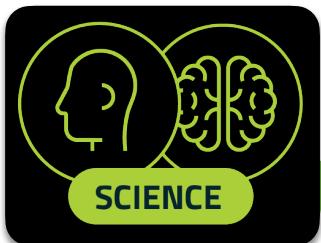
A specific deep ANN (a neurally-mechanistic scientific hypothesis!)



Yamins, Hong, ... DiCarlo **NeurIPS (2013), PNAS (2014)**

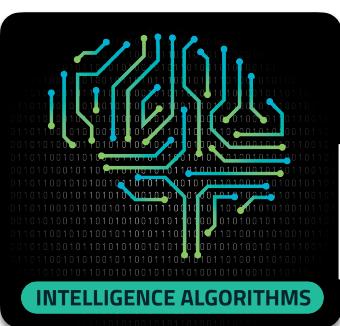
A specific deep ANN (a neurally-mechanistic scientific hypothesis!)



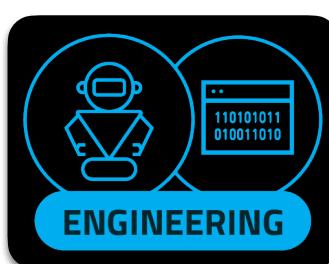
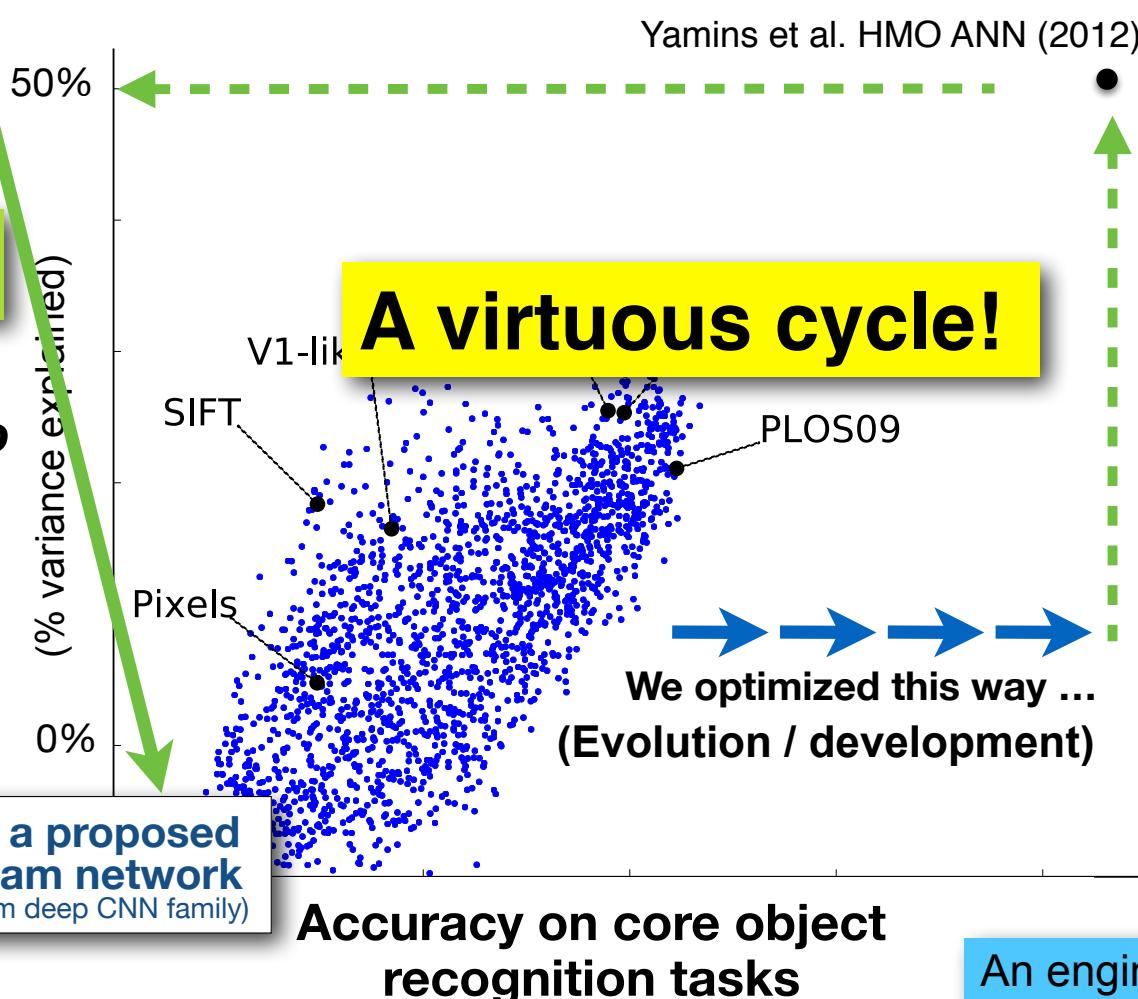


A neuroscience goal

**Match score
of the *in silico*
“IT” neurons
with primate
IT neurons**

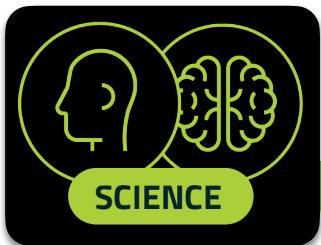


Each dot is a proposed
ventral stream network
(here sampled from deep CNN family)



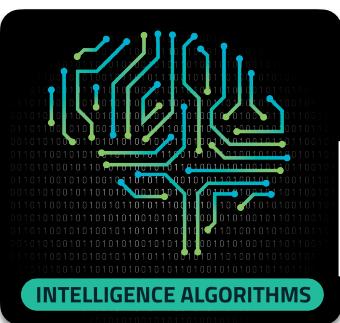
An engineering / AI goal

Models
Adapted from Yamins, Hong, Solomon,
Seibert and DiCarlo PNAS (2014)

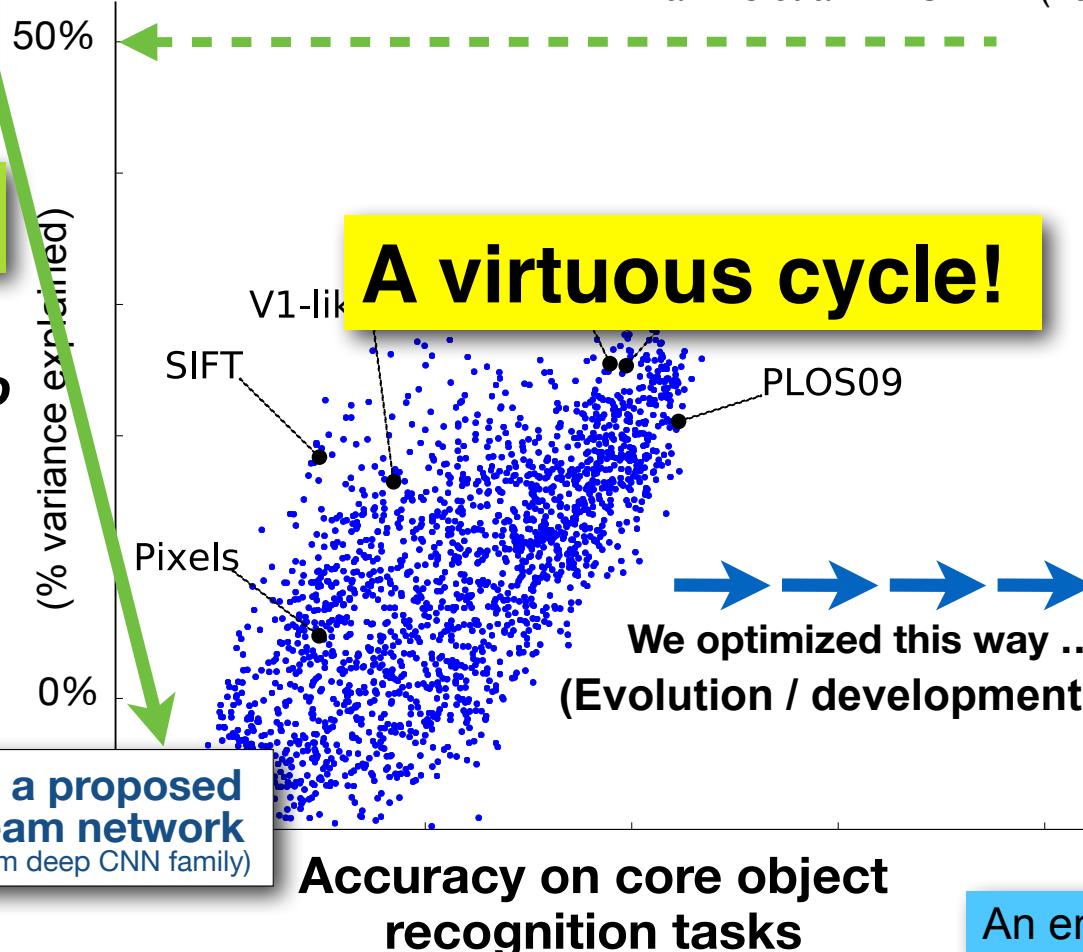


A neuroscience goal

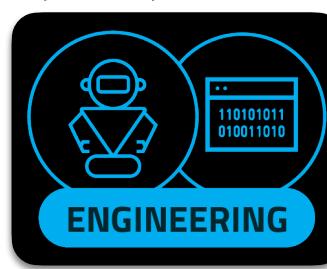
**Match score
of the *in silico*
“IT” neurons
with primate
IT neurons**



Each dot is a proposed
ventral stream network
(here sampled from deep CNN family)



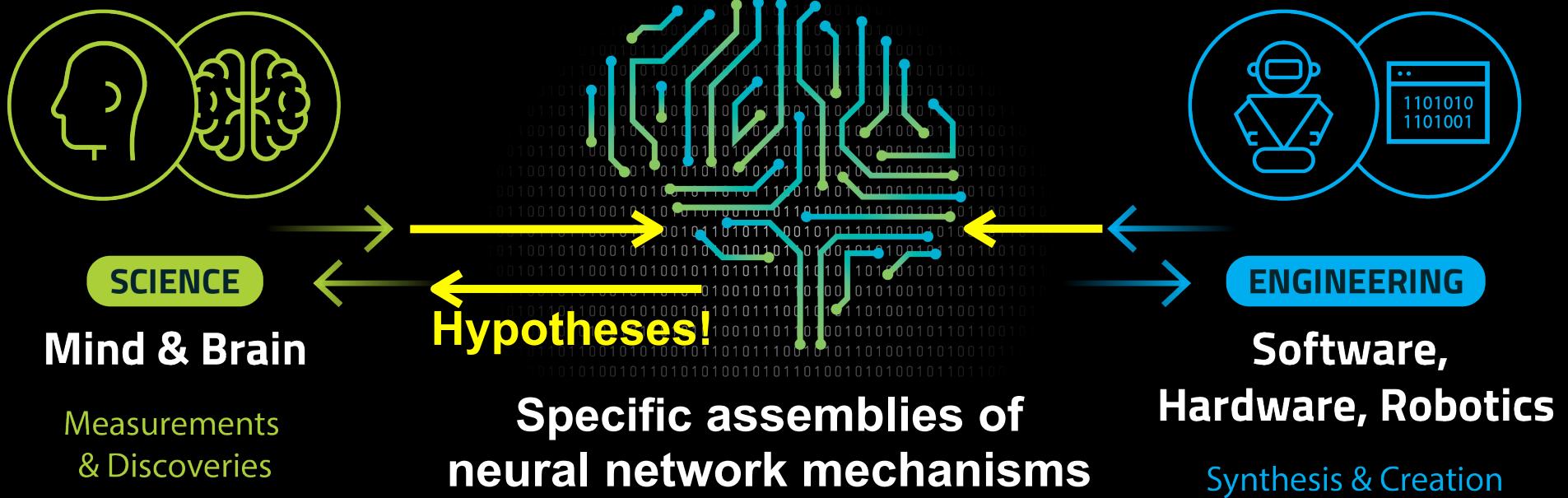
An engineering / AI goal



Can brain science just wait for engineers to build even more accurate models?

Models
Adapted from Yamins, Hong, Solomon,
Seibert and DiCarlo PNAS (2014)

An implicit collaboration!



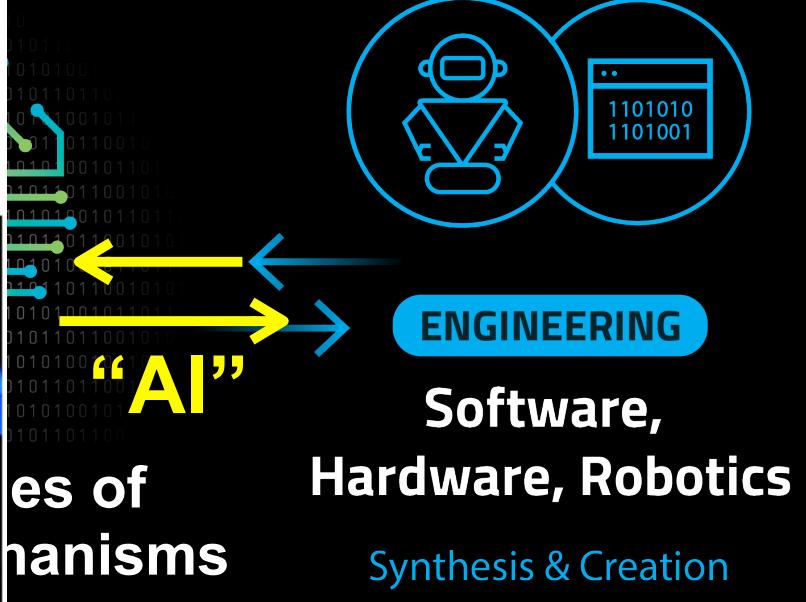
This approach is leading to rapid progress in other areas of brain science:

- Vision** ([Kriegeskorte](#), Oliva, Konkle, Ganguli, Kanwisher, Tsao, ...)
- Audition** (McDermott & Yamins)
- Somatosensation** (Hartmann & Yamins)
- Decision making** (Sussillo & Newsome, Freedman, ...)
- Motor planning and control** (Jazayeri, Batista, Churchland, ...)
- Navigation** (Fiete, ...)

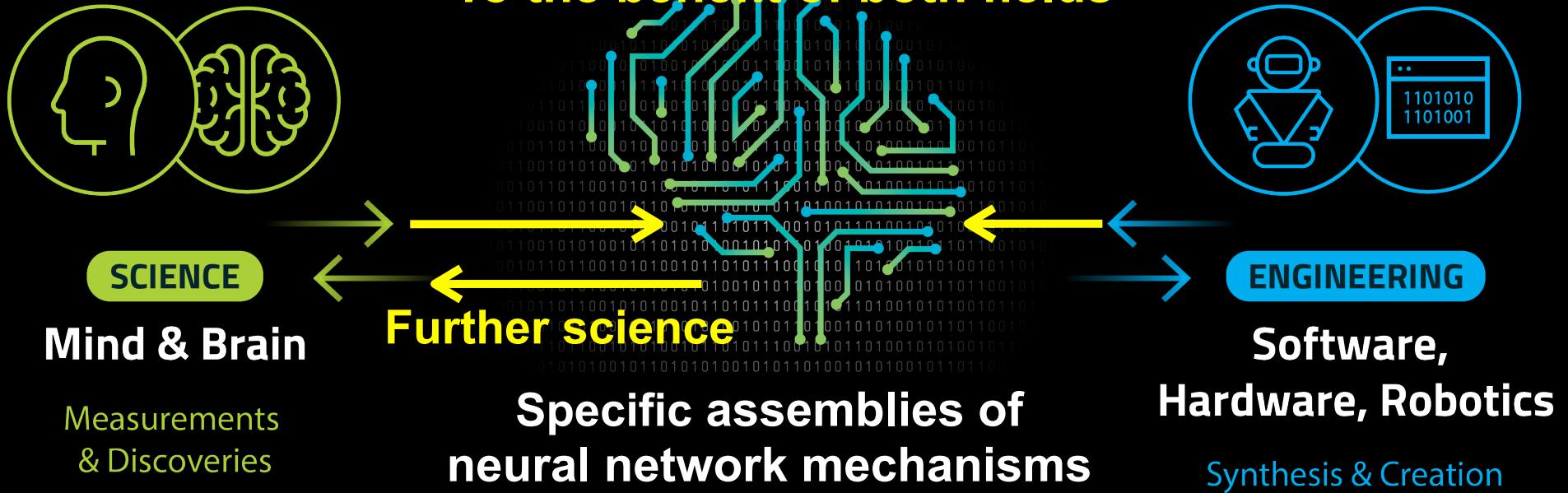
An implicit collaboration! To the benefit of both fields



Using an artificial intelligence technique inspired by theories about how the brain recognizes patterns, technology companies are reporting startling gains in fields as diverse as computer vision, speech recognition and the identification of promising new molecules for designing drugs.



An implicit collaboration! To the benefit of both fields

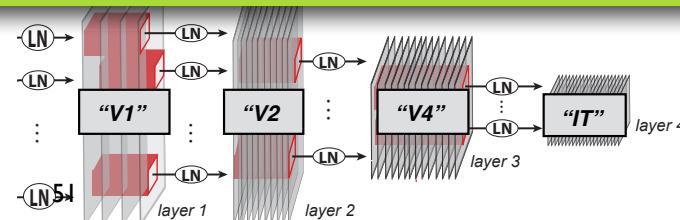


Current leading mechanistic models of the ventral visual stream...



**But, no ANN
model aces all
of our brain and
behavioral tests.**

Particular ANNs can now reasonably accurately
explain / predict the workings of the ventral
visual stream (single-neuron-level & behavioral-level)



Current leading mechanistic models of the ventral visual stream...

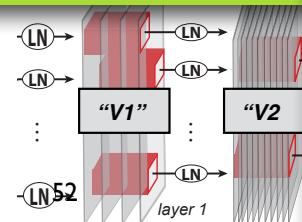


Summary: clear progress, but the current scientific hypotheses (models) are still incomplete (i.e. demonstrably inaccurate in important ways)

But, no ANN

model aces all of our brain and behavioral tests.

Particular ANNs can now explain / predict the whole visual stream (single-neuron



Rank ▼	Model	average	V1: FreemanZU	V2: FreemanZU	V4: Majaj2015	IT: Majaj2015	IT-temporal: Kellam	behavior: Rajah	ImageNet
1	CORnet-S <i>Kubilius et al., 2018</i>	.539	.429	.430	.781	.730	.316	.545	.747
2	resnet-50_v2 <i>He et al., 2015</i>	.509	.478	.532	.780	.710	X	.553	.756
2	resnet-101_v1 <i>He et al., 2015</i>	.509	.455	.524	.775	.738	X	.561	.764
4	densenet-169 <i>Huang et al., 2016</i>	.508	.445	.537	.786	.736	X	.543	.759
4	resnet-101_v2 <i>He et al., 2015</i>	.508	.466	.527	.784	.713	X	.555	.774
6	resnet-152_v1 <i>He et al., 2015</i>	.507	.459	.528	.779	.740	X	.533	.768
6	densenet-201 <i>Huang et al., 2016</i>	.507	.454	.533	.777	.738	X	.537	.772
8	resnet-50_v1 <i>He et al., 2015</i>	.506	.456	.528	.782	.747	X	.526	.752

Current leading mechanistic models of the ventral visual stream...



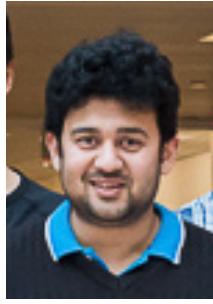
Summary: clear progress, but the current scientific hypotheses (models) are still incomplete (i.e. demonstrably inaccurate in important ways)

What can brain scientists do with the current best mechanistic hypotheses?

How should we improve these mechanistic hypotheses?



Pouya Bashivan



Kohitij Kar

RESEARCH ARTICLE SUMMARY

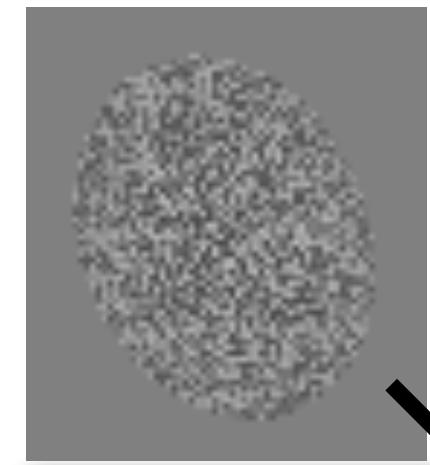
SCIENCE

NEUROSCIENCE

Neural population control via deep image synthesis

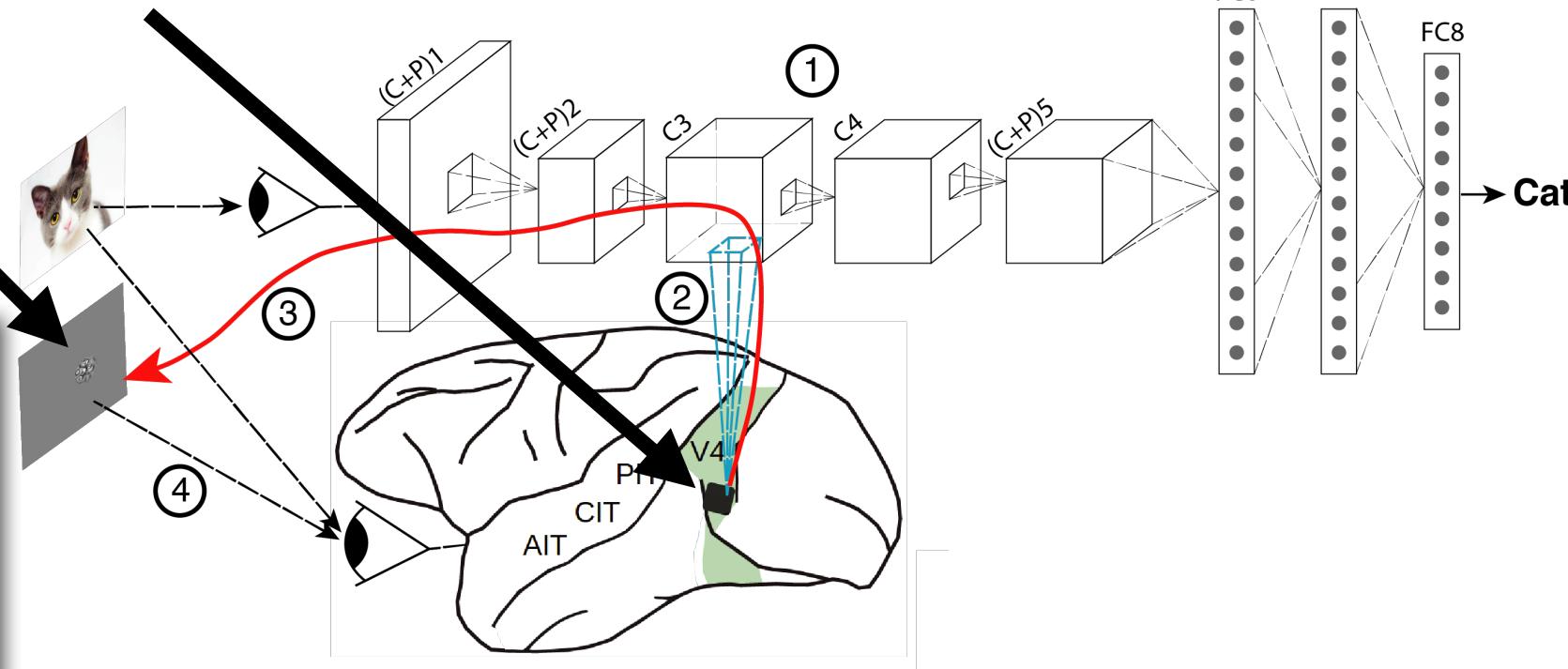
Pouya Bashivan*, Kohitij Kar, James J. DiCarlo

What can brain scientists do with the current best mechanistic hypotheses?



We used our *in silico* ventral stream models to design precise patterns of light energy on the eyes to try to “set” different internal states of the brain.

A control goal: set a desired neural state (of a target brain region)

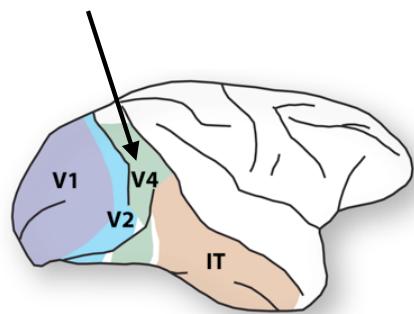


We found that we could use these model-designed images to successfully set neural activity states deep in the brain !

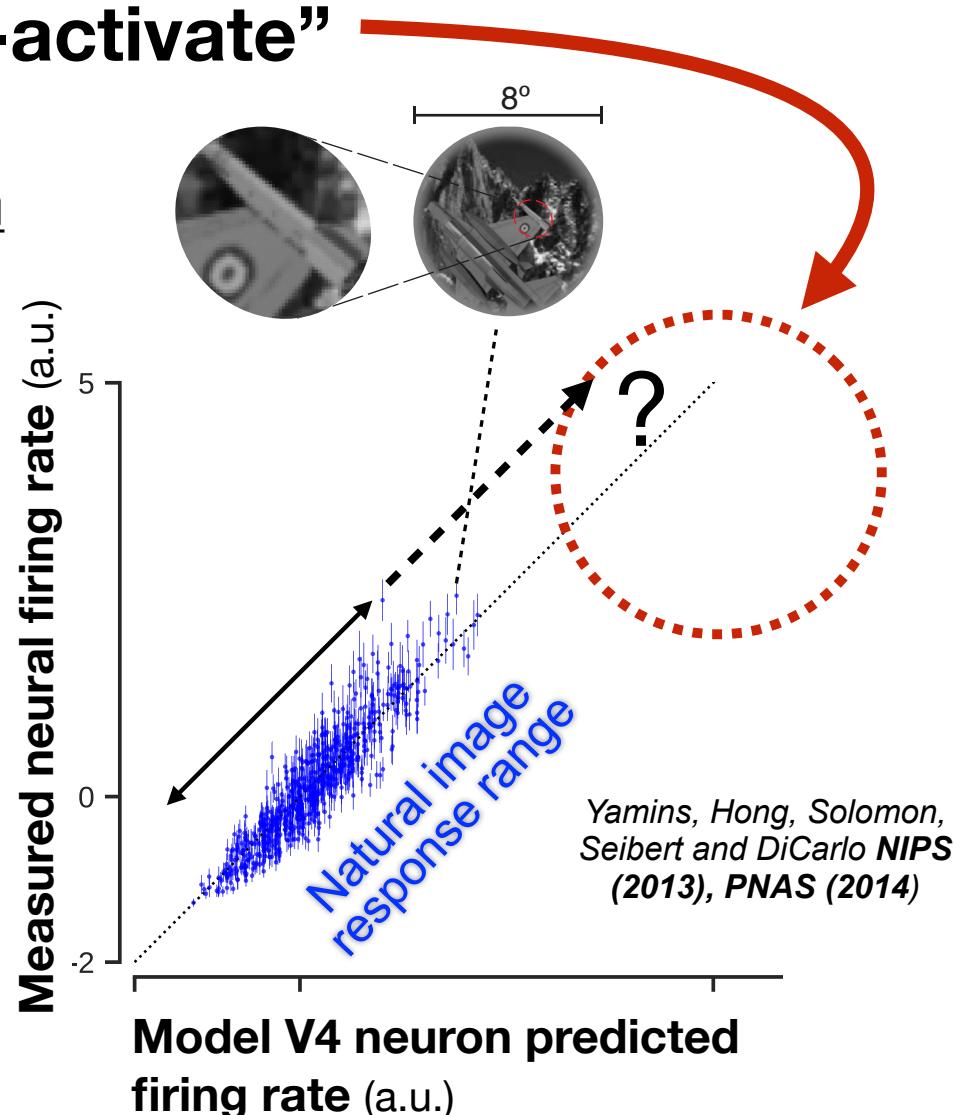
Control goal 1: “super-activate”

Drive any single neural site’s activity beyond the maximum response observed thus far.

Responses of an example V4 neural site:
(a mid-level visual area)



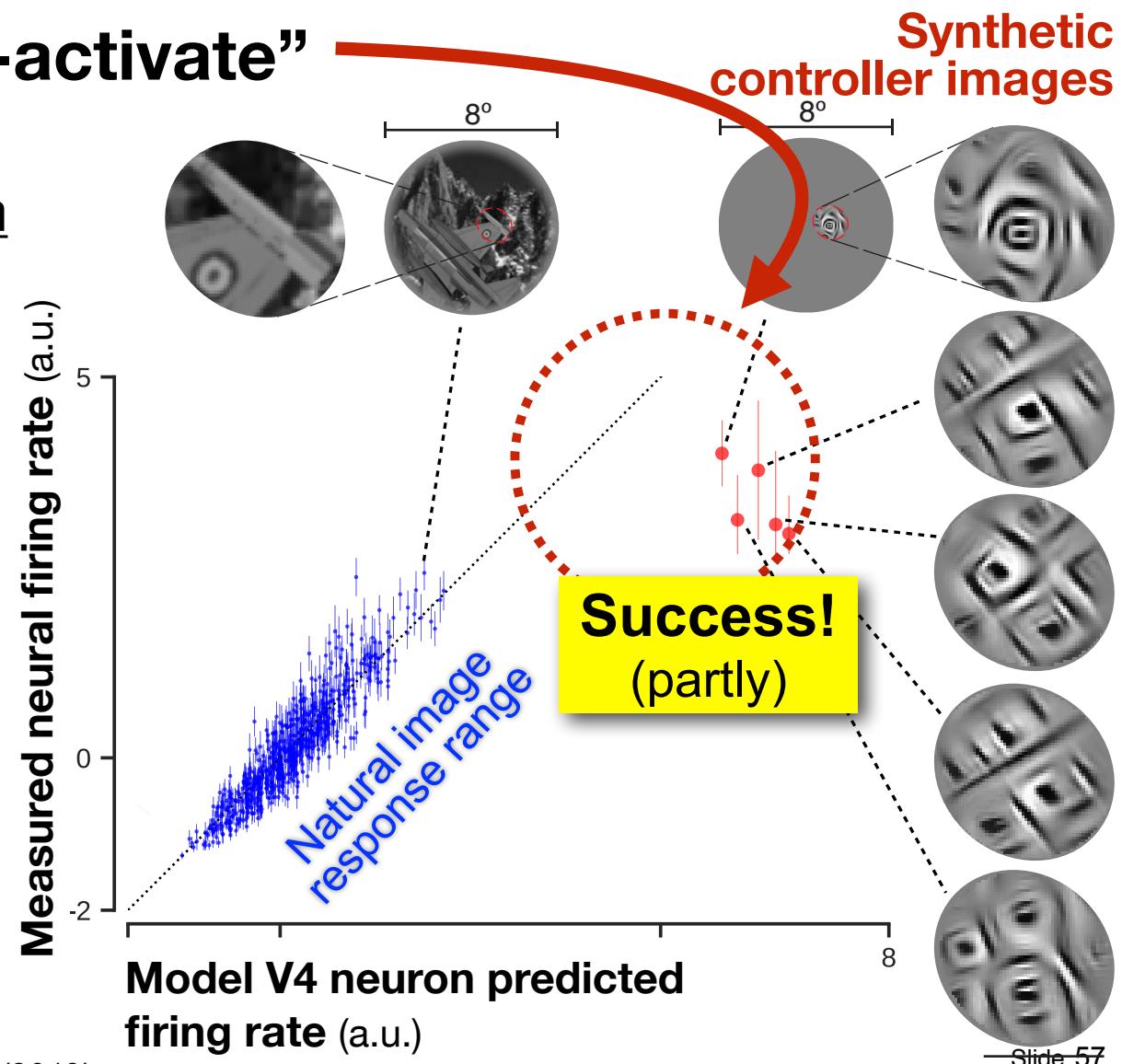
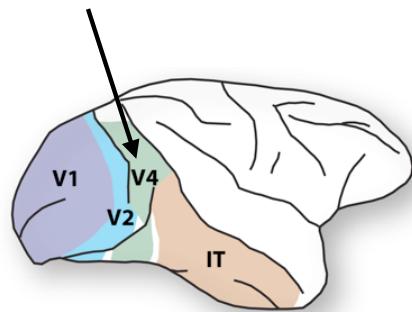
Bashivan, Kar and DiCarlo **CCN** (2018), **Science** 364 (2019)



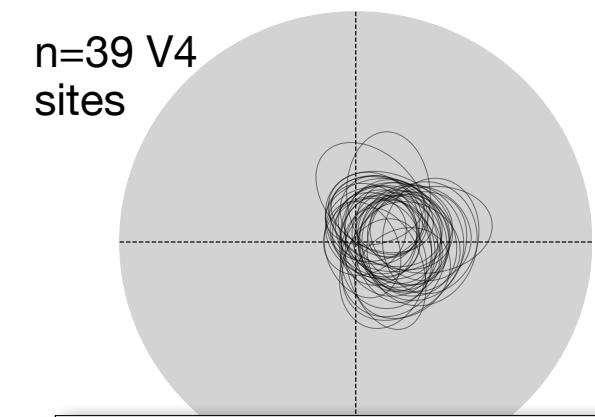
Control goal 1: “super-activate”

Drive any single neural site’s activity beyond the maximum response observed thus far.

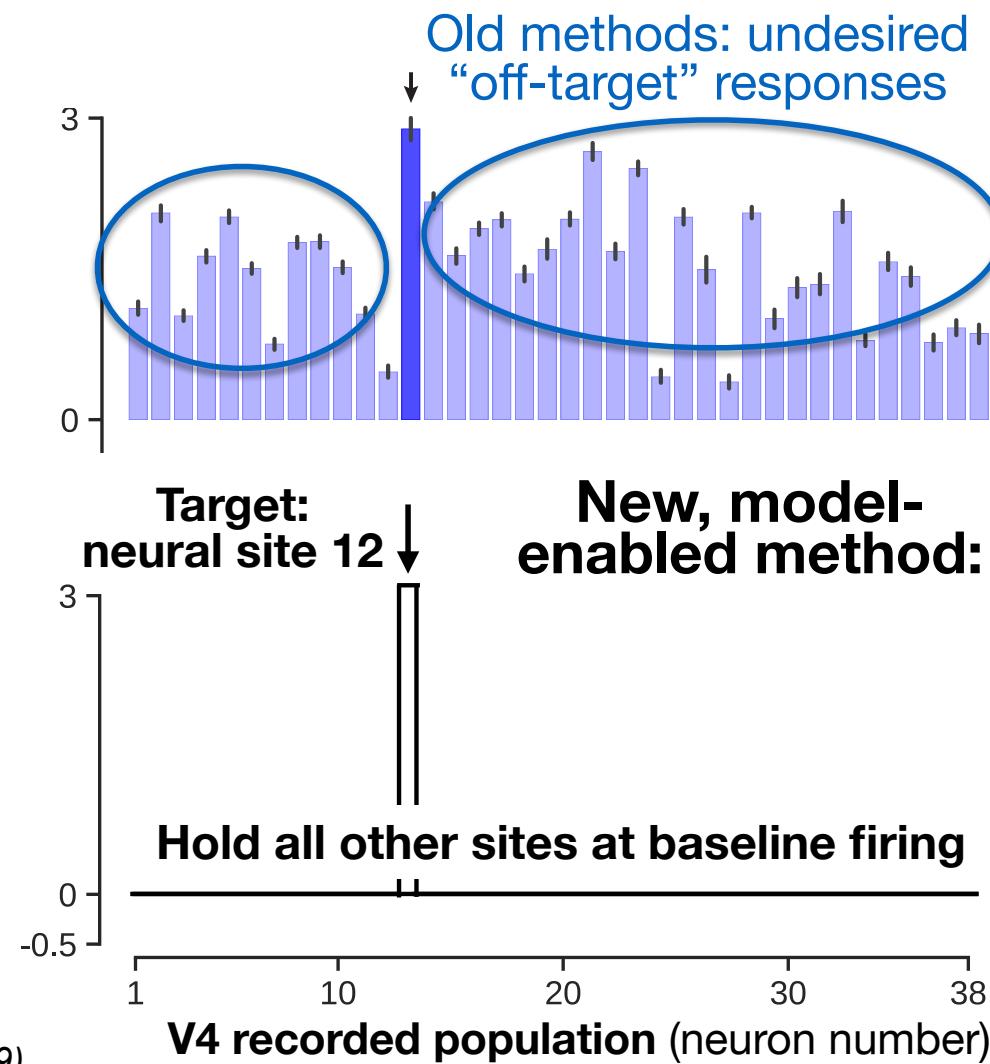
Responses of an example V4 neural site:
(a mid-level visual area)



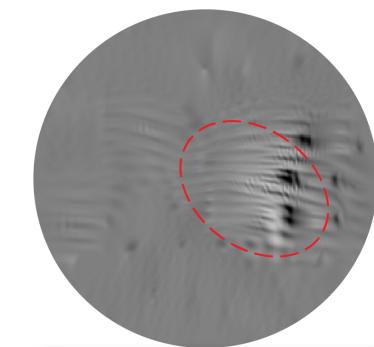
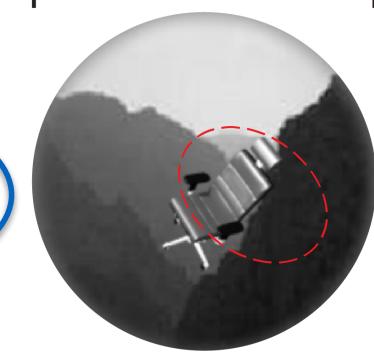
Application: control the neural population state deep in the brain



Example goal:
drive one target neural
site to high activity,
while driving all other
neural sites low.

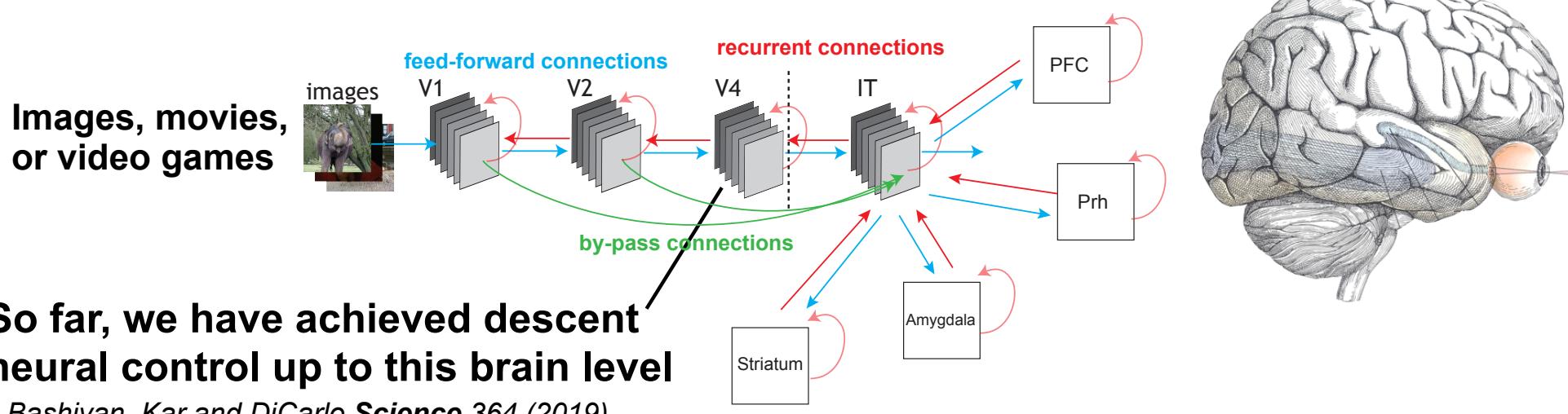


8°



A new application superpower for scientists?:

The ability to control patterns of neural activity deep in the brain by precisely designing the visual input (images, movies)



So far, we have achieved descent neural control up to this brain level

Bashivan, Kar and DiCarlo *Science* 364 (2019)

As we further improve our models, this superpower will further improve!

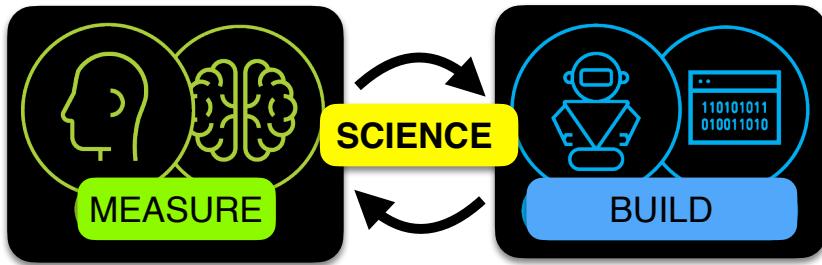
Current leading mechanistic models of the ventral visual stream...



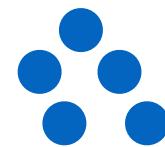
Summary: clear progress, but the current scientific hypotheses (models) are still incomplete (i.e. demonstrably inaccurate in important ways)

**What can brain scientists do
with the current best
mechanistic hypotheses?**

**How should we improve these
mechanistic hypotheses?**



**Test/break
current
hypotheses**



**Current best ventral
stream models**

**Use those
breaks to
build better
hypotheses**

**How should
mechanisms
be used?**

Daballa*, Marques*, Schrimpf, Geiger, Cox & DiCarlo **NeurIPS** (2020)

**Two major lines of ongoing
experimentally-driven work:**

- 1. Study & improve the early components of these ventral stream hypotheses**

**White box
adversarial
attack**

ANN: "Terrier" + Perturbation → ANN: "Church"

in

Amygdala
Striatum

Summary take home messages

1. **Background:** The ventral visual stream produces an IT neural population representation that carries linearly decodable, image generalizable solutions for all (tested) core object recognition tasks.
2. Optimizing deep artificial neural network (ANN) architectures for core recognition tasks leads to internal neural representations in those ANNs that are remarkably similar to the internal neural representations of the ventral visual stream. (*NIPS* 2013, *PNAS* 2014)
 - This result (above) includes IT “face neurons”.
 - This result (above) is consistent with, but does not imply, that the brain learns by classical backpropagation.
3. These same (optimized) ANN models can be used to guide the construction of novel synthetic images to super-activate ventral stream neurons and control sub-populations of neurons (*CCN*, 2018; *Science*, 2019)

Summary take home messages

4. Nevertheless, these same (optimized) ANNs are not yet functionally identical to the ventral visual stream. (e.g. *J Neuroscience*, 2018; “Brain-Score” *bioRxiv* 2018)
5. One difference is the lack of recurrent circuits, and recent IT neurophysiology suggests that fast-acting, automatically-evoked recurrent circuits enable the ventral stream’s superior performance on many images (*Nature Neuro.*, 2019).
6. We and our collaborators are building a series of new models that incorporate more biological constraints — thus far, these model show computer vision gains in efficiency (depth) and gains in robustness to image perturbation.

dicarlo@mit.edu

