The Science of Natural Intelligence: Reverse engineering primate visual perception

AAAI Symposium, Computational Principles of Natural and Artificial Intelligence, March 2017

James DiCarlo MD, PhD

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





The Science of Natural Intelligence: Reverse engineering primate visual perception

AAAI Symposium, Computational Principles of Natural and Artificial Intelligence, March 2017



~3 pounds ~20 watts

Humans have "strong perception." Can we reverse engineer that?



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

We started (~2000) by trying to reverse engineer object detection and categorization

Car Person Building **Tree** Sign Lamp post

Other latent variables pertaining to each object: position, size, pose, etc. Image adapted from MIT Street Scenes Database (Courtesy of Tommy Poggio)

Constraints from brain and cognitive sciences... Central ~10 degrees



Constraints from brain and cognitive sciences... ~200 ms snapshots



Object detection/categorization as solved by primates

Core object perception

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

Not surprisingly, the primate brain excels at core object perception

Core object perception

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

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



Human object categorization accuracy as a function of image viewing time





A stepwise approach to reverse engineering NI

Core visual object perception (central 10 deg, first ~200 ms)



Where do currently* engineered systems fall down (relative to the human primate)? Which human primate measurements are likely to be most informative? (behavior, blood flow, neural activity, anatomy, neural perturbation, subcellular, genetics, etc.)

Can we get those measurements (during system operation)? Q3

How do we forward engineer from such measurements?

Where do (did?) engineered systems fall down (relative to humans)?

Identity preserving image variation



Poggio, Ullman, Grossberg, Edleman, Biederman, etc.

Q1

DiCarlo and Cox, **TICS** (2007), Pinto, Cox, and DiCarlo, **PLoS Comp Bio** (2008), DiCarlo, Zoccolan and Rust, **Neuron** (2012)

2009: Machines vs. humans on our benchmarks

Q1



Data merged here: 48 basic-level tasks (8 labels x 6 level of variation)

2 Decision to gather data from the non-human primate



Ventral visual stream

We know which brain regions house the most critical computations.

We know the system anatomy at a course grain.

We have models of the elemental computations.

We can systematically measure and manipulate every stage of the processing stream, at the level of spiking neural activity, at msec resolution, in behaving subjects. (not currently possible in the human brain)

Adapted from Motter and Mountcastle 1981

Decision

and action

Memory

Decision to gather data from the non-human primate

Q2



Adapted from Motter and Mountcastle 1981











Upshot: learned weighted sums of IT features achieve high performance.



In the language of machine learning: the IT neural population was (and still is) a remarkably powerful set of features.



Methods advance: large scale neuronal recording along the ventral stream







From IT features to performance (behavior). This is easy!

Camel

Dog Rhino

V4

IT neural population patterns

The IT feature set convey's biology's solution to intelligent object sensing.

The IT feature set is a high performance basis for lowdimensional, linear read of object category and identity. (Hung, Kreiman, Poggio, DiCarlo 2005; Rust and DiCarlo 2010)

The IT feature set is computed in less that 200 ms. (Hung, Kreiman, Poggio, DiCarlo 2005; Rust and DiCarlo 2010; Majaj, Hong, Soloman and DiCarlo, 2015)

The IT feature set is a GENERAL BASIS — it immediately supports rapid, linear-read, learning of any new object category and identity. AND it supports report of object position, pose,

scale, etc. (Hong, Yamins, Majaj, DiCarlo 2016)

LGN

RGC

pixel

An engineering spec of the above (LAWS of RAD IT) accurately predicts the difficulty of all tested object recognition tasks in humans and monkeys. (Majaj, Hong, Soloman, and DiCarlo, 2015; Rajalinghan, Schmidt and DiCarlo 2015).



IT





Forward engineering (model building) within Q3known constraints on the ventral stream Ventral stream is a "deep" stack of areas Each neuron is well approximated as a Each area conveys a retinotopic map linear filter + output non-linearity + pooling In each area, the same set of normalization operations is applied at each location in the map (operating on different inputs). ↓↓↓↓↓↓↓↓ V1Those operations depend only on a local set of outputs from the previous area. Multiple filter types are applied Linear Gain Output at each location (e.g. different oriented Gabor filters in V1) control nonlinearity operator Carandini & Heeger, 1994 00000

IT

V4

RGC

pixel

LGN

Forward engineering (model building) within known constraints on the ventral stream

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

Today, this family of models is called "Deep convolutional neural networks" (Deep CNN's)









Forward engineering (model building) within known constraints on the ventral stream

How do we determine which model in the deep CNN family, if any, is the actual mechanism of the ventral stream?

Strategy: Use selection methods to find specific models (i.e. parameter settings) in this model family. Problem:

What to select for? Models that are good at tasks that we hypothesize that the ventral stream evolved and/or developed to solve. Neuroscience data suggest: the task is "invariant" object recognition.

How to select? Biology does not yet tell us, so we used engineering optimization methods.

ach layer

(* not directly determined by neurobiology)









Remarkably ability of these models to explain and predict IT features







Yamins, Hong, Solomon, Seibert and DiCarlo **PNAS (2014**) In 2012, this model elevated deep CNN models to a leadership role in computer vision

Error rate on object categorization



Q3 Forward engineering summary:

neuroscience-constrained architectural family (deep CNN) +
cognitive science-derived task (e.g. invariant categorization) +
engineering optimization (ImageNet, s. gradient descent, etc.)
produces a decent approximation of nature's neural mechanism of intelligent object sensing (in primates).

But these models do not fully predict nature's solution.

And primate behavior (and IT) are still better than current deep CNNs



Invariant categorization performance

Differences between state-of-the-art deep CNNs and Primates

Where do currently engineered systems fall down (relative to humans)?





We tested thousands of images and discovered hundreds of images that are <u>reliably</u> solved by (human AND non-human) primate brains with only 100 ms viewing duration, but not solved by current CV systems (red lines).



Each dot is an image



(*results are very similar for human performance)



Differences between state-of-the-art deep CNNs and Primates

Challenge images

Control images





Acknowledgements

brain+cognitive sciences

nitive sciences

Current lab members:

Pouya Bashivan Elias Issa Kohitij Kar Hyodong Lee Micheal Lee Shay Ohayon Jon Prescott-Roy Rishi Rajalingham Kailyn Schmidt Darren Seibert Christopher Shay Chris Stawarz



Alumni:

Arash Afraz David Cox Chou Hung Gabriel Kreiman Nuo Li Najib Majaj Nicolas Pinto Nicole Rust Ethan Soloman Dan Yamins Davide Zoccolan

Key collaborating labs:

Ed Boyden (MIT) David Cox (Harvard) Bob Desimone (MIT) Tomaso Poggio (MIT) John H.R. Maunsell (U Chicago) J. Anthony Movshon (NYU) Nancy Kanwisher (MIT)

- NIH NEI
- Simons Foundation
- ONR MURI
- IARPA
- McGovern Institute