Untangling object recognition in the ventral visual processing stream

James DiCarlo

The McGovern Institute for Brain Research Department of Brain and Cognitive Sciences Massachusetts Institute of Technology, Cambridge MA, USA



Visual object and face recognition

MASSACHUSETTS INSTITUTE OF TECHNOLOGY PROJECT MAC

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

Courtesy of A

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.



Visual object and face recognition



- Accurate
- Fast
- Tolerant to variation
- Effortless
- Critical to survival

Our mission: Understand how the brain constructs a neuronal representation that underlies object recognition

Focus the problem: core object recognition

The non human primate model

	Dorsal visual stream
We know where that image representation lives in the primate (Inferior temporal cortex, IT).	
<i>We can study that representation (a precursors) at the level of neurona</i>	and its I spikes.











Our primary tools to attack this problem





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









- Position
- Size
- Pose
- Illumination
- Clutter
 - Other objects
 - Background scene

Natural variability

How does the brain recognize each object across this wide range of natural variability?



It must create an **image representation** that is selective for object identity, yet tolerant ("invariant") to this variability.

Explicit object representation

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







Actual pixel representation

(~ retinal image representation)



individual 2



ineffective separating hyperplane



individual 1

(Due to identity-preserving image variation.)

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

Understanding how the brain solves object recognition



Background: IT neurons are rapidly selective (among images)

e.g. Gross, Desimone, Albright, Rolls, Tanaka, Logothetis, Miyashita, Sheinberg, Connor, etc. etc.....



Background: object information is explicit in IT

(*n* ~ 350 IT sites)



Object image



Understanding how the brain solves object recognition



Understanding how the brain solves object recognition



The temporal contiguity of natural experience is a gold mine for learning





Hypothesis: the ventral visual stream uses the temporal contiguity of natural visual experience to **learn** to construct explicit object representation

Pose



same object



Foldiak, 1991; Wallis & Rolls, 1997; Wallis & Bulthoff, 2001; Wiskott & Sejnowski, 2002; Fiser and Aslin, 2002; Perry, Rolls, Stringer, 2006; Wyss, Konig, Verschure, 2006; Sprekeler, Michaelis, Wiskott, 2007; Masquelier & Thorpe, 2007; Masquelier, Serre, Thorpe, Poggio, 2008 etc.

Start simple: focus on position tolerance ("invariance")

Natural visual experience and position tolerance ("invariance")



Retinal image:

Unsupervised retinal experience:



short time



Learn position tolerance?



Cox, Meier, Oertelt and DiCarlo. Nature Neuroscience (2005)



Cox, Meier, Oertelt and DiCarlo. Nature Neuroscience (2005)

Retinal (position) space



Cox, Meier, Oertelt and DiCarlo. Nature Neuroscience (2005)



IT neurons will tend to confuse these objects across these positions

Perceptual prediction:

subjects will tend to confuse these objects across these positions



Experiment 1 (human subject, within subject design)



Background: Role of IT neurons in position invariant object representation



Background: Role of IT neurons in position invariant object representation



Retinal Position















Experimental design (within-neuron, longitudinal design) Ρ Ν Real-time eye tracking and stimulus control **Unsupervised** ! 120 Spikes / s 40 Position (deg) Exposure phase Test Screen phase 100 swaps, ~15 min



Prediction



Results: IT single unit activity



Results: IT single unit activity



Results: pooled population data



Effects of experience in visual worlds with specific alterations in temporal contiguity of object images

• Monkey IT neurons: highly specific change in

position tolerant object selectivity

Li and DiCarlo, Science (2008)



• Very strong:

- effect size is as large or larger than long-term IT plasticity studies (e.g. Baker et al 2003; Sigala & Logothetis 2000; Kobetake et al 1998; Cox and DiCarlo, 2008)

- significant at individual IT sites

Can study learning "online"

More important facts about this result ...

- <u>Cannot</u> be explained as attention or adaptation
- Likely reflects changes in <u>feed-forward</u> processing
- Is object/feature specific
- IT changes last at least weeks (indirect, preliminary)
- May share underlying mechanism with "paired associates" learning (Miyashita 1988, 1991; Erikson et al 1999; Messinger et al. 2001)
- ... but is unsupervised and on time scale of natural vision

"Unsupervised temporal tolerance learning" (UTL)

Inference: temporal contiguity of visual experience "instructs" explicit ("invariant") object representation



Summary

Adult learning --> same mechanisms as infant?

Direct population comparison of the V4 and IT representation

How well can each population (V4 vs. IT) retain object selectivity?

Result: the IT representation improves upon the V4 representation

Summary

"Unsupervised temporal tolerance learning" (UTL)

- Experiments: characterization of UTL
 - size, pose tolerance ? saccade? attention? animal perception? stability? development?

"Unsupervised temporal tolerance learning" (UTL)

- Experiments: characterization of UTL
 - size, pose tolerance ? saccade? attention? animal perception? stability? development?
- Search for cortical untangling transforms

Nicolas Pinto

Pinto, Cox & DiCarlo, PLoS Comp Bol (2008), COSYNE (2008) Pinto, DiCarlo and Cox, ECCV (2008); Pinto, DiCarlo & Cox (submitted)

High-throughput search for a cortical untangling transform

3. Supervised testing of each variant

"Unsupervised temporal tolerance learning" (UTL)

- Experiments: characterization of UTL
 - size, pose tolerance ? saccade? attention? animal perception? stability? development?
- Search for cortical untangling transforms

day 1

"Unsupervised temporal tolerance learning" (UTL)

- Experiments: characterization of UTL
 - size, pose tolerance ? saccade? attention? animal perception? stability? development?
- Test real-world performance of these ideas

4. View of the solution 3. New empirical results

2. Data transformations

1. Focus on the crux problem

Acknowledgements

Chou Hung

Gabriel Kreiman

David Cox

Nicolas Pinto

Nicole Rust

MIT:

Tomaso Poggio Nancy Kanwisher

-

Paul Aparicio Jennifer Deutsch Elias Issa Najib Majaj Marie Maloof Davide Zoccolan

Ben Andken Philip Meier Nadja Oertelt Hans Op de Beeck Dan Oreper Alex Papanastassiou **MGH:** Wim Vanduffel

Caltech: Christof Koch Rodrigo Quiroga

Harvard/HHMI: John Maunsell

• NIH NEI

- NIH NIMH
- DARPA / ONR
- Pew Trust
- McKnight Foundation
- McGovern Institute

