# From feedforward vision to natural vision:

**The impact of free viewing and clutter on monkey inferior temporal object representations**

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

# The core problem of object recognition



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

How does the brain recognize each object across this wide range of conditions?
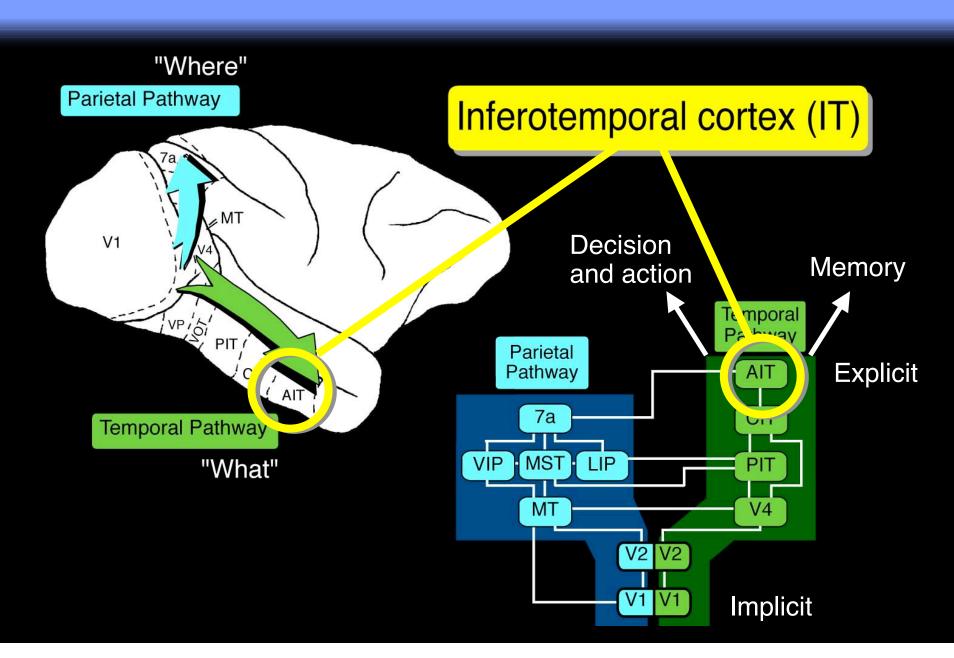
One needs an image representation that is selective for object identity, yet tolerant to such transformations.
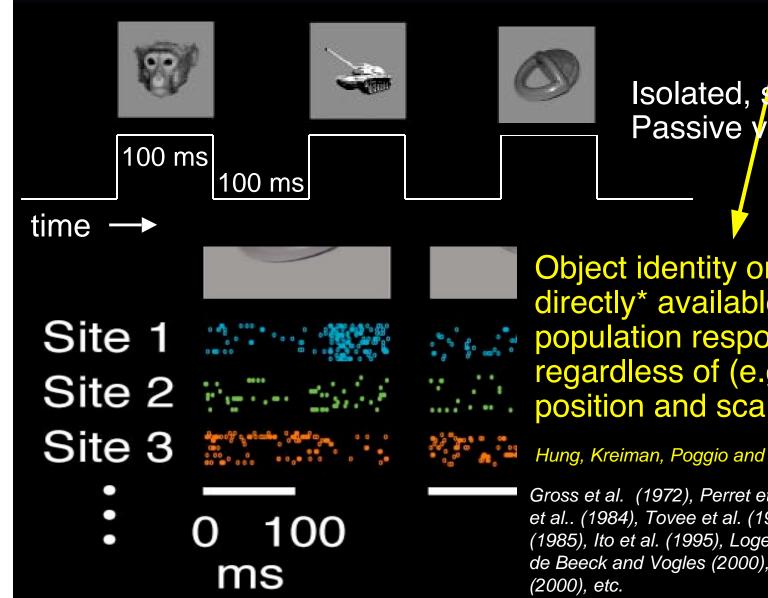
# Rhesus monkey model

We have some idea of where we can find such an image representation (IT).

We can study it at the most appropriate level of abstraction (neuronal spikes).

# Monkey visual system

# AIT contains a rapidly evoked, explicit object representation



100 ms

100 ms

time →

Isolated, single objects. Passive viewing.

Site 1
Site 2
Site 3

0   100
ms

Object identity or category is directly* available in the population response, regardless of (e.g.) object position and scale.

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

*Gross et al.  (1972), Perret et al. (1982), Desimone et al.. (1984), Tovee et al. (1984), Schwartz et al (1985), Ito et al. (1995), Logethetis et al. (1996), Op de Beeck and Vogles (2000), DiCarlo and Maunsell (2000), etc.*

# Feedforward* representation (The Core)

*First evoked pattern of IT activity
when an image is presented to the eye

The Core is fast.

The Core is powerful.

The Core is not yet understood.

Mechanisms ?    Role in "natural vision" ?
(Is it generalizable?)

# The Core and "natural vision"

What is "natural vision" ?

"You know it when you see it."

# The Core and "natural vision"

# The Core and "natural vision"

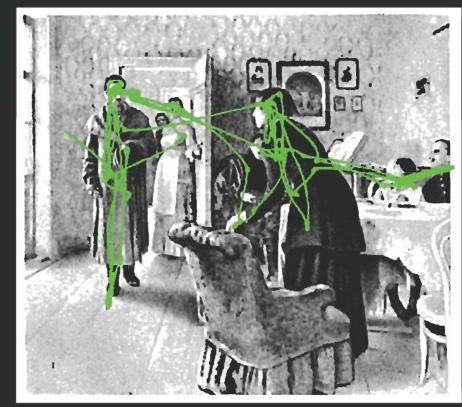How does "natural vision" challenge the basic model
of core vision?



1) **Eye movements** ( "free viewing" )

► 2) **Clutter / Scene / Context**: objects appear among other objects and on backgrounds

3) **Goal directed** (e.g. feature and spatial attention, motor preparation to act, arousal)
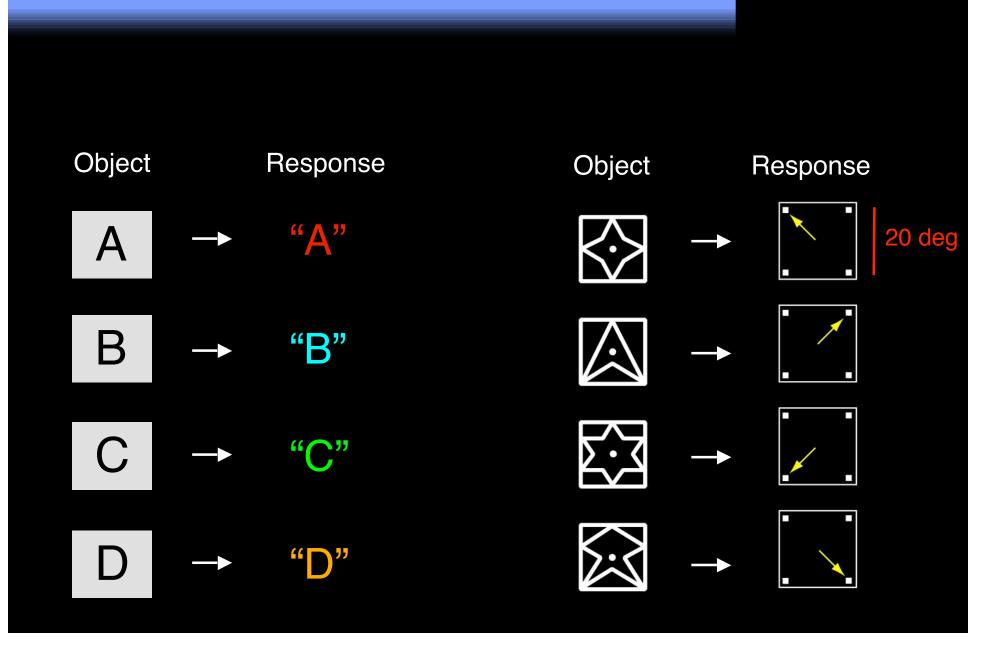
# Natural vision:  Eye movements
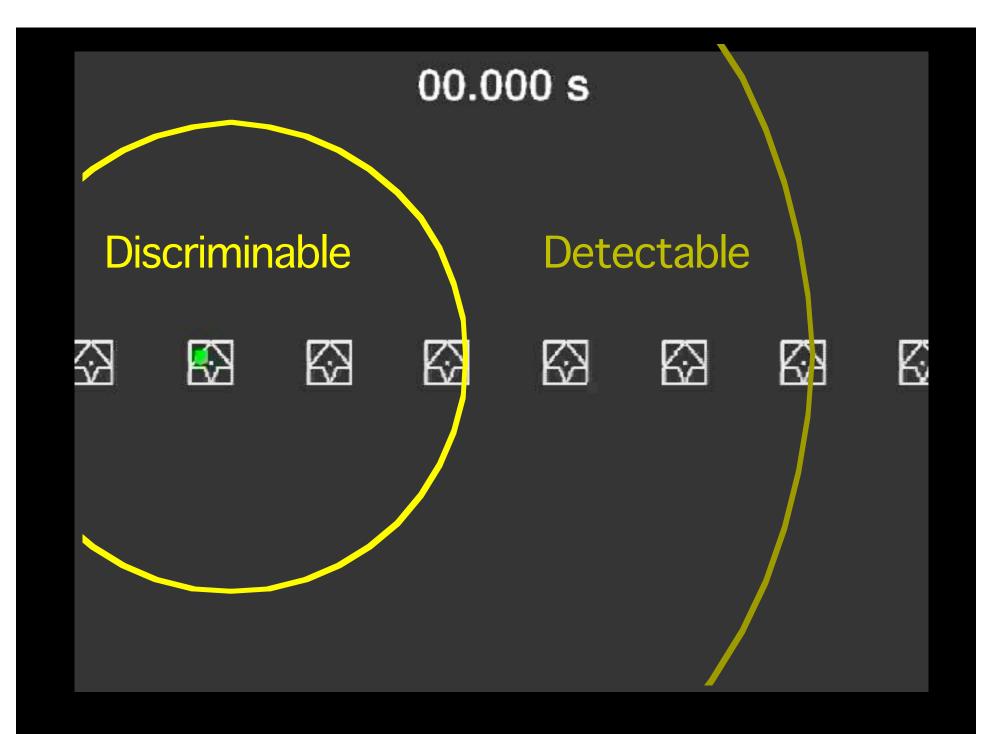
In the lab …

In the real world …



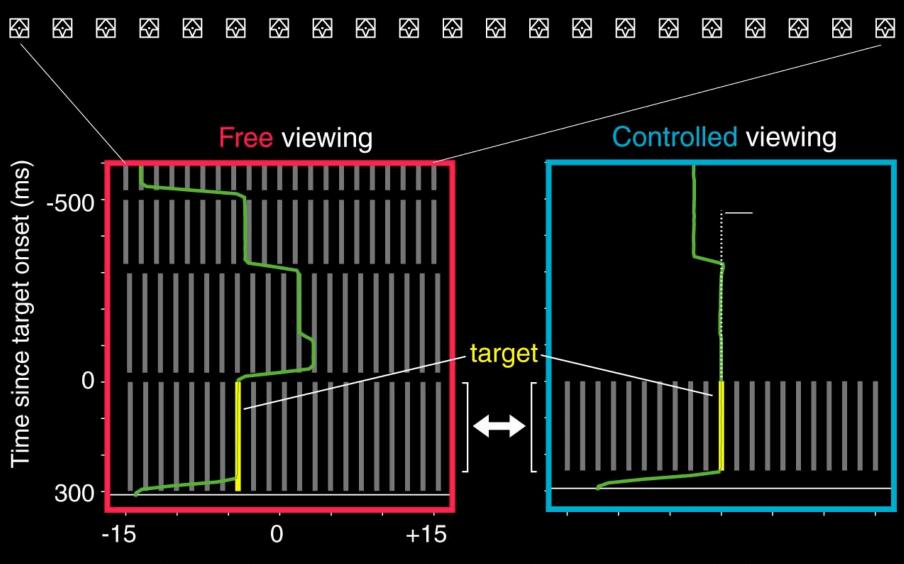*(Adapted from Yarbus, 1967)*
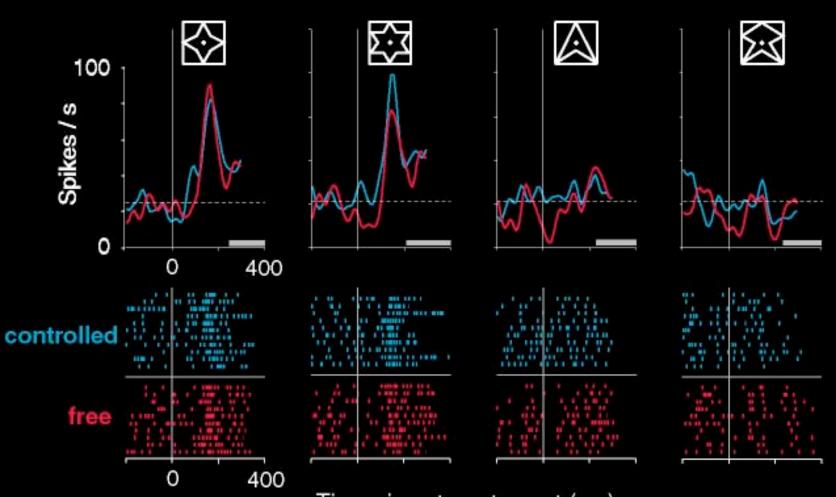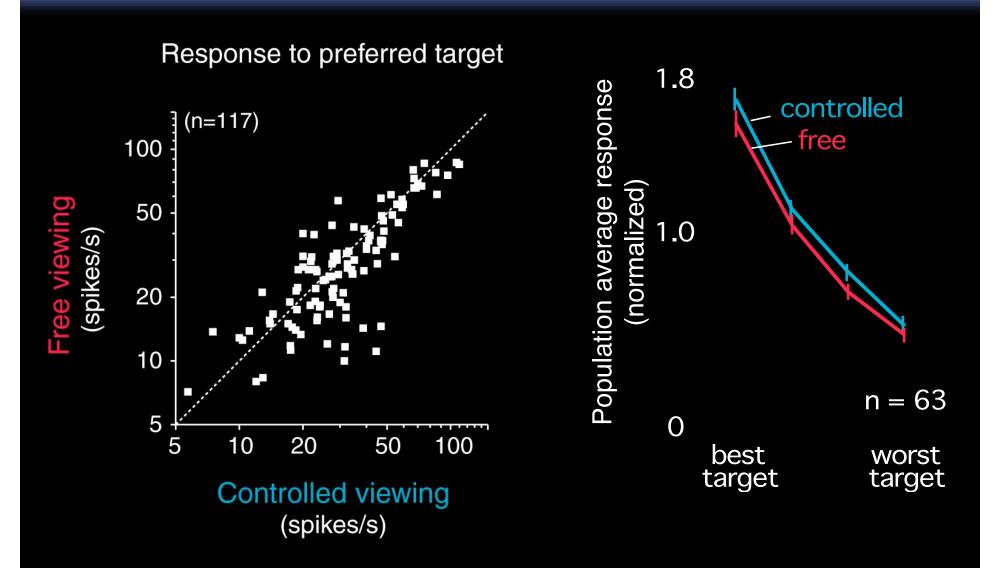
# Object identification task

00.000 s

Correct: 0

Free viewing    Controlled viewing

Time since target onset (ms)

-500

0

300

target

Horizontal eye position (degs relative to screen center)

-15    0    +15

# Example IT neuron

# IT Population summary



Response to preferred target

(n=117)

Free viewing (spikes/s)

Controlled viewing (spikes/s)

Population average response (normalized)

controlled
free

best target — worst target

n = 63

# IT responses are nearly identical in controlled and free viewing conditions

*DiCarlo and Maunsell, Nature Neuroscience, 3: 814-821 (2000)*

*DiCarlo and Maunsell, J Neurophysiology (2005)*

# The Core and "natural vision"

How does "natural vision" challenge the basic model
of core vision?



1) **Eye movements** ( "free viewing" )

2) **Clutter / Scene / Context**: objects appear among other objects and on backgrounds

3) **Goal directed** (e.g. feature and spatial attention, motor preparation to act, arousal)

Not much to worry about here.
*DiCarlo and Maunsell, 2000*
*Sheinberg and Logothetis, 2001*

# Natural vision: Clutter, scene, and context

In the real world…

In the lab…

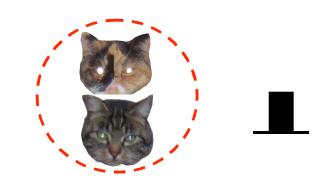# Natural vision:  Clutter, scene, and context

IT Receptive Field

# Long term goal: Understand IT in clutter

**IT responses to object are typically reduced when additional objects are presented**

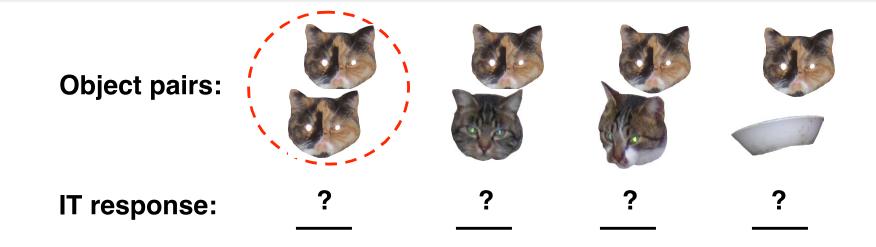*(Sato, 1989; Miller et al., 1993; Rolls and Tovee, 1995; Chelazzi et al., 1998; Missal et al., 1999)*

# First open questions …



Object pairs:

IT response: ? ? ? ?

- Any **systematic relationship** between:

  – response to an object pair
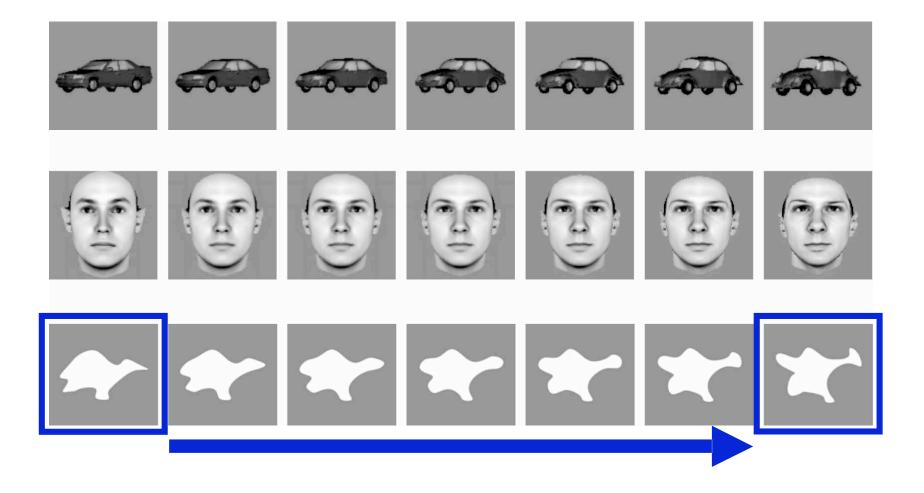  – responses to the constituent objects?

# Experimental design
## overview

- Davide Zoccolan and David Cox
- Recorded IT neuronal responses to the presentation of:

  – **Single** objects
  – **Pairs** of objects
  – **Triplets** of objects

  – In **three monkeys**
  – Using **two complementary experiments**

# Experiment 1

# Experiment 2

# Stimulus conditions

## EXPERIMENT 1



**Single objects**

2 deg

**Object pairs**

## EXPERIMENT 2

# Stimulus conditions

## EXPERIMENT 1



Single objects

2 deg

Object triplets

## EXPERIMENT 2

# Core response: Rapid visual presentation



- Stimuli presented at **5 per sec**
- Passive viewing

**104 neurons** recorded in three monkeys

# Example IT neuron

# Example IT neuron

# Example IT neuron



*Zoccolan, Cox and DiCarlo, 2005*

# Population analysis



**Pairs** (*n* = 79)   *r* = 0.92

**Triplets** (*n* = 48)   *r* = 0.91

Response to multiple objects (spikes/s)

Sum

Average

Sum of responses to single objects (spikes/s)

*Zoccolan, Cox and DiCarlo, 2005*

# Summary: The Core and multiple objects

Under the conditions described here:

- An "average rule" is a very good predictor of the response of individual IT neurons
  (explains **~63%** of response variance → *r ≈ 0.8*)

- => The response pattern of The Core can be predicted by the response pattern to each constituent object

- => useful for supporting the simultaneous representation of multiple objects

# The Core and "natural vision"

How does "natural vision" challenge the basic model of core vision?



1) **Eye movements** ( "free viewing" )

2) **Clutter / Scene / Context**: objects appear among other objects and on backgrounds

3) **Goal directed** (e.g. feature and spatial attention, motor preparation to act, arousal)

Not much to worry about here.
*DiCarlo and Maunsell, 2000*
*Sheinberg and Logothetis, 2001*

Very important challenge.
Beginnings of a systematic understanding.
*Zoccolan, Cox and DiCarlo, 2005*