# CS375 / Psych 249:
## Large-Scale Neural Network Models for Neuroscience
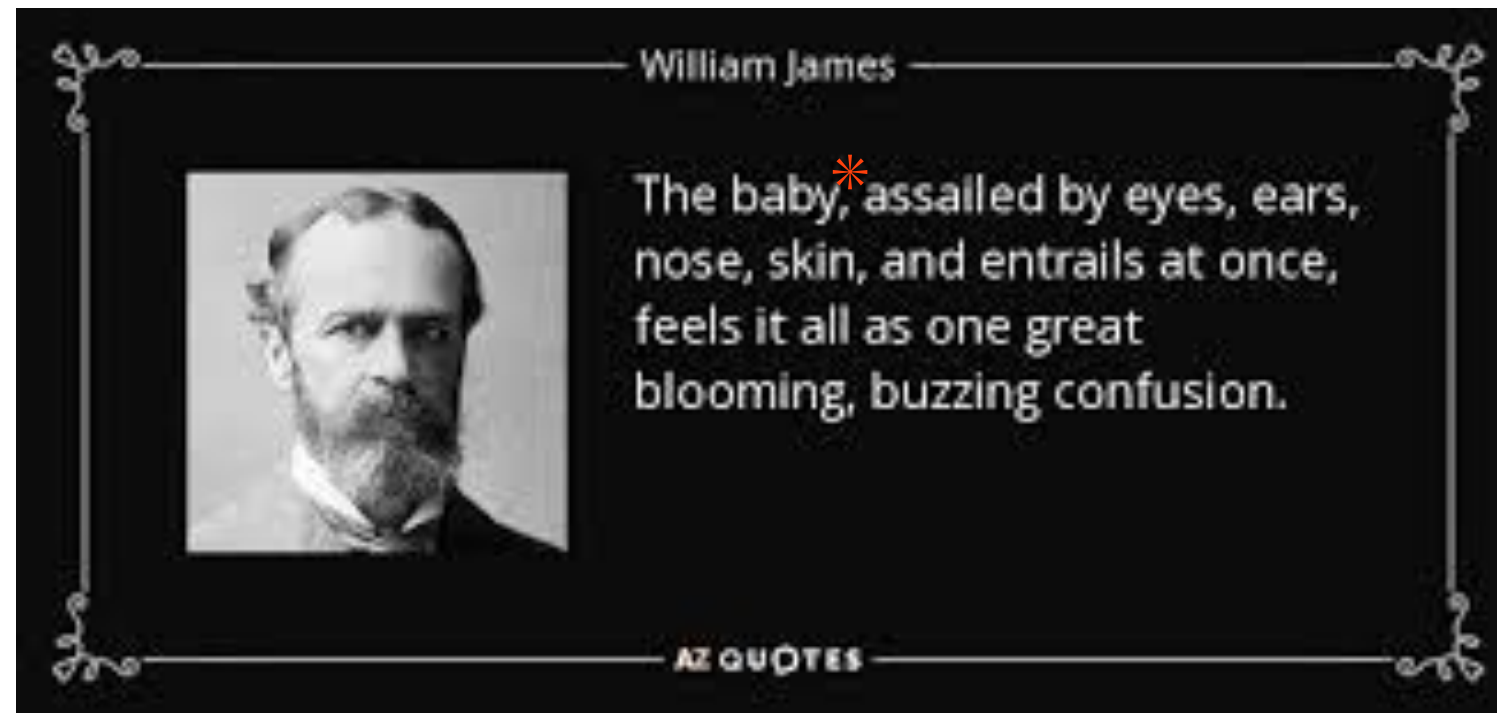
Lecture 2: The biological inspiration of CNNs

*2025.01.07*

Daniel Yamins

Departments of Computer Science and of Psychology
Stanford Neuroscience and Artificial Intelligence Laboratory
Wu Tsai Neurosciences Institute
Stanford University

# Problem: Entity Extraction

Understanding complex, noisy data streams is a critical part of cognition.



William James

The baby,* assailed by eyes, ears, nose, skin, and entrails at once, feels it all as one great blooming, buzzing confusion.

AZ QUOTES

Without sophisticated parsing and entity extraction, the world would be "as one great blooming, buzzing confusion" (for babies or otherwise).
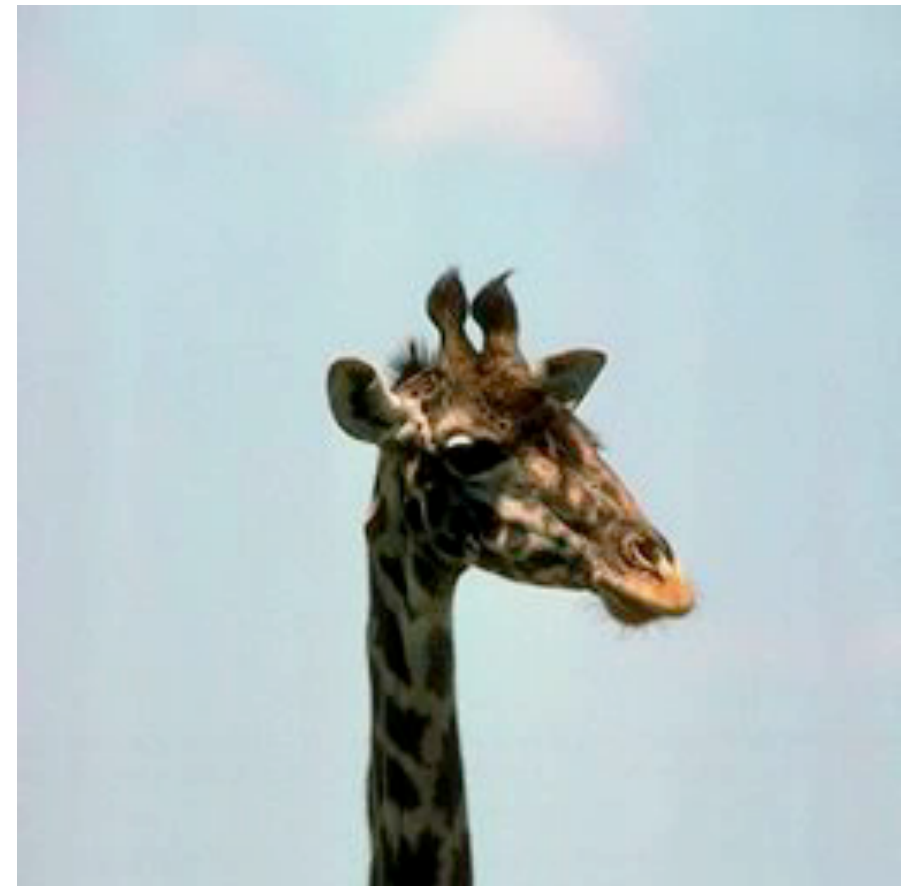
*actually not clearly true for babies …

Why is the problem hard computationally?

1. Nonlinear misalignment between physical and behavioral dimensions
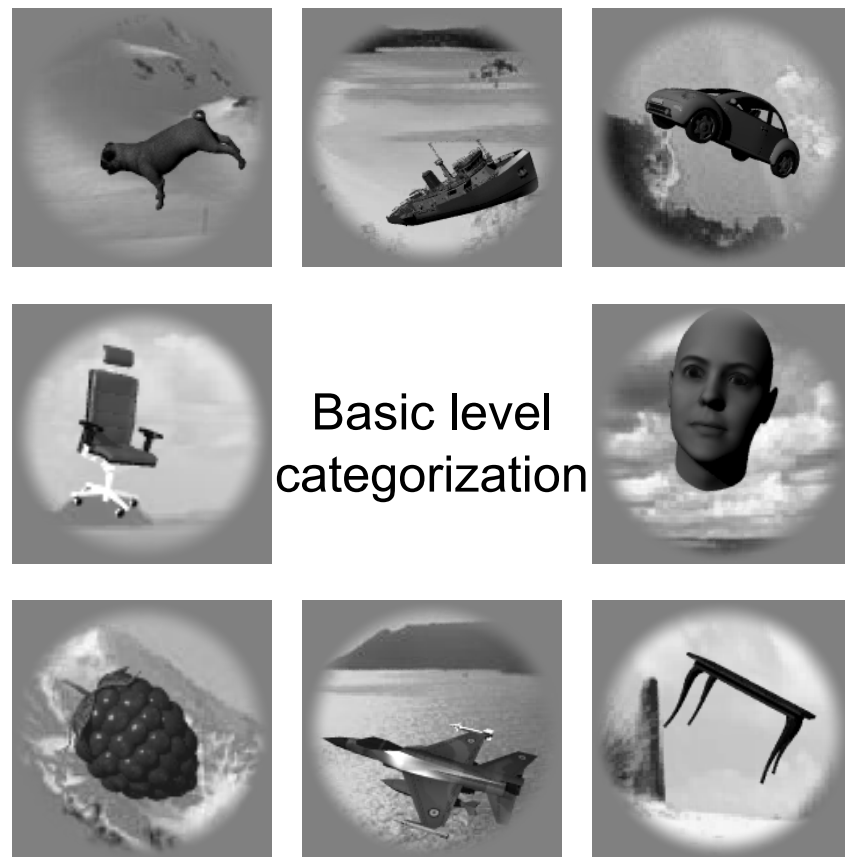
Why is the problem hard computationally?

1. Nonlinear misalignment between physical and behavioral dimensions

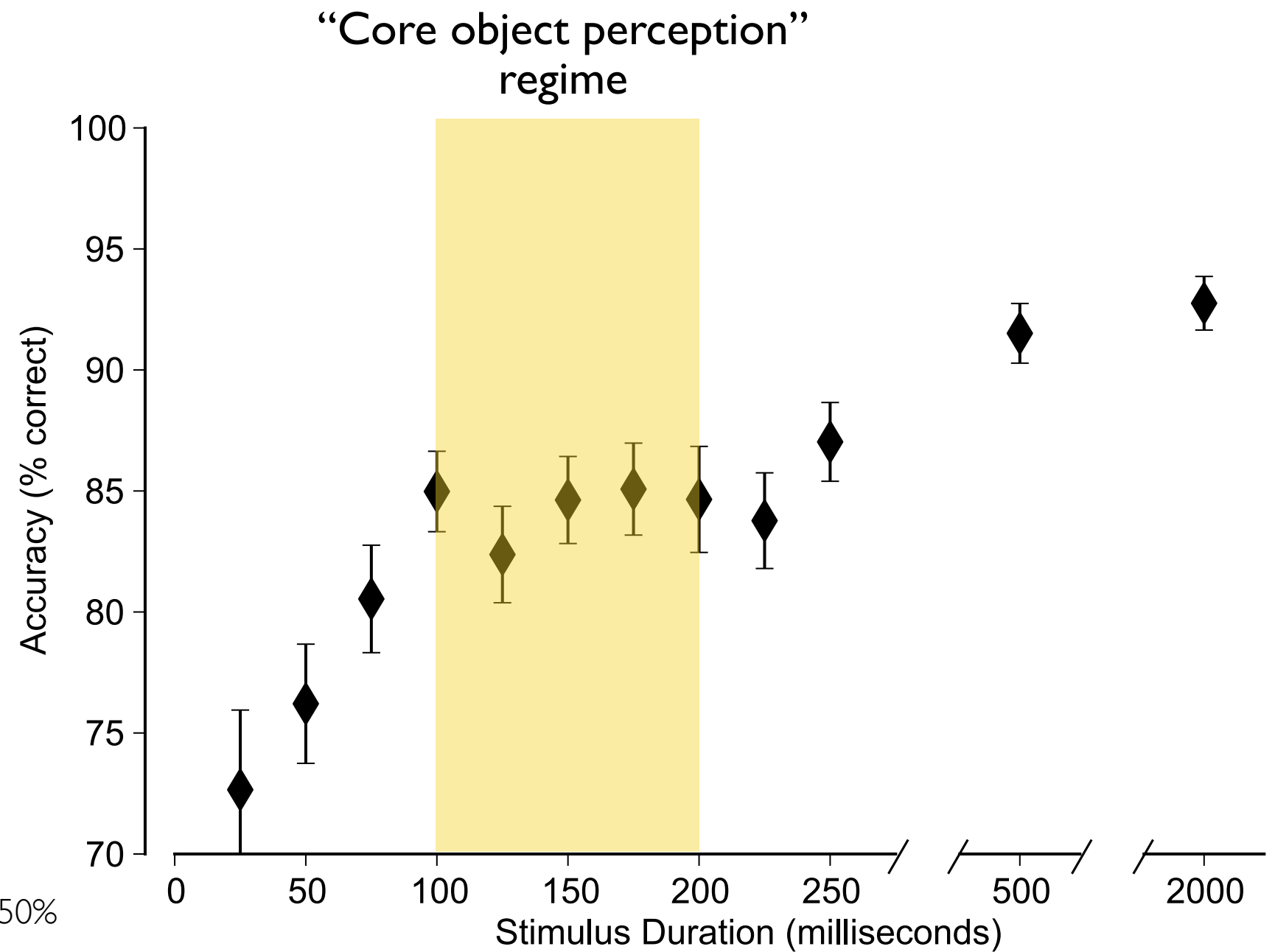2. Needs to be done ***fast***, and thus, presumably, massively in parallel
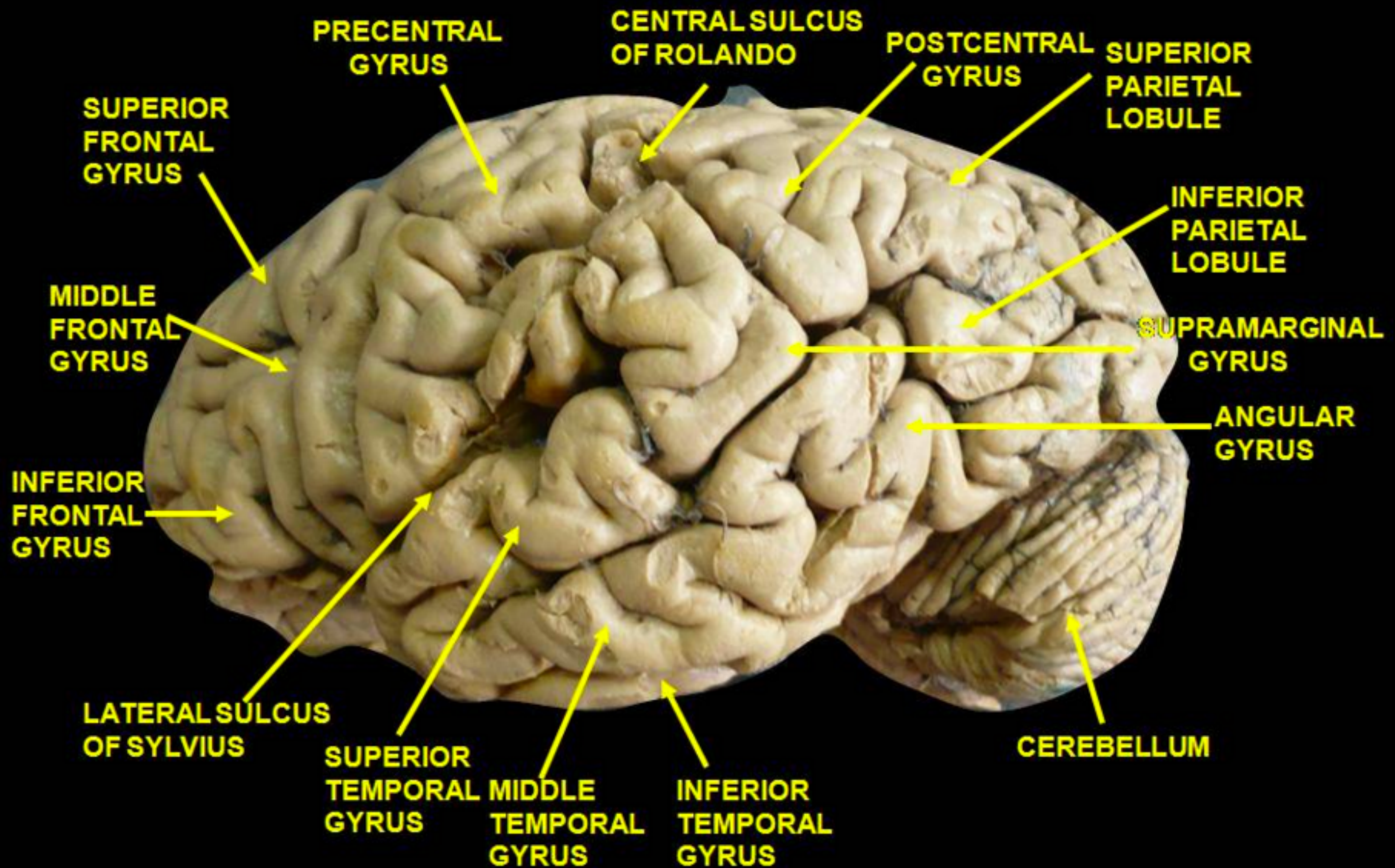
# Problem: Entity Extraction
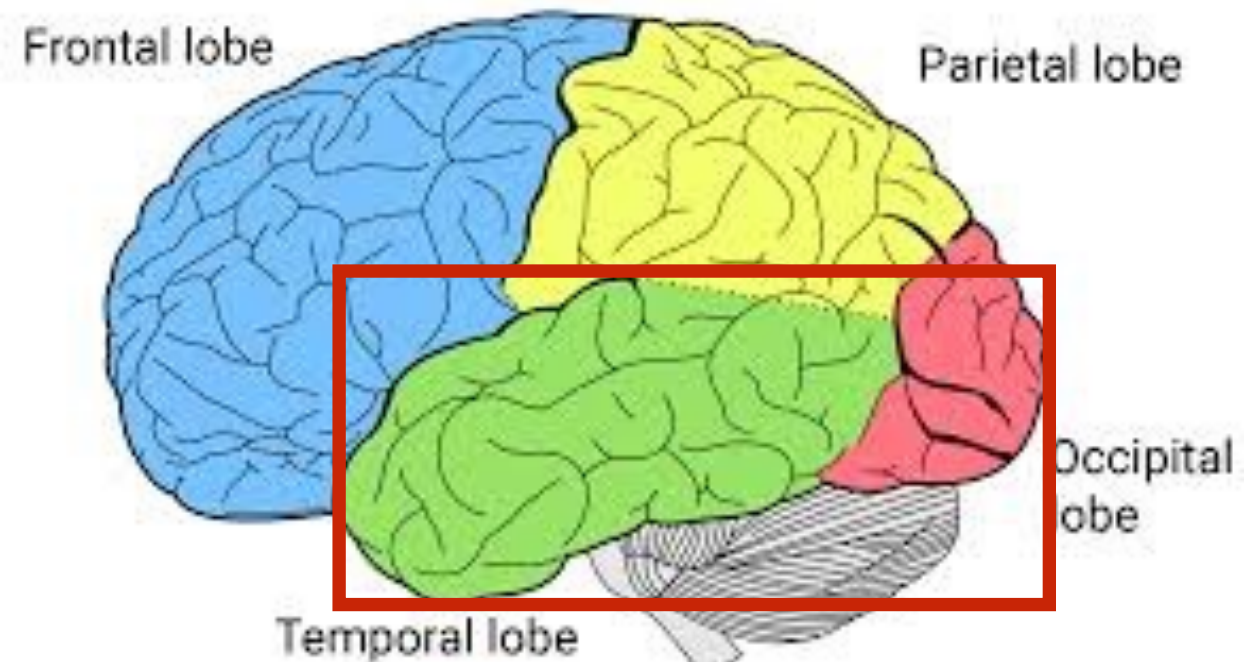
‣ Sensory processing
  - visual, auditory, somatosensory recognition (occipital, temporal)
  - navigation (hippocampus?)

‣ motor command production & execution (motor cortex)

‣ memory, decision making and planning (hippocampus, prefrontal cortex)

‣ language

‣ emotions, theory-of-mind

Frontal lobe   Parietal lobe

Occipital lobe

Temporal lobe

▸Sensory processing
  - visual, auditory, somatosensory recognition (occipital, temporal)
  - navigation (hippocampus?)

▸ motor command production & execution (motor cortex)

▸ memory, decision making and planning (hippocampus, prefrontal cortex)

▸ language

▸ emotions, theory-of-mind

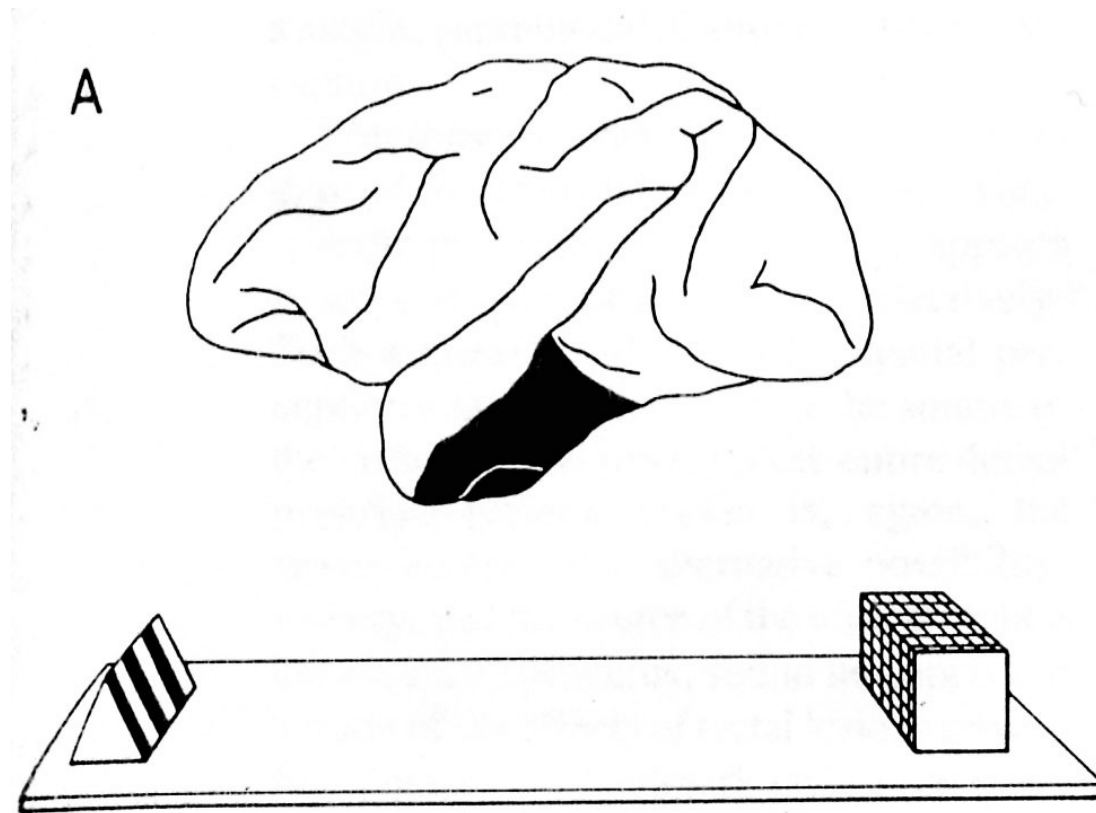Frontal lobe

Parietal lobe

Occipital lobe

Temporal lobe

# Background: Ventral visual stream

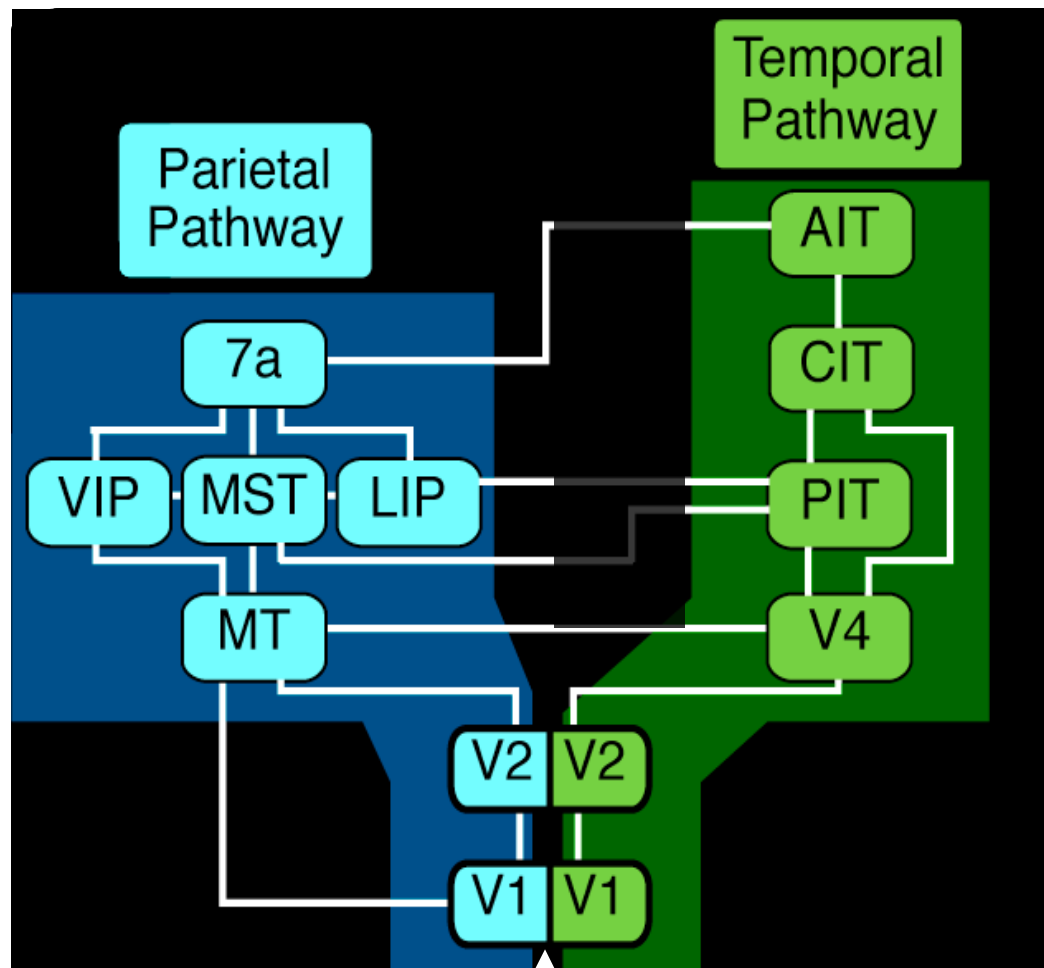Mishkin & Ungerleider, 1982



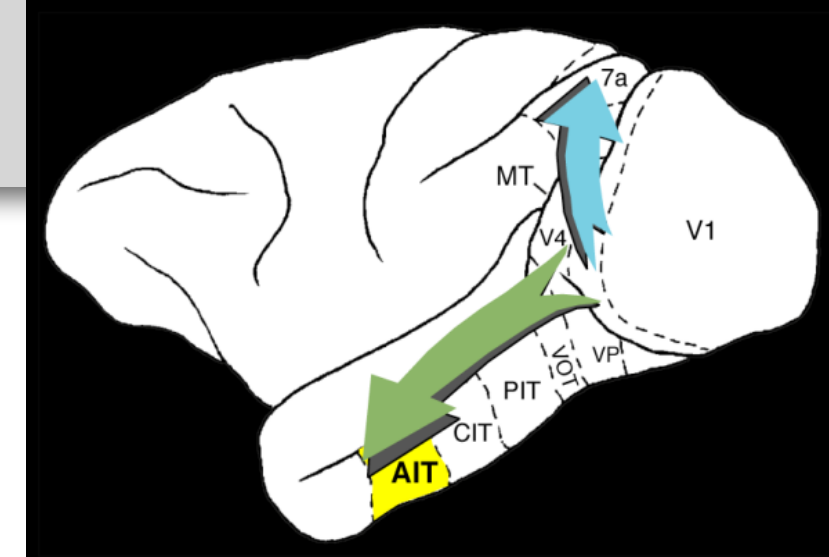**"what" / ventral / occipitotemporal**

*Lesions in IT cortex produce
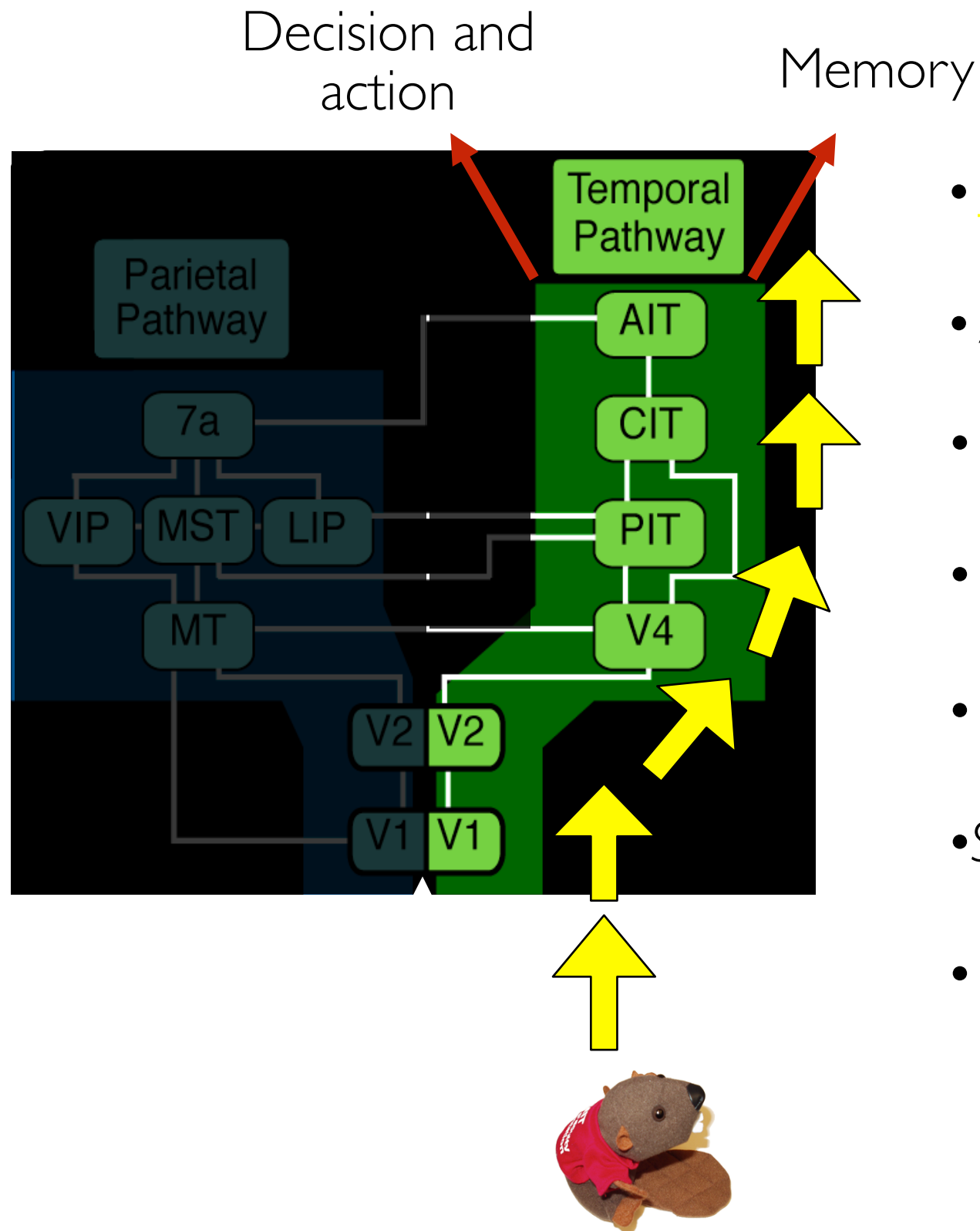deficits in shape discrimination tasks
(Gross et al, 1973, Mishkin 1982)*

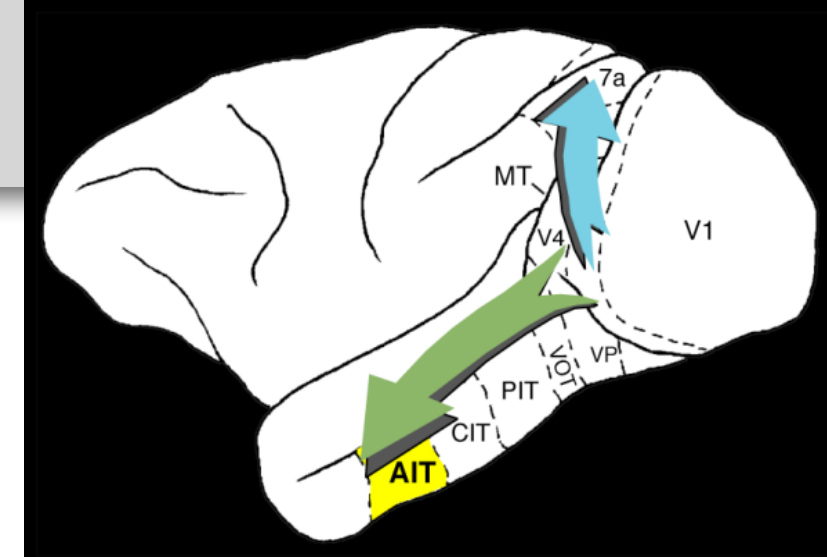**"where" / dorsal / parietal**

*Lesions in parietal cortex produce
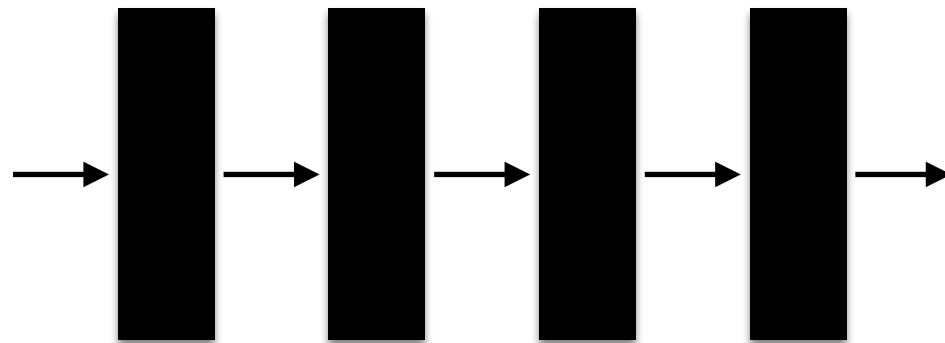deficits in landmark task
(Pohl et al. 1973)*

Decision and action

Memory

- *Tolerance to identity-preserving transforms*

- *Ability to support visual recognition*

- *Correlation with perceptual report*

- *Sensitivity to behavioral state* (e.g. attention)

- *Visually-evoked latency*

- *Selectivity to visual "feature" conjunctions*
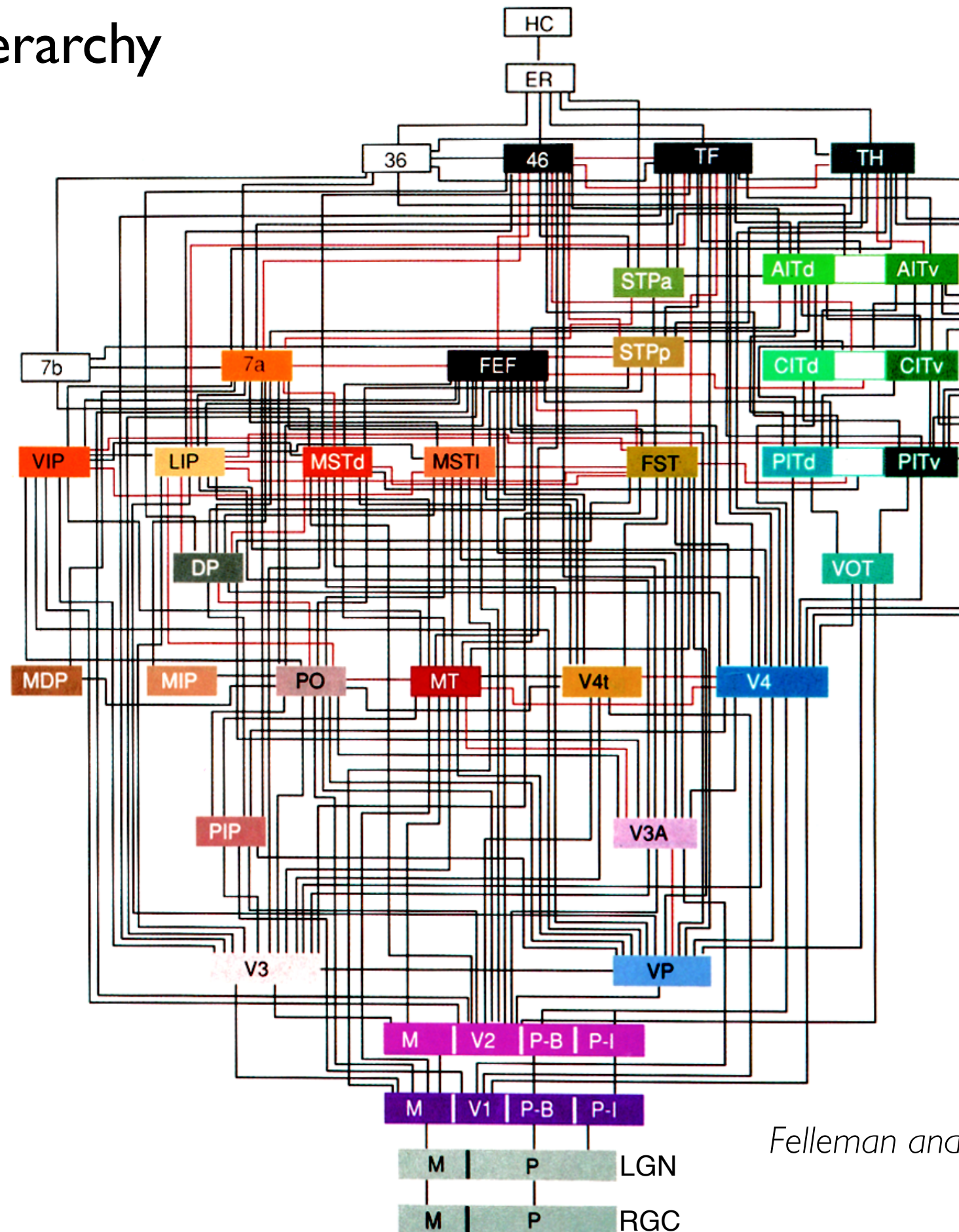
- *Effects of experience* (plasticity)

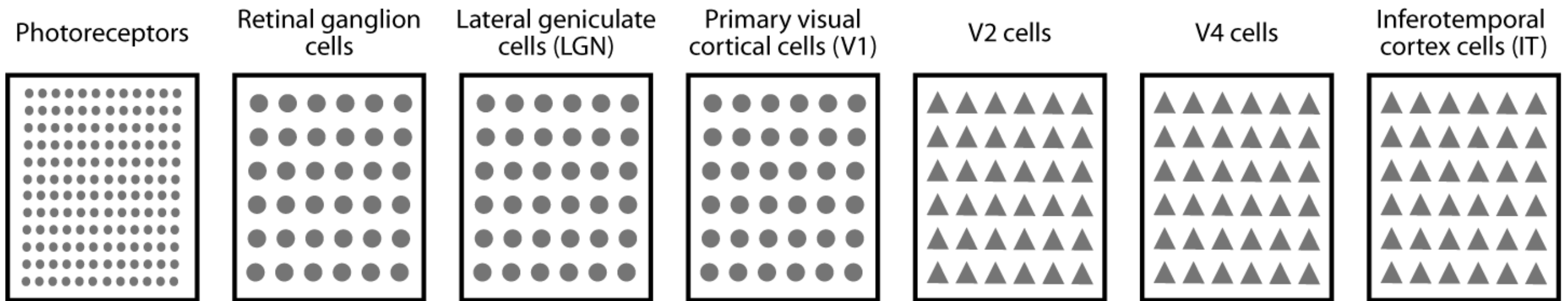sensory cascade in
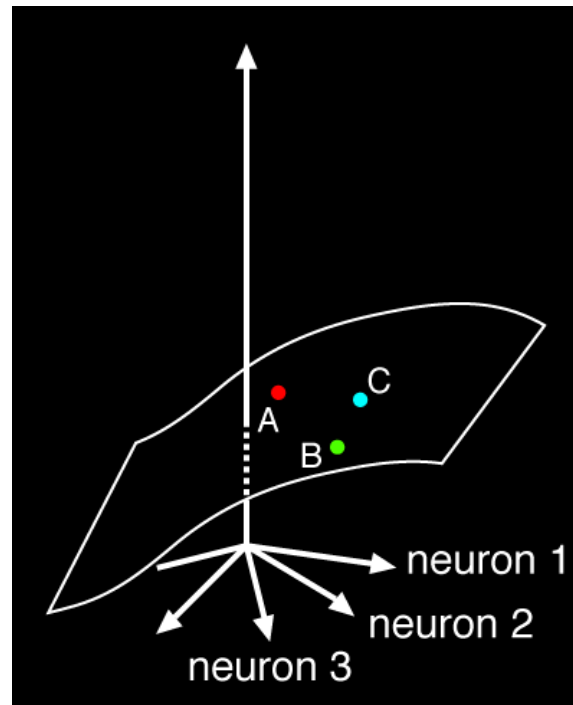visual (mostly-) cortex

Madame Curie!

# Visual area hierarchy



*Felleman and Van Essen, 1991*

# How does the brain represent the visual world?



Photoreceptors | Retinal ganglion cells | Lateral geniculate cells (LGN) | Primary visual cortical cells (V1) | V2 cells | V4 cells | Inferotemporal cortex cells (IT)
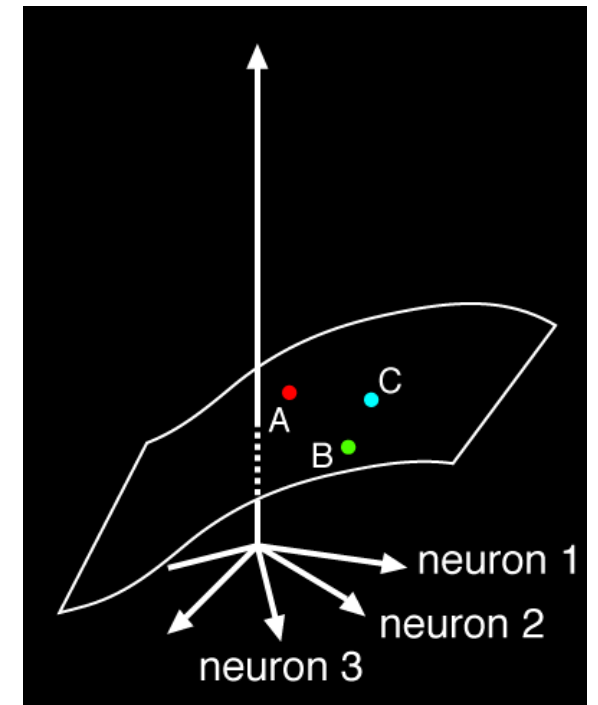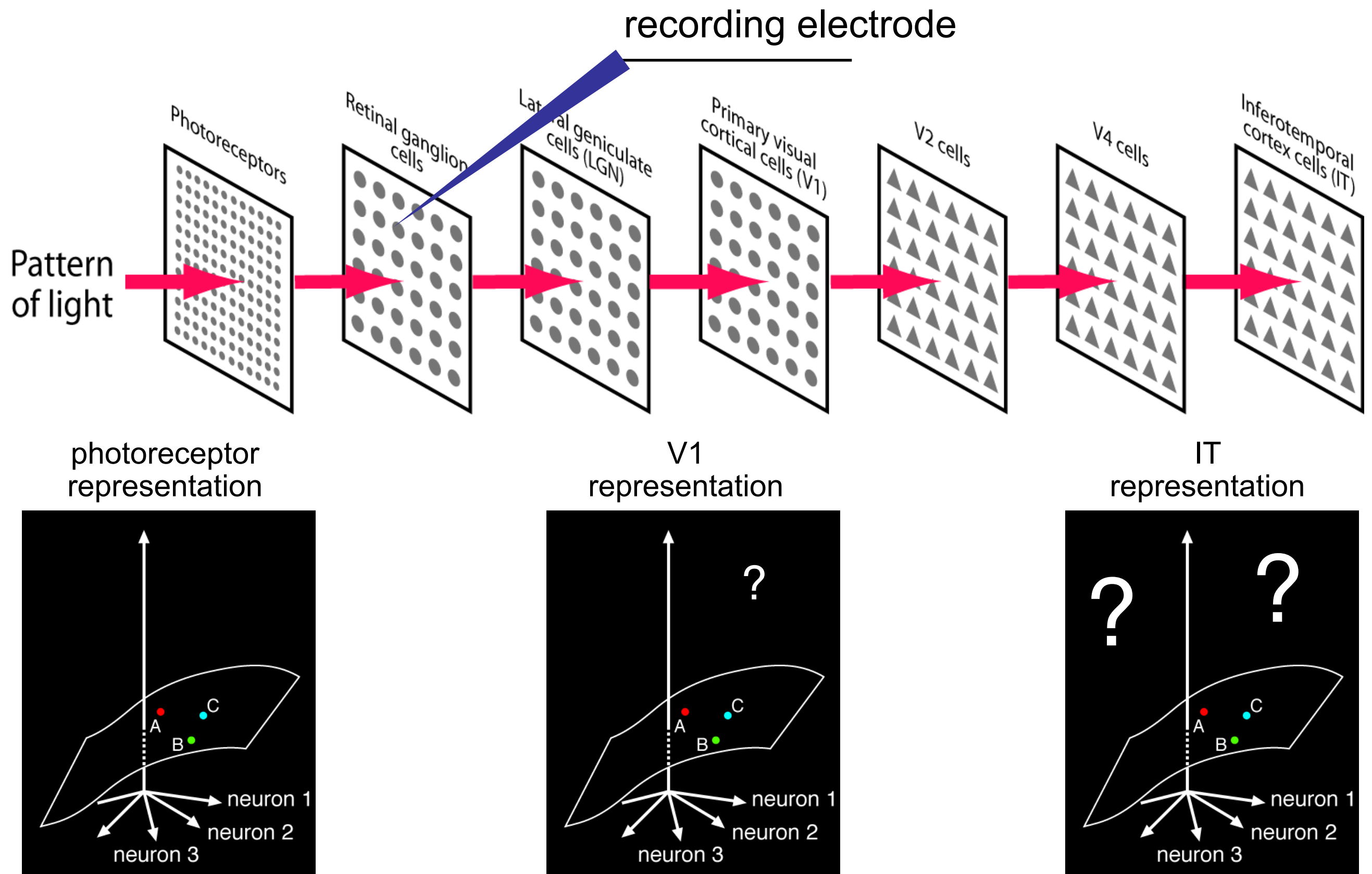
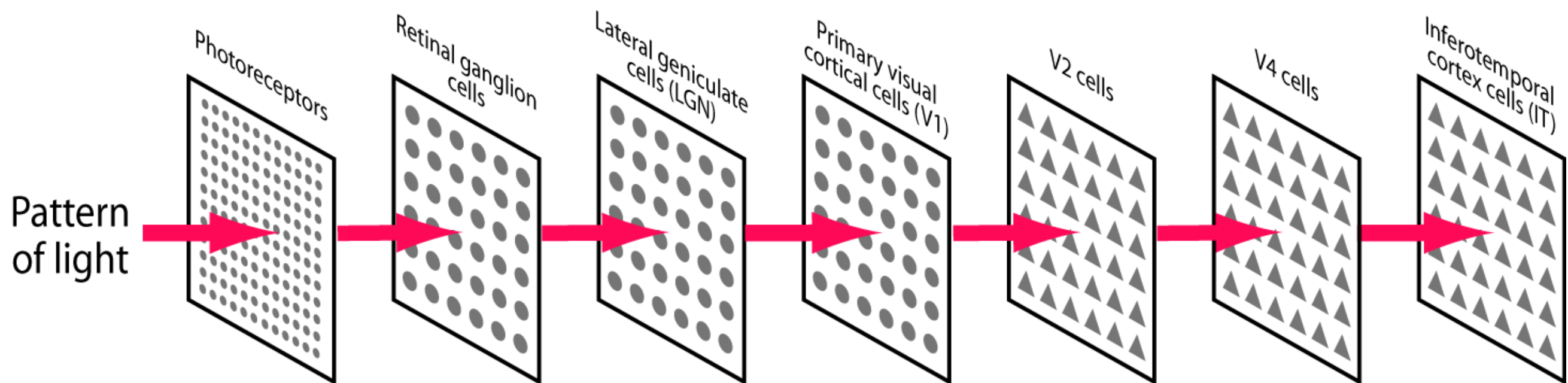photoreceptor representation

V1 representation

IT representation

# How does the brain re-represent the visual world?

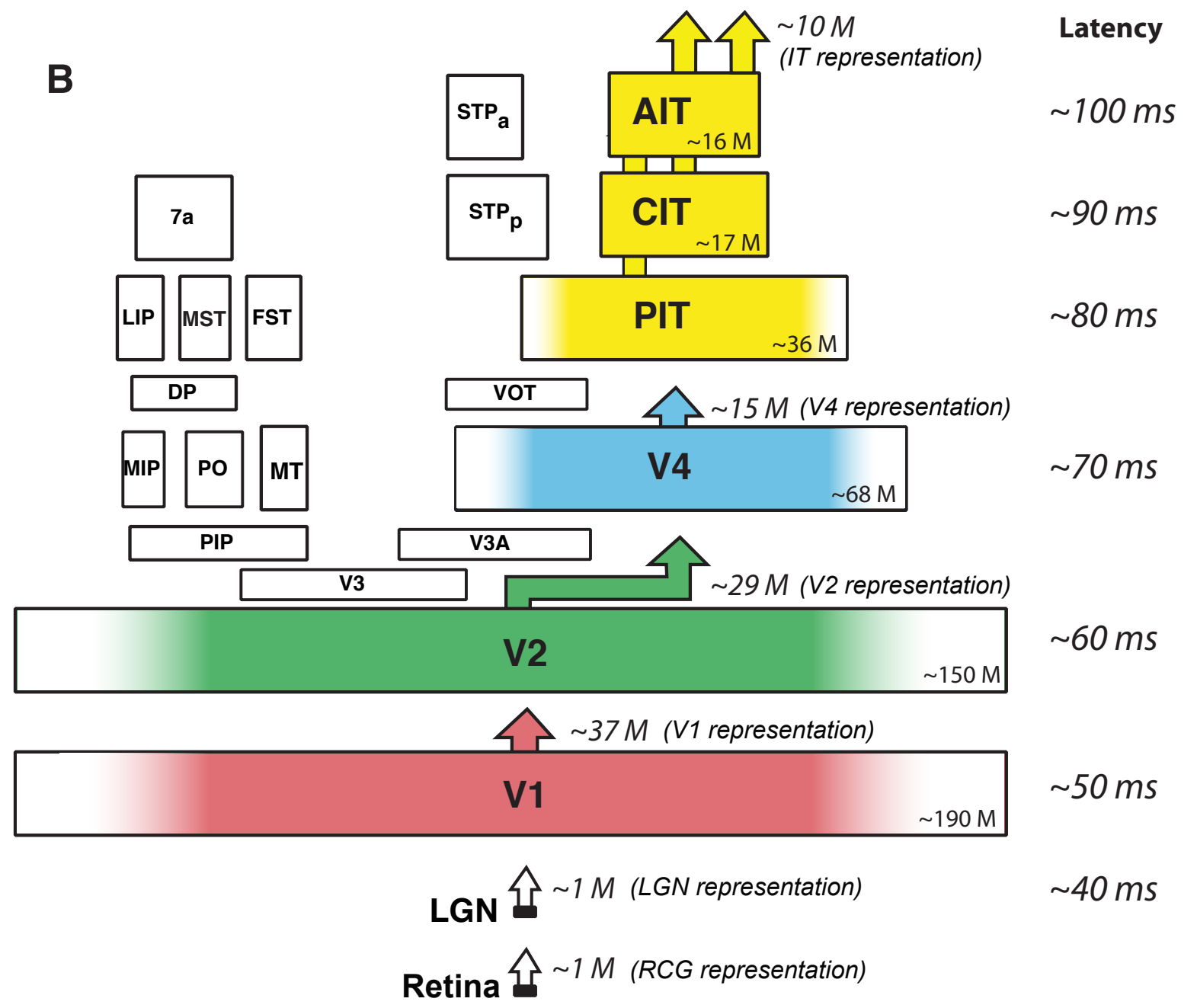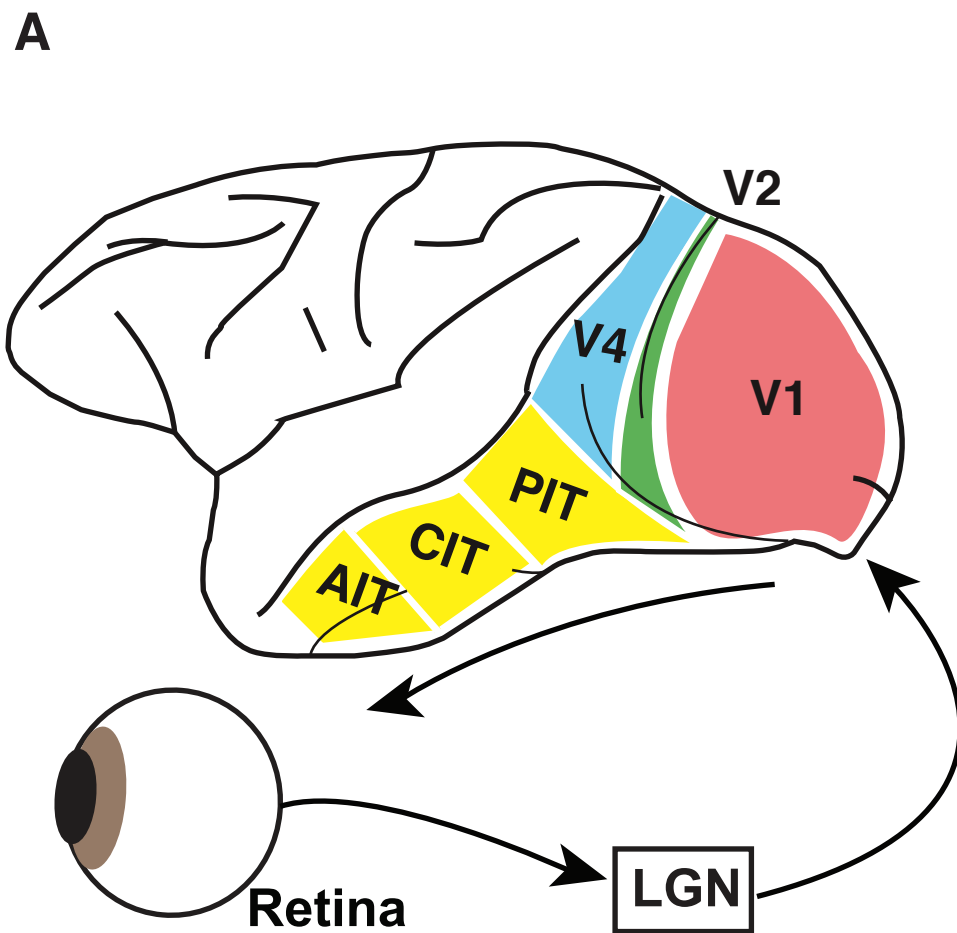# Four important pieces on information

1) Neuronal selectivity generally increases as we move up the cortical hierarchy

2) Receptive field (RF) size generally increases as we move up the cortical hierarchy

3) Selectivity pattern is typically apparent at the time first spikes are elicited by a visual stimulus ("feedforward" assumption)

4) There is hierarchy of times at which first spikes are detected.

**A**

**B**

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)
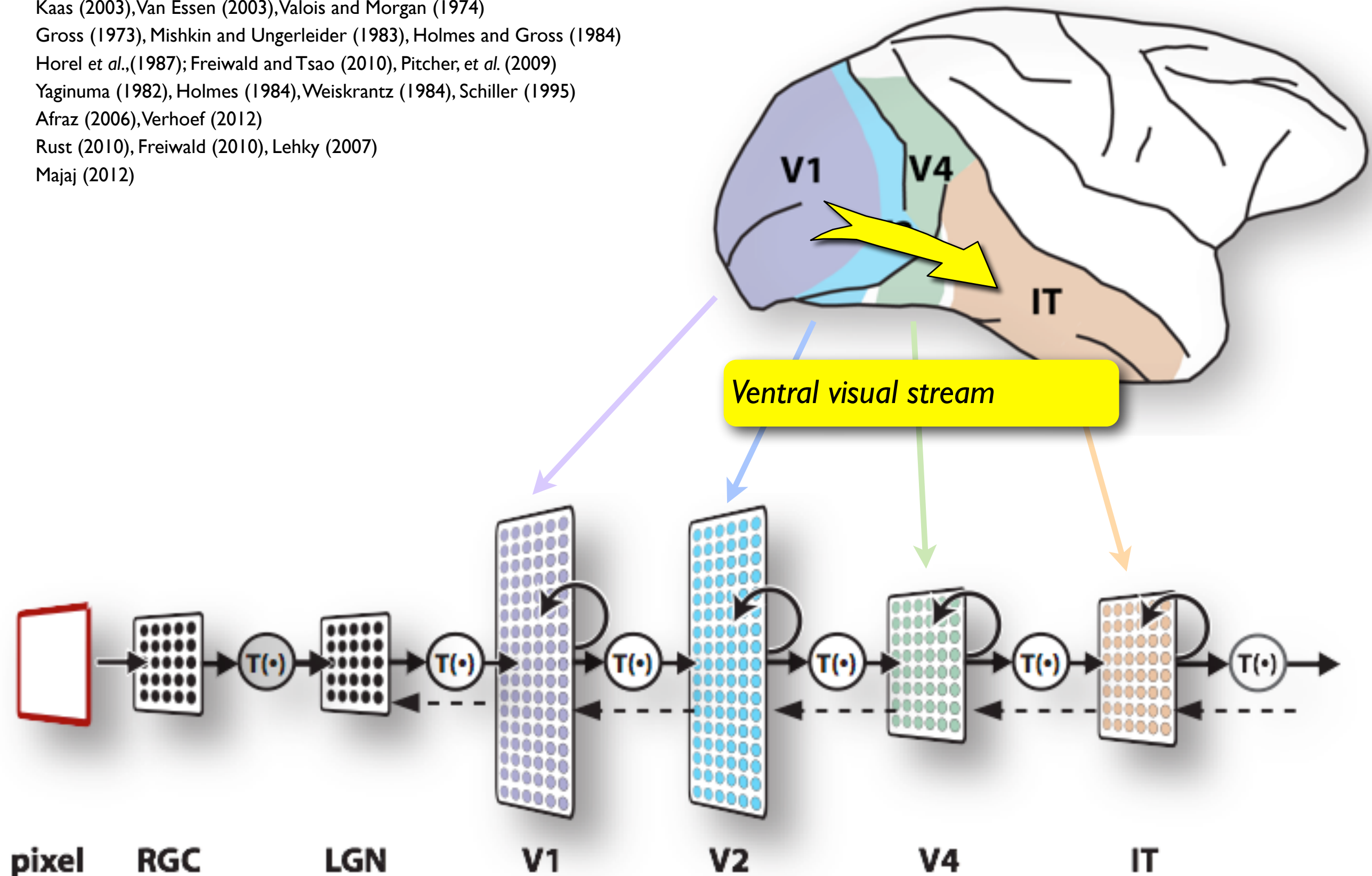
*Adapted from DiCarlo et al. 2012*

# Background: Ventral visual stream

Kaas (2003), Van Essen (2003), Valois and Morgan (1974)
Gross (1973), Mishkin and Ungerleider (1983), Holmes and Gross (1984)
Horel *et al.*,(1987); Freiwald and Tsao (2010), Pitcher, *et al.* (2009)
Yaginuma (1982), Holmes (1984), Weiskrantz (1984), Schiller (1995)
Afraz (2006), Verhoef (2012)
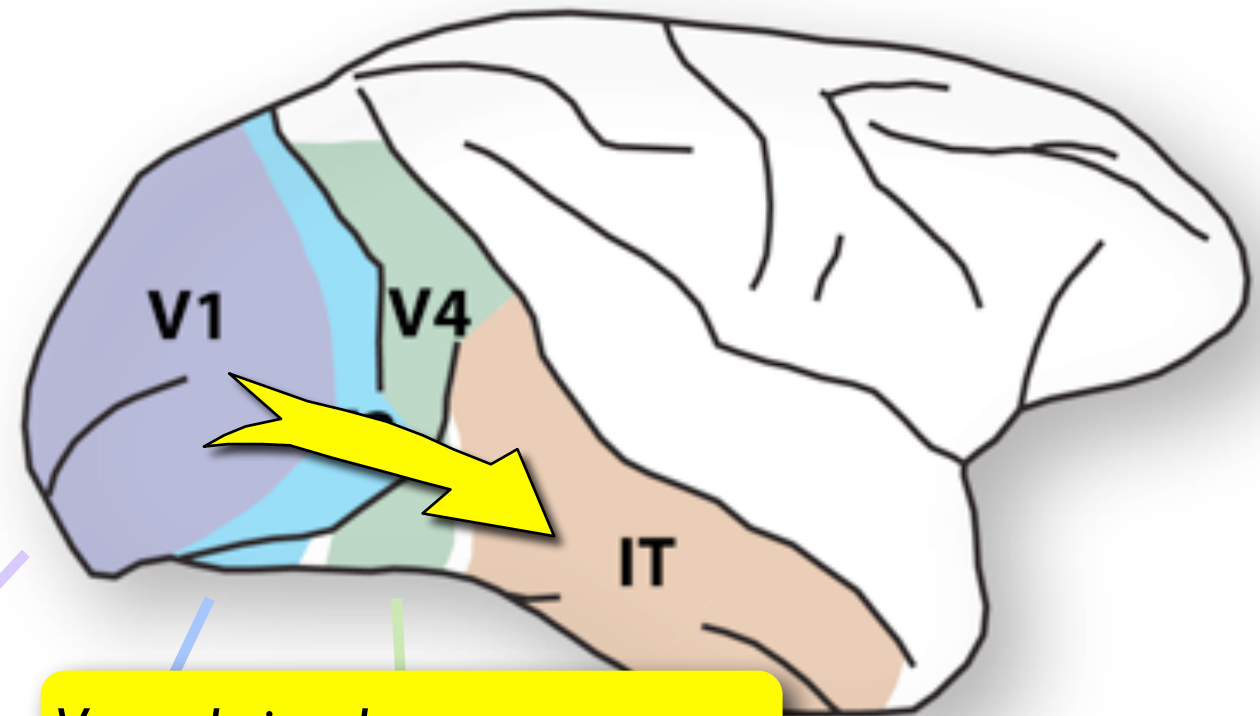Rust (2010), Freiwald (2010), Lehky (2007)
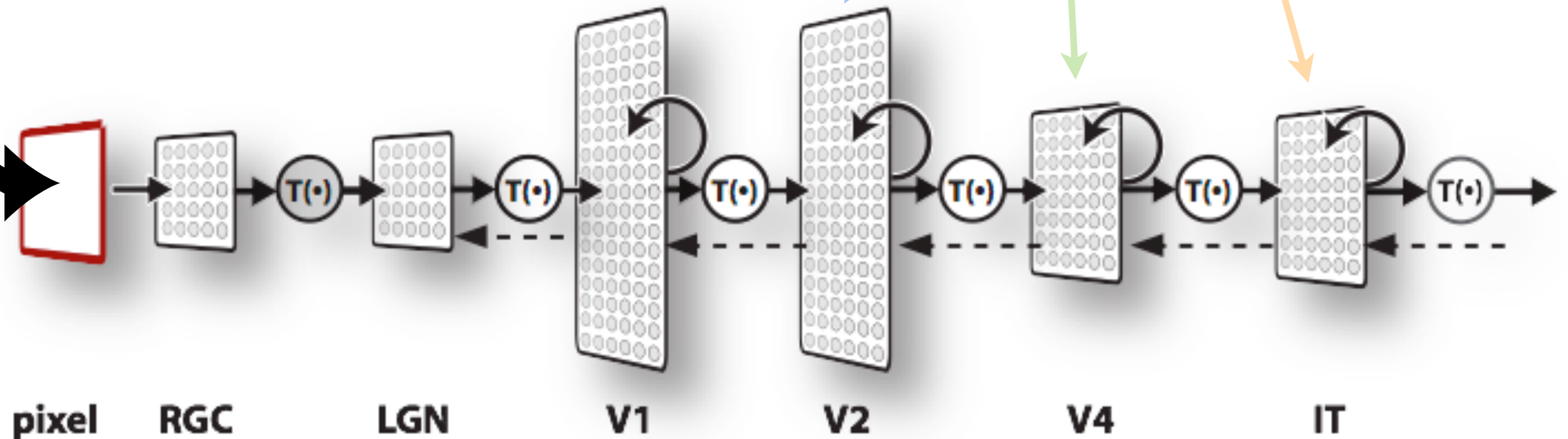Majaj (2012)

*rhesus macaque (macaca mulatta)*

V1  V4  IT

Ventral visual stream

**pixel   RGC   LGN   V1   V2   V4   IT**

*rhesus macaque (macaca mulatta)*

V1

V4

IT

**Ventral visual stream**

pixel    RGC    LGN    V1    V2    V4    IT

*rhesus macaque (macaca mulatta)*

V1  V4  IT

Ventral visual stream

pixel  RGC  LGN  V1  V2  V4  IT

*rhesus macaque (macaca mulatta)*

V1    V4

IT

Ventral visual stream

pixel    RGC    LGN    V1    V2    V4    IT

Stimulus → *representation* → Neurons → *read-out* → Behavior

pixels

100ms Visual Presentation

RGC    LGN    V1    V2    V4    PIT    CIT    AIT

DOG    T(·)    ?    ?    ?

????

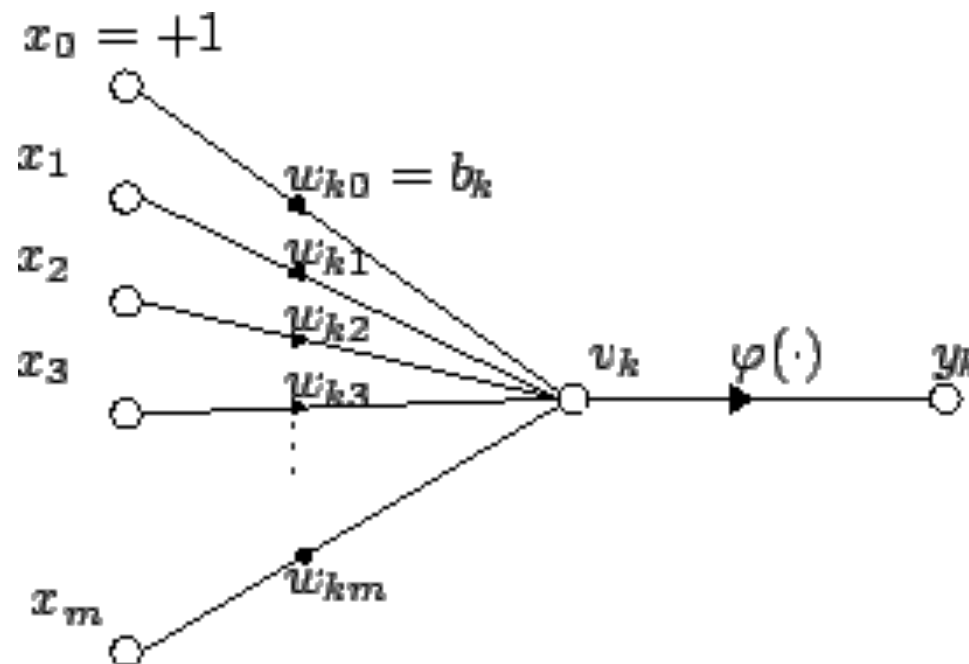Madame Curie!

## McCulloch and Pitts (1943)



$$y_k = \phi \left( \sum_{j=0}^{m} w_{kj} x_j + b_k \right)$$

$$\phi : \mathbb{R} \longmapsto \mathbb{R}$$

some nonlinear activation function
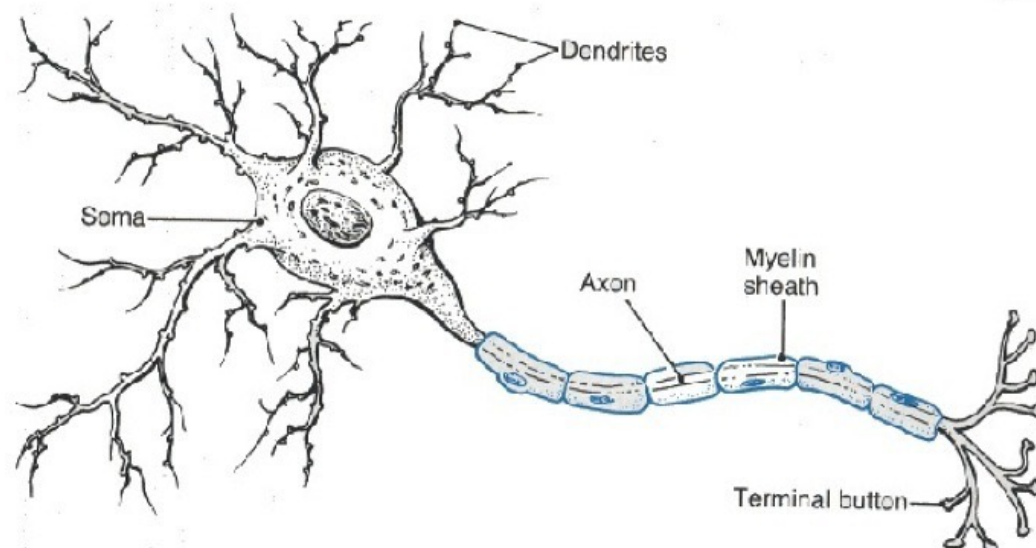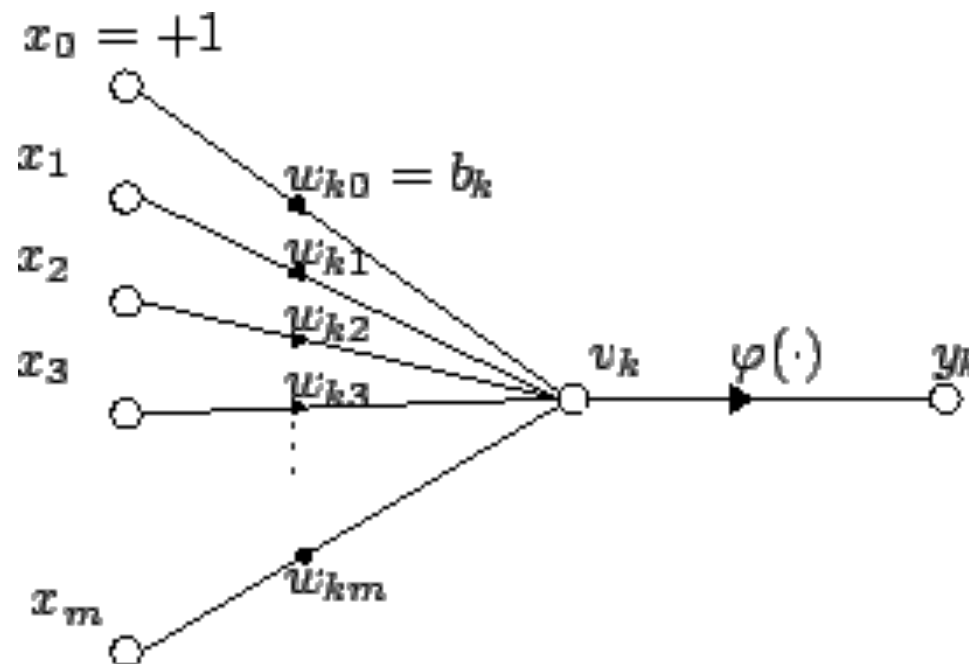


$$w_{kj} \in \mathbb{R}^{m+1}$$

"synaptic strengths"

$$b_j \in \mathbb{R}$$

"biases"

## McCulloch and Pitts (1943)



$$y_k = \phi \left( \sum_{j=0}^{m} w_{kj} x_j + b_k \right)$$

$x_0 = +1$

$x_1$    $w_{k0} = b_k$

$x_2$    $w_{k1}$

$x_3$    $w_{k2}$

     $w_{k3}$

$v_k$   $\varphi(\cdot)$   $y_k$

$x_m$    $w_{km}$

**and what's the connectivity?**

??? 
$$\phi : \mathbb{R} \longmapsto \mathbb{R}$$
some (nonlinear) activation function
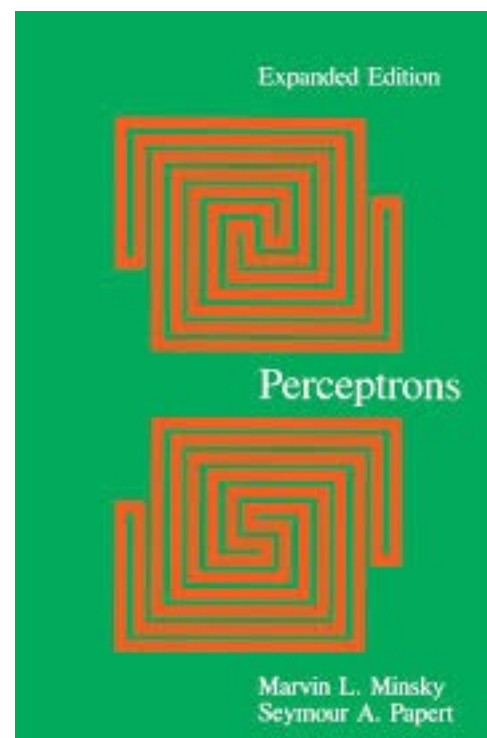
??? 
$$w_{kj} \in \mathbb{R}^{m+1}$$
"synaptic strengths"

$$b_j \in \mathbb{R}$$
"biases"

Minsky & Papert (1969)



$$y_k = \phi \left( \sum_{j=0}^{m} w_{kj} x_j + b_k \right)$$

$$\phi : \mathbb{R} \longmapsto \mathbb{R}$$

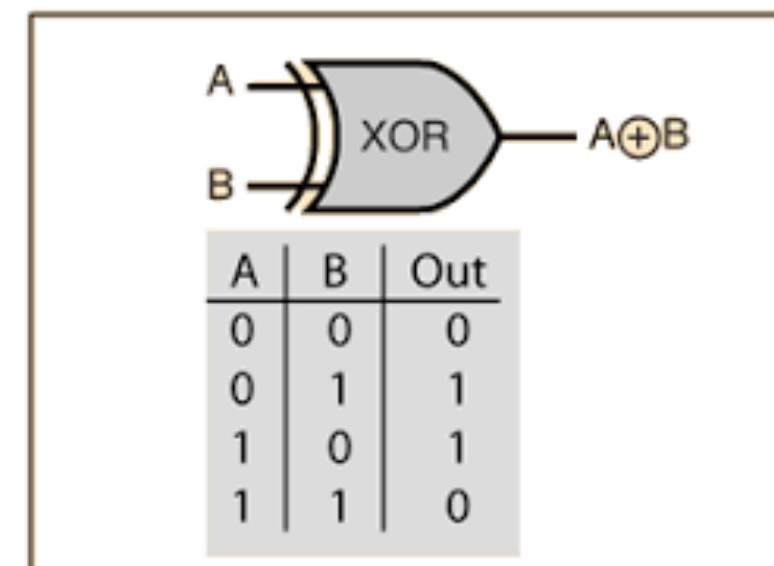**1. better have more than one layer**

**2. better be actually nonlinear**

$$\sum_{i=0}^{N} v_i \phi(w_i^T x + b_i) \quad \text{at least} \ \textit{(and which, according to the UAT, is enough)}$$

**cause otherwise ... ain't no XOR**

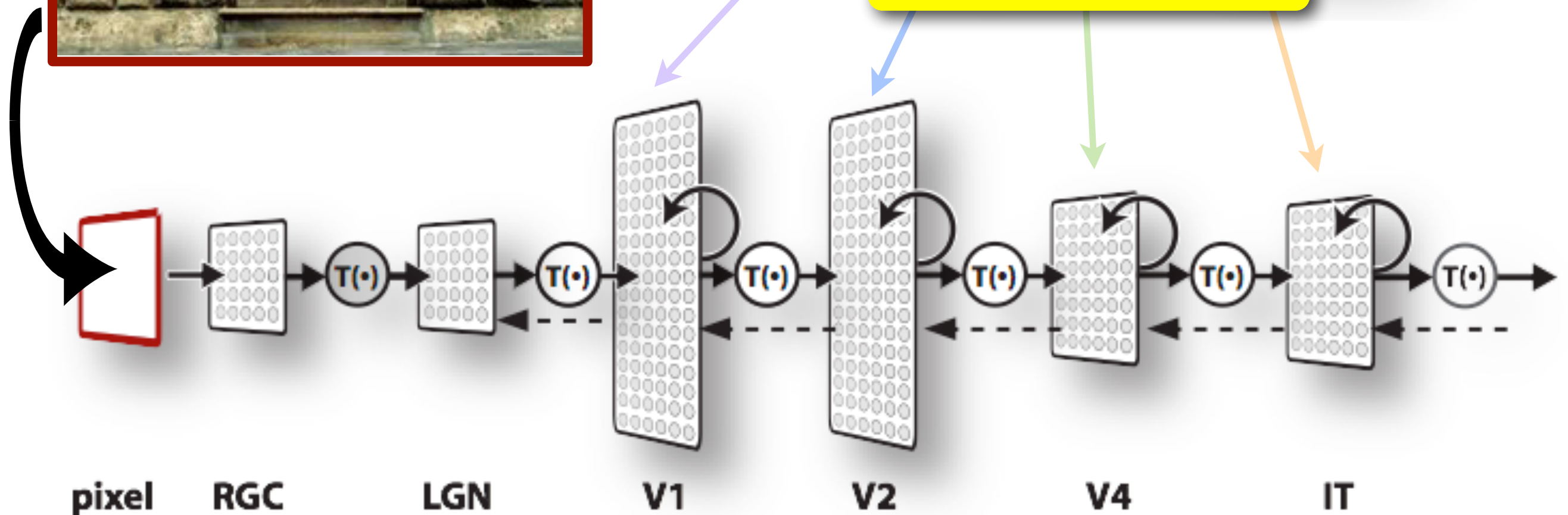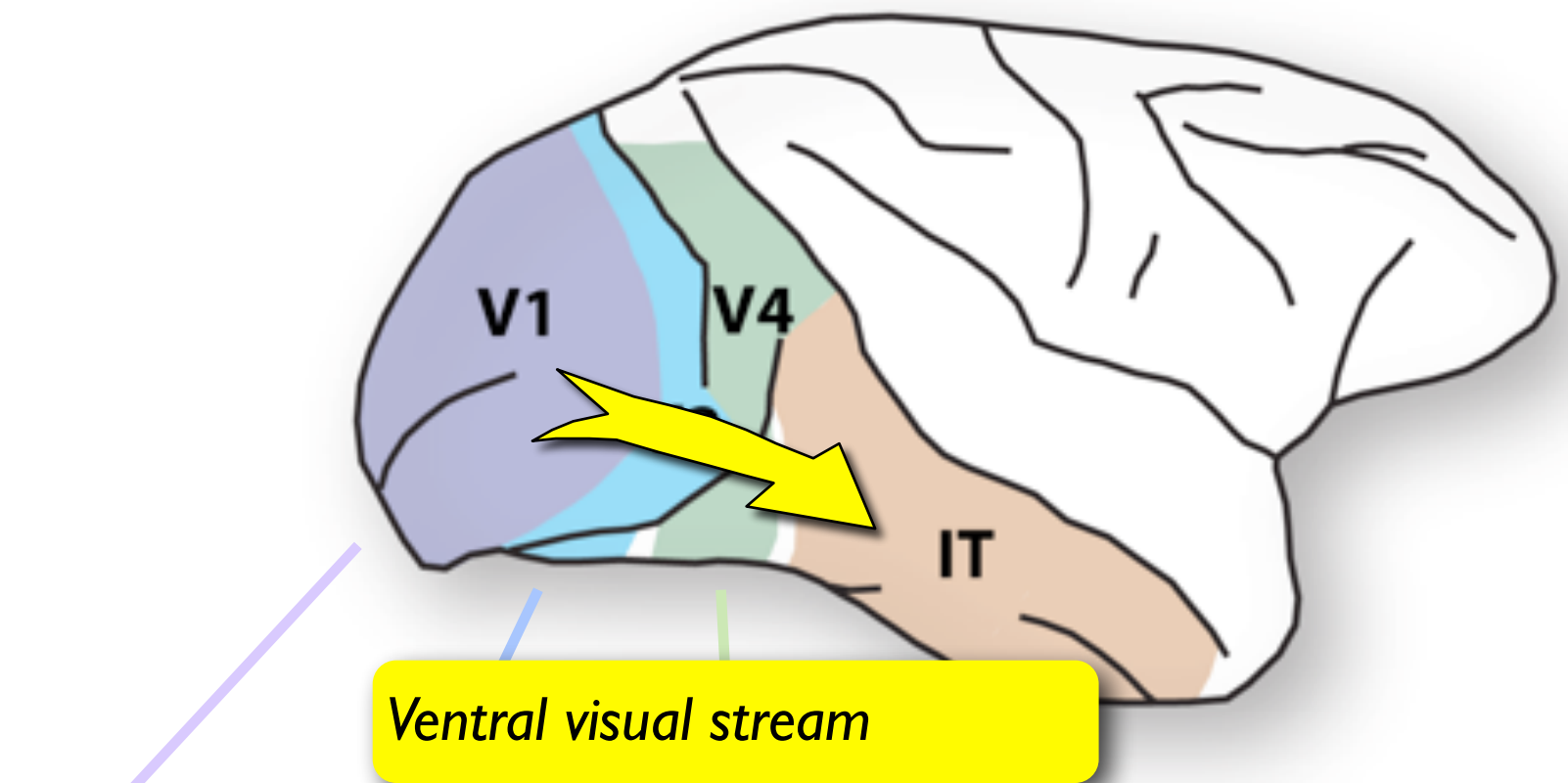**and what's the connectivity?**

**Limitations of Perceptrons**

- *Minsky & Papert published (1969) "Perceptrons" stressing the limitations of perceptrons*

- *Single-layer perceptrons cannot solve problems that are linearly inseparable (e.g., xor)*

- *Most interesting problems are linearly inseparable*

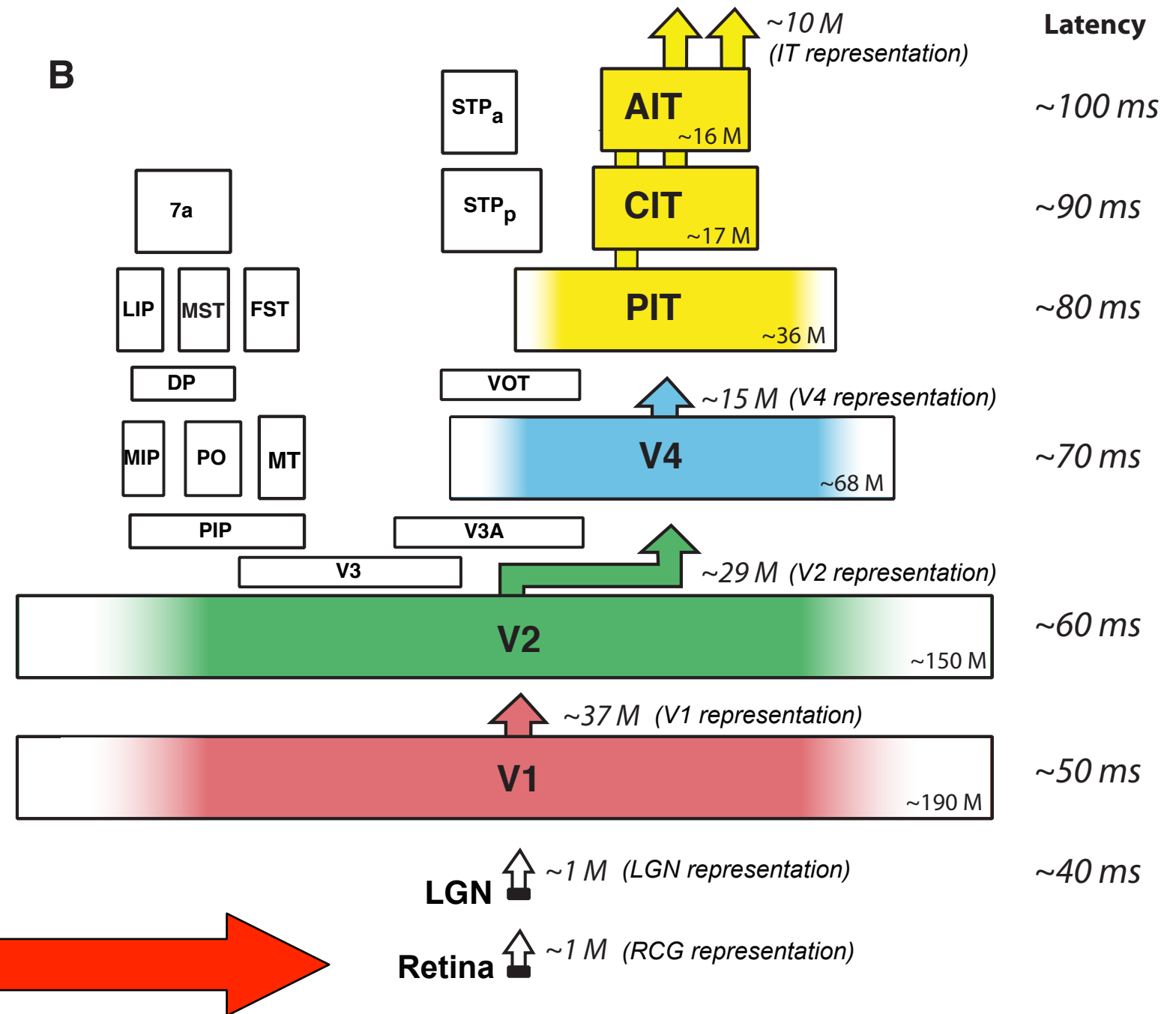- *Kills funding for neural nets for 12-15 years*

UNIVERSITY of VIRGINIA

Maybe a bit apocryphal …. but I can definitely say from personal experience that MIT CSAIL felt very "anti-neural networks" as late as 2012

# Ventral Stream = Connected series of brain areas
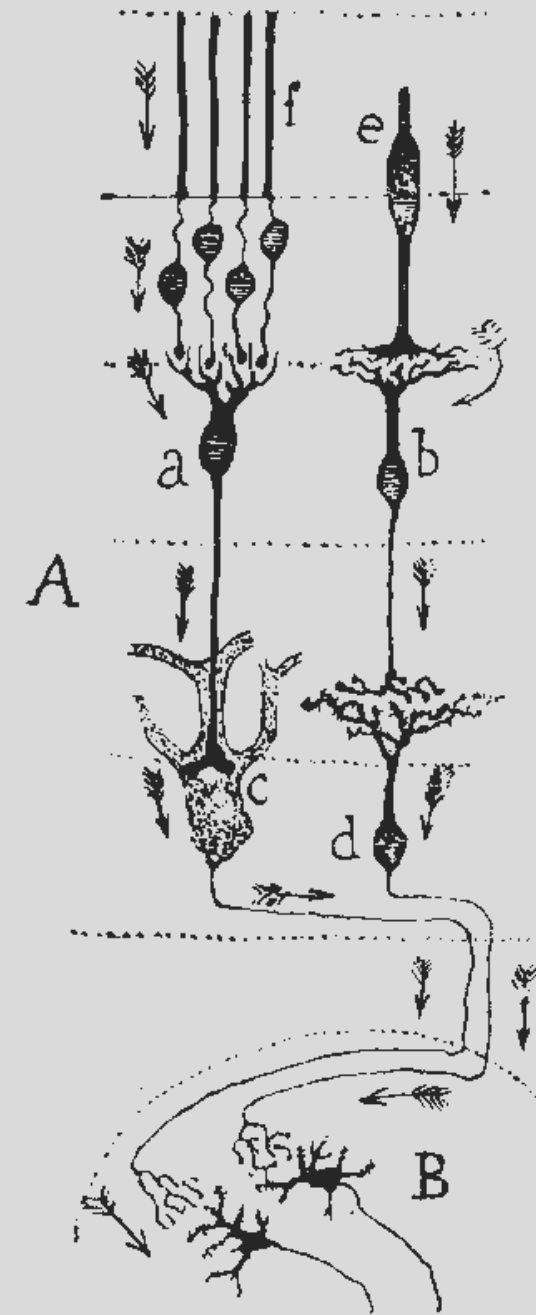
neuroanatomy +  neurophysiology tell us:



*Ventral visual stream*

pixel    RGC    LGN    V1    V2    V4    IT

**B**

~10 M
(IT representation)

Latency

~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)

**You are here.**

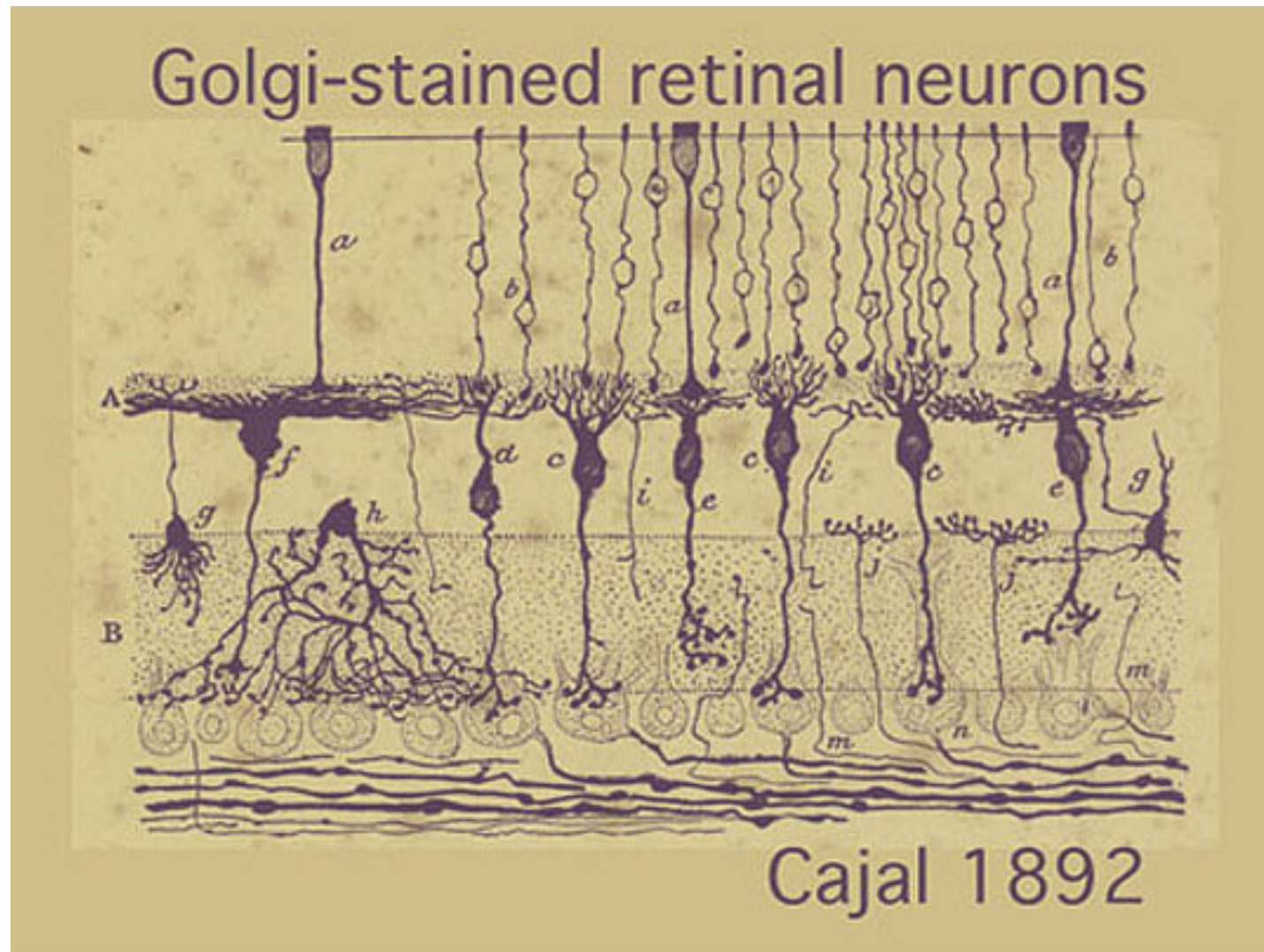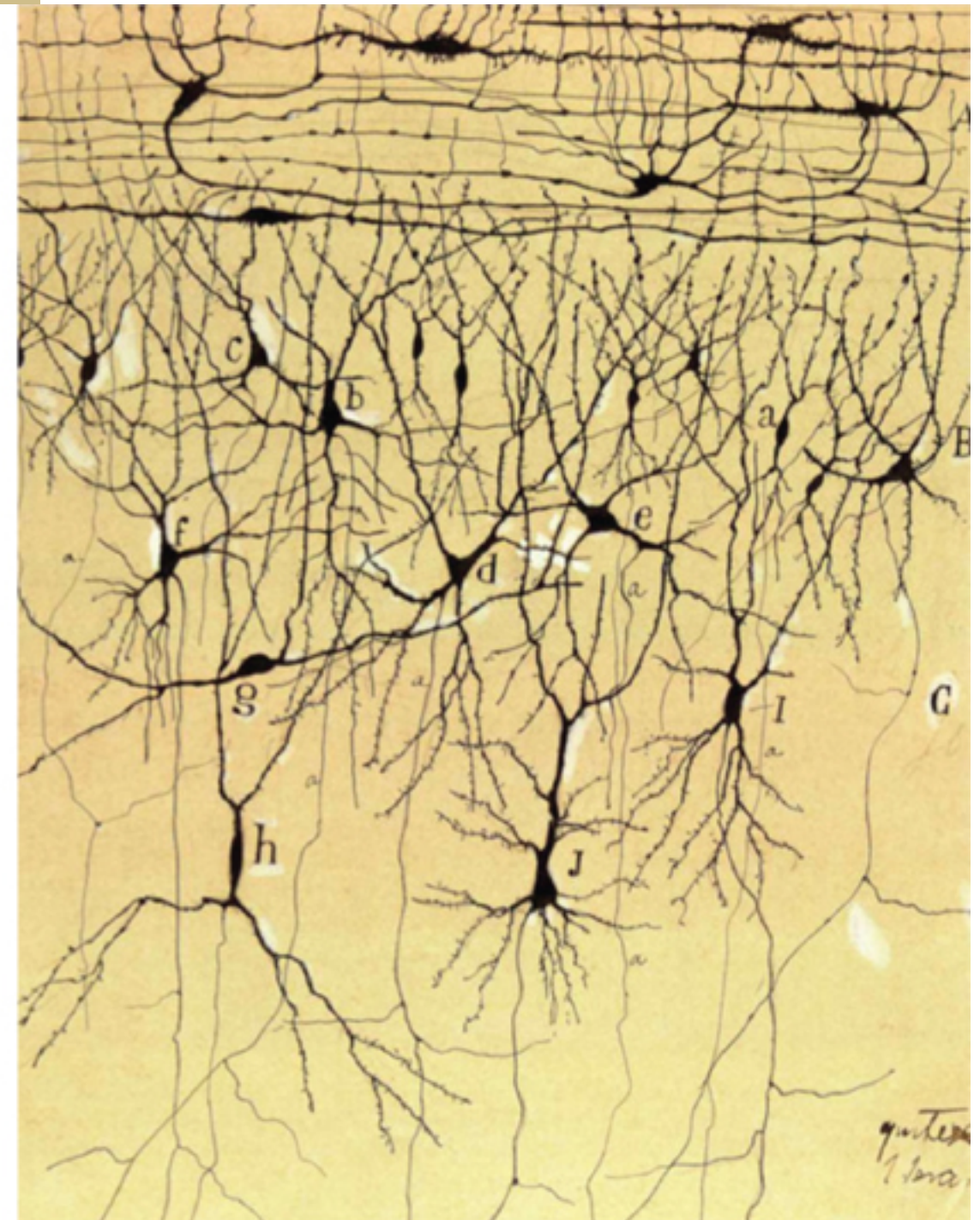*Adapted from DiCarlo et al. 2012*

**Ramon y Cajal** *from Rodieck (1973)*

Fig. 2. A drawing done by Cajal to show some of the neurons of the retina in vertical section.

Golgi-stained retinal neurons
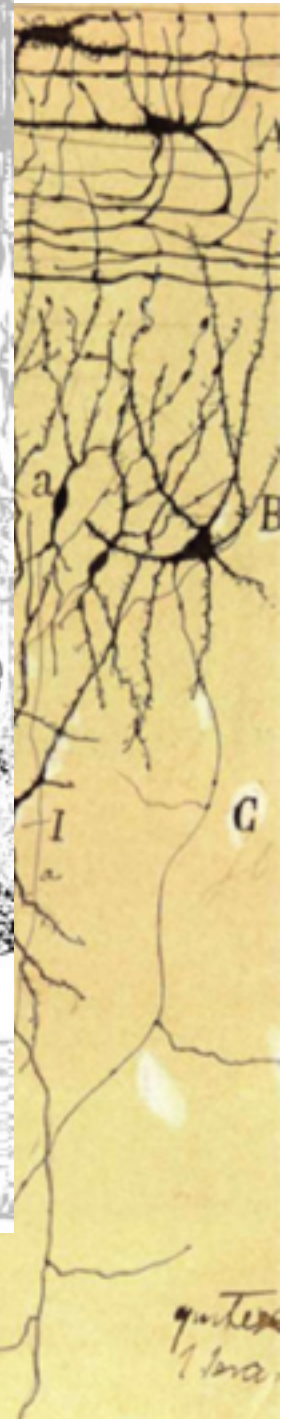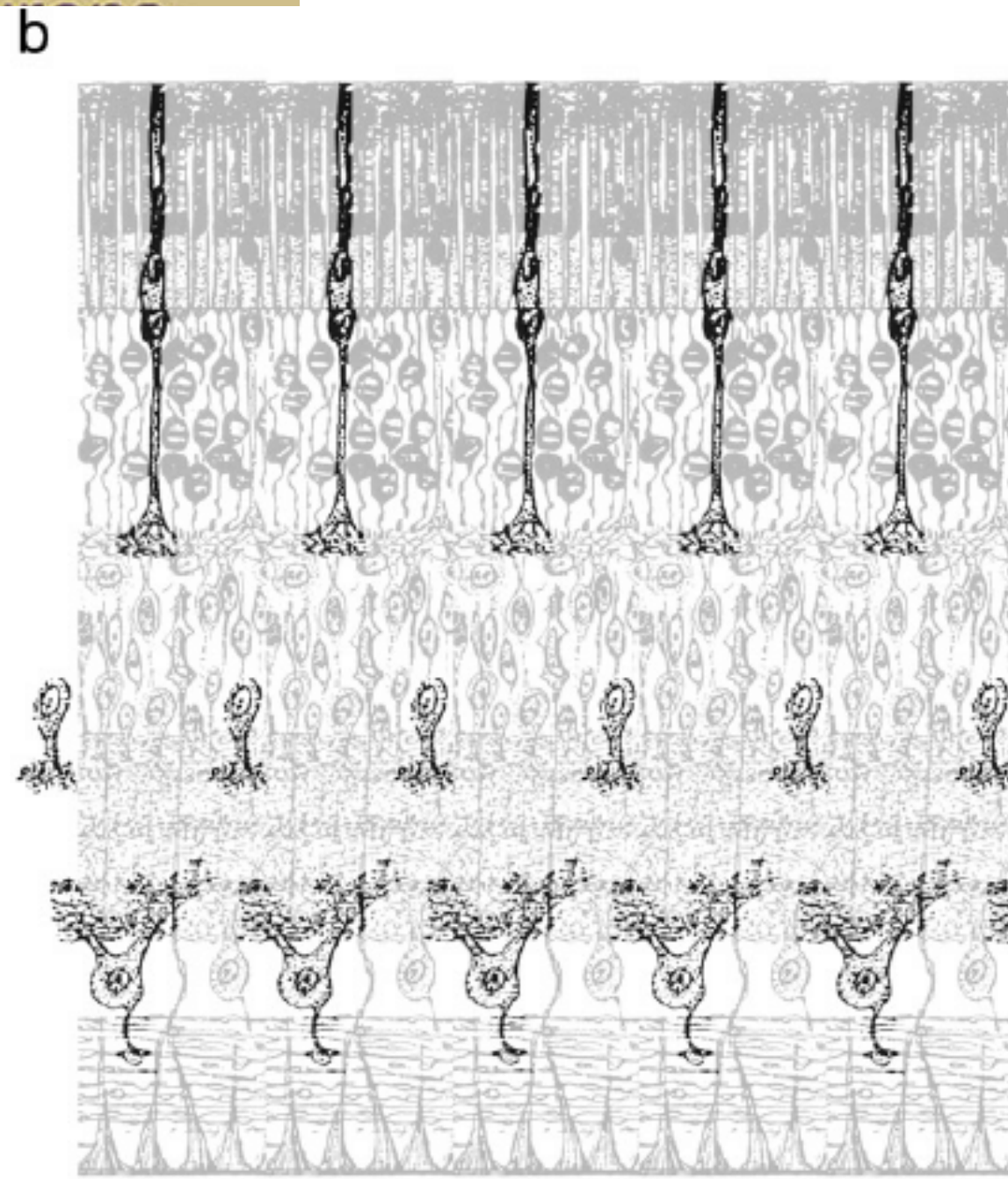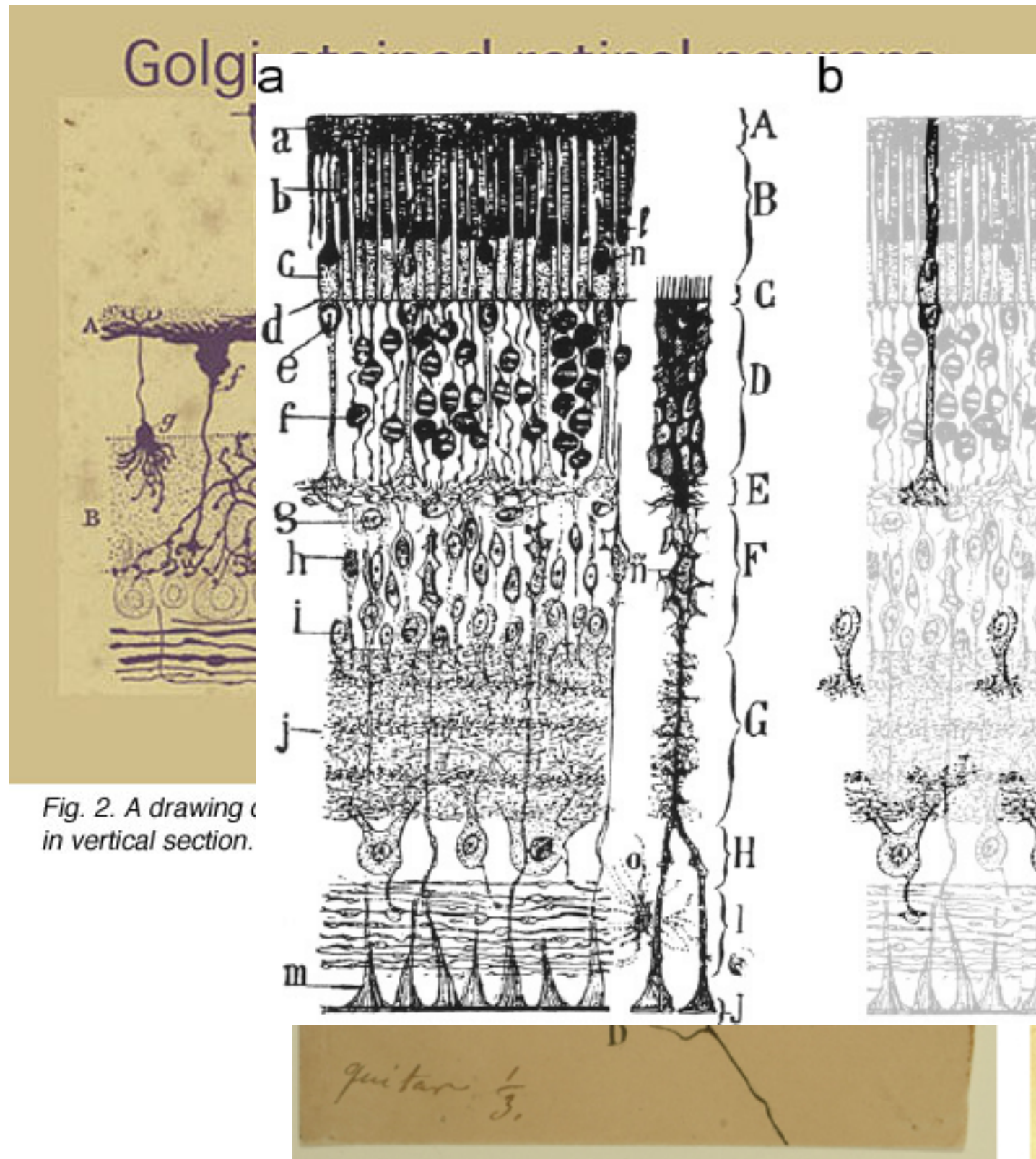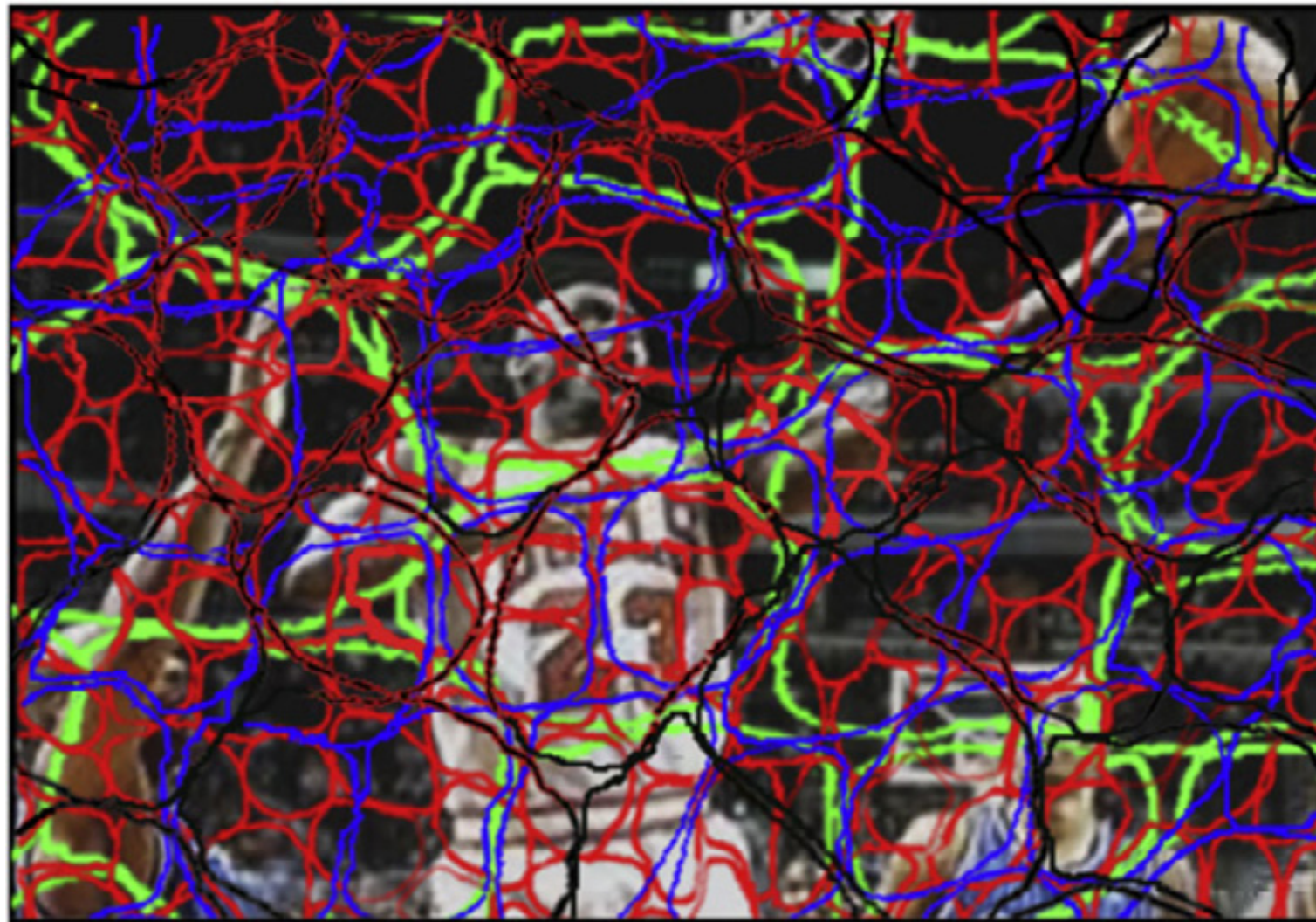
Fig. 2. A drawing don[...]
in vertical section.

Golgi stained retinal neurons

Fig. 2. A drawing in vertical section.

(a)   (b)   (c)   (d)

*Masland (2012)*

http://museum.eyewire.org

## SPATIAL CONVOLUTION



(a)     (b)     (c)     (d)

**Input Digital Image**

| 9 | 5 | 7 |
| 4 | 6 | 2 |
| 6 | 7 | 9 |

Target Pixel

**Convolution Mask**

| 0 | -1 | 0 |
| -1 | 5 | -1 |
| 0 | -1 | 0 |

**Mask Overlay on Target Pixel**

| 0 x 9 | -1 x 5 | 0 x 7 |
| -1 x 4 | 5 x 6 | -1 x 2 |
| 0 x 6 | -1 x 7 | 0 x 9 |

**Multiply Mask Times Pixel Intensities**

| 0 | -5 | 0 |
| -4 | 30 | -2 |
| 0 | -7 | 0 |

```
-5
-2
-7
-4
+30
____
12
```

| | | |
| 12 | | |
| | | |

**Output Grey Value of Target Pixel is Sum of Products**

1. The convolution mask is overlated on the original image so that the center pixel of the mask is matched with a pixel location on the image (Target Pixel- to be convolved).

2. Each pixel value in the original image is multiplied by the corresponding value in the overlying mask..

3. The grey value of the target pixel is replaced by the sum of all the products in the second step.

4. The operation is repeated for each pixel in the original image (the mask scans the entire image) and each pixel is replaced by the weighted average of its 3 x 3 neighbors.
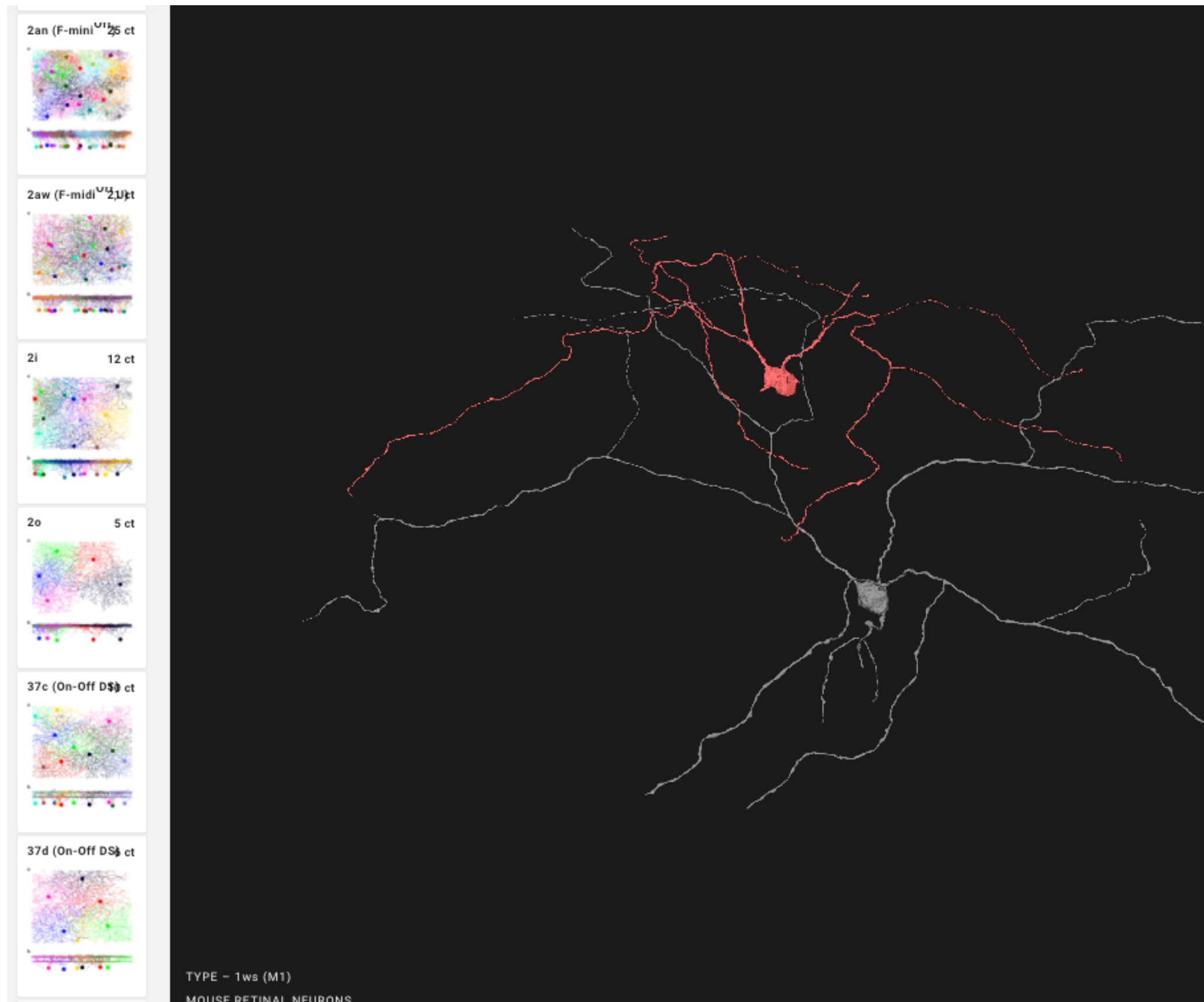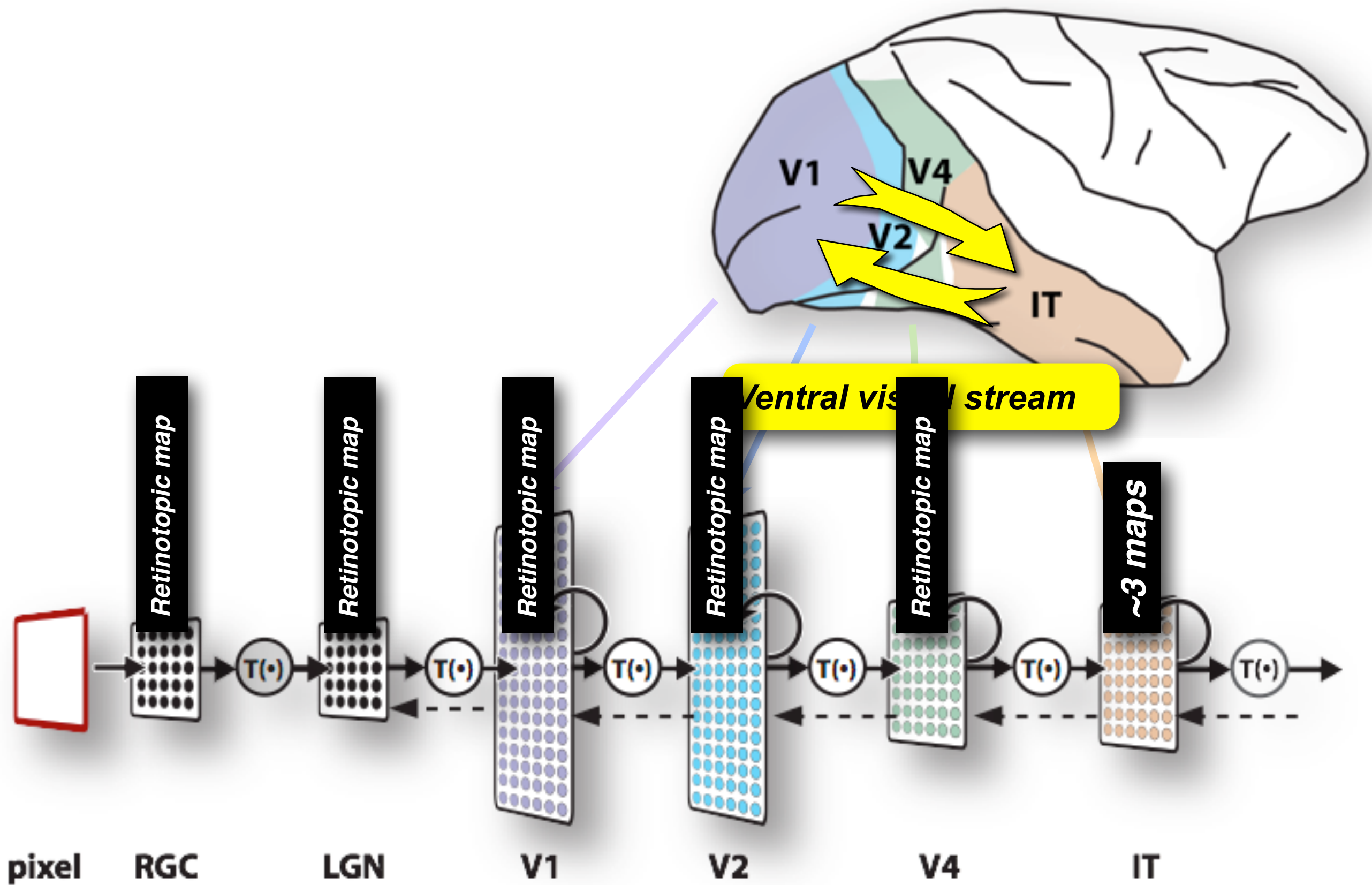
cell types
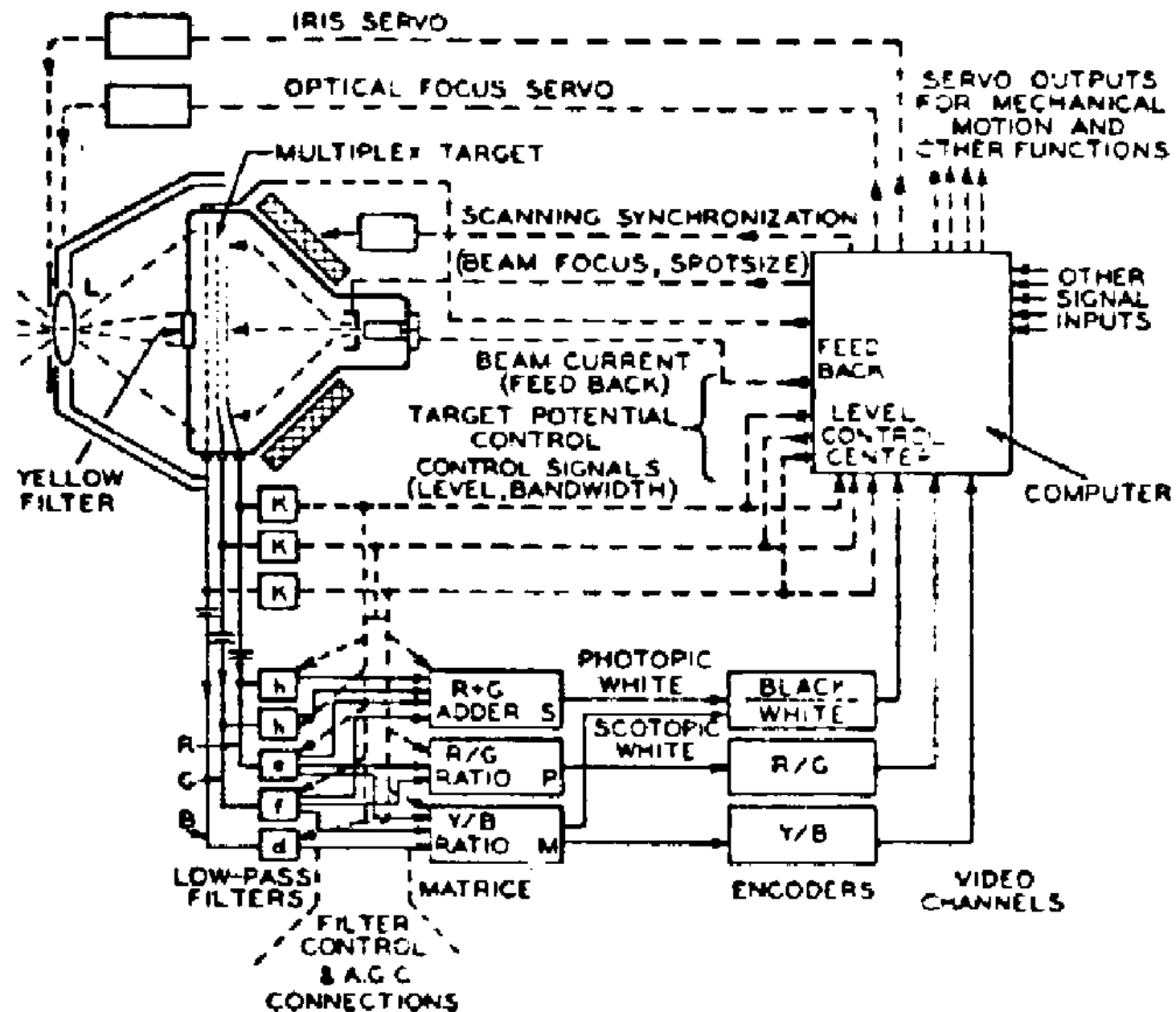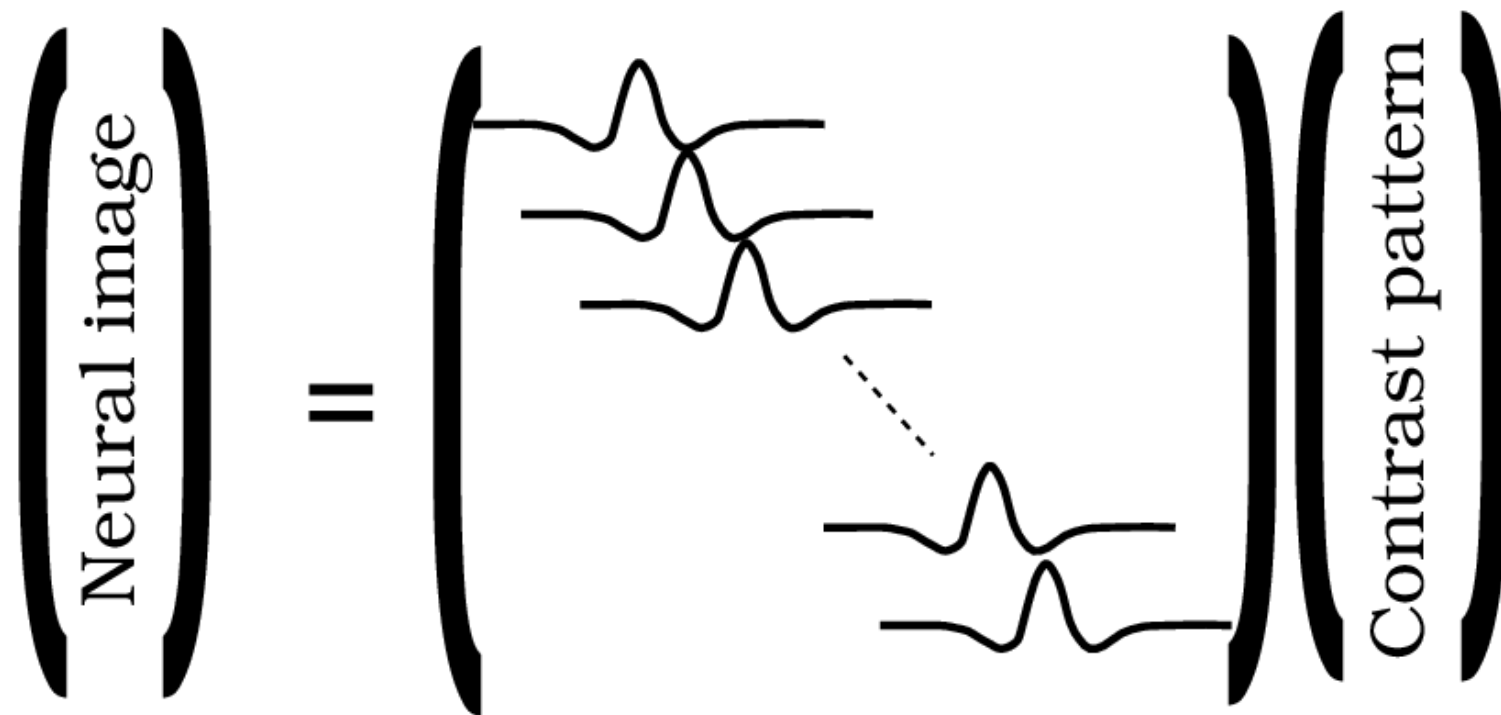like different
filters in a
filterbanks
.
.
.
but which
filters?

V1

V4

V2

IT

Ventral visual stream

Retinotopic map

Retinotopic map

Retinotopic map

Retinotopic map

Retinotopic map

~3 maps

T(•)    T(•)    T(•)    T(•)    T(•)    T(•)    T(•)

pixel    RGC    LGN    V1    V2    V4    IT

characterizing a *transfer function . . .*



*Schade, 1956 from Wandell 1996*

characterizing a *transfer function …*
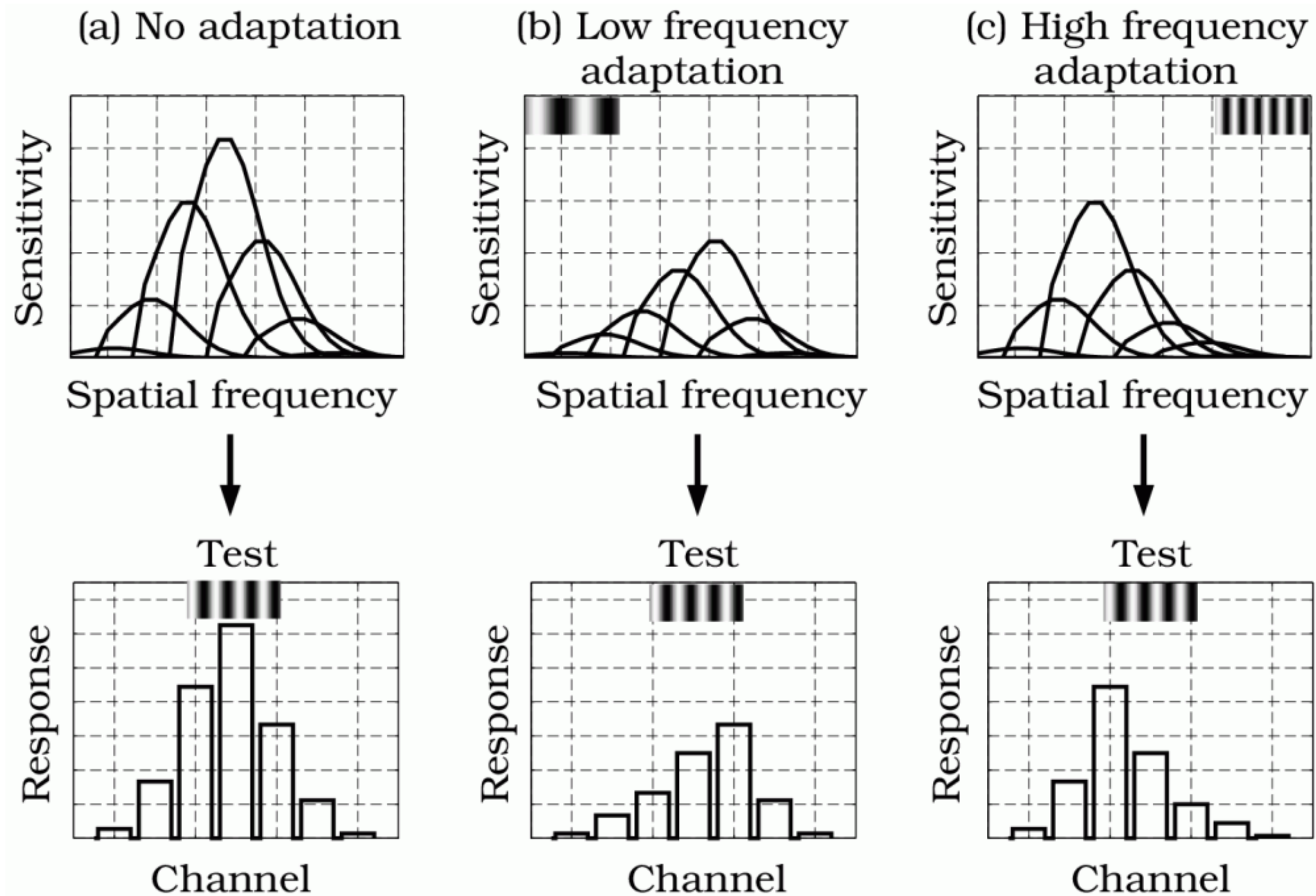


*Christina Enroth-Cugell*

*John Robson*

… and thus, presumably, doing linear systems (e.g fourier) analysis

# Origins in the Retina



(a) No adaptation

(b) Low frequency adaptation

(c) High frequency adaptation

*from Wandell 1996*

# Origins in the Retina

## THE CONTRAST SENSITIVITY OF RETINAL GANGLION CELLS OF THE CAT

By CHRISTINA ENROTH-CUGELL and J. G. ROBSON*

*From the Biomedical Engineering Center, Technological Institute, Northwestern University, Evanston, Illinois, U.S.A.† and the Department of Physiology, Northwestern University Medical School, Chicago, U.S.A.*

(*Received 19 April 1966*)

*Christina Enroth-Cugell*

*John Robson*

1. Spatial summation within cat retinal receptive fields was studied by recording … responses of ganglion cells to grating patterns

2. Summation over the receptive fields of some cells (X-cells) was found to be **approximately linear**, while for other cells (Y-cells) summation was **very non-linear**.

3. The mean discharge frequency of Y-cells … was greatly increased when grating patterns drifted across their receptive fields.

4. In X-cells …  it was found that the contrast sensitivity function, **could be satisfactorily described by the difference of two Gaussian functions**.

5. This finding supports the hypothesis that the sensitivities of the antagonistic centre and surround summating regions of ganglion cell receptive fields fall off as Gaussian functions of the distance from the field centre.
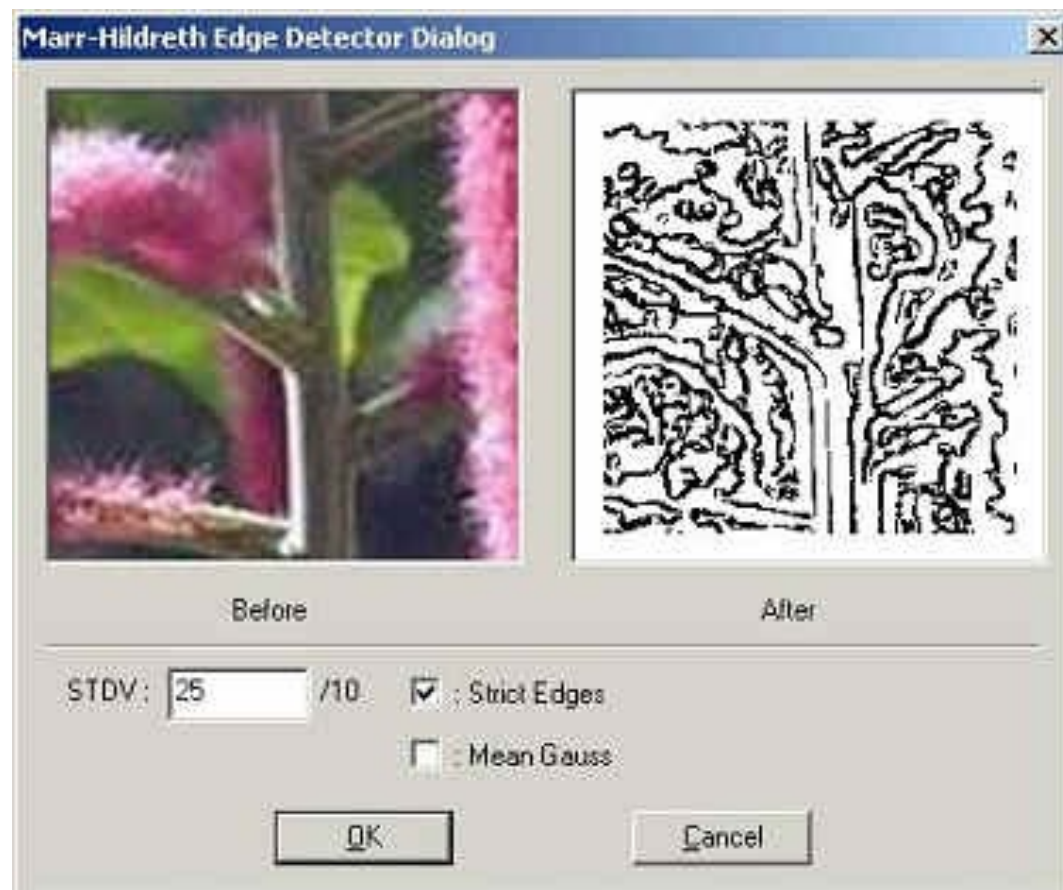
## Theory of edge detection

By D. Marr and E. Hildreth

*M.I.T. Psychology Department and Artificial Intelligence Laboratory,
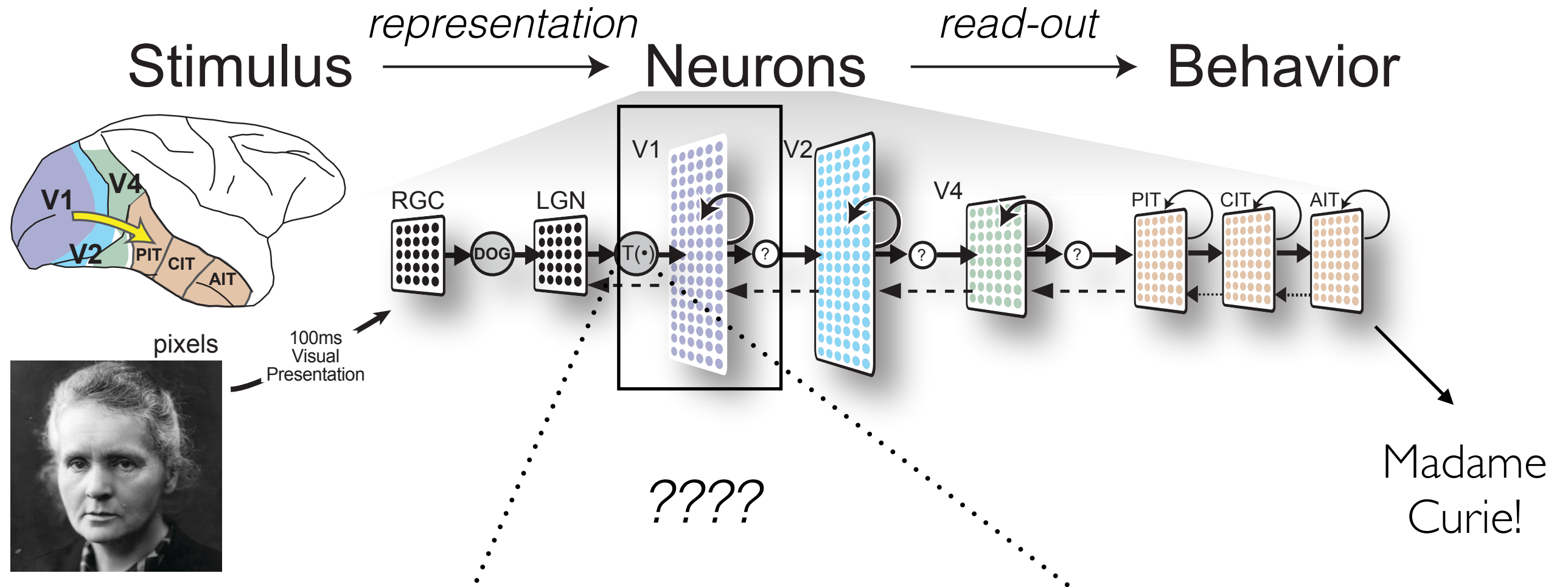79 Amherst Street, Cambridge, Massachusetts 02139, U.S.A.*

$$\vec{\nabla}^2 G(x,y) * Im(x,y)$$

$$\sim DoG$$

**Marr-Hildreth Edge Detector Dialog**

Before    After

STDV: 25    /10    ☑ : Strict Edges
                   ☐ : Mean Gauss

OK    Cancel

Adapted from DiCarlo et al. 2012

# RECEPTIVE FIELDS, BINOCULAR INTERACTION AND FUNCTIONAL ARCHITECTURE IN THE CAT'S VISUAL CORTEX

BY D. H. HUBEL AND T. N. WIESEL

*From the Neurophysiology Laboratory, Department of Pharmacology Harvard Medical School, Boston, Massachusetts, U.S.A.*
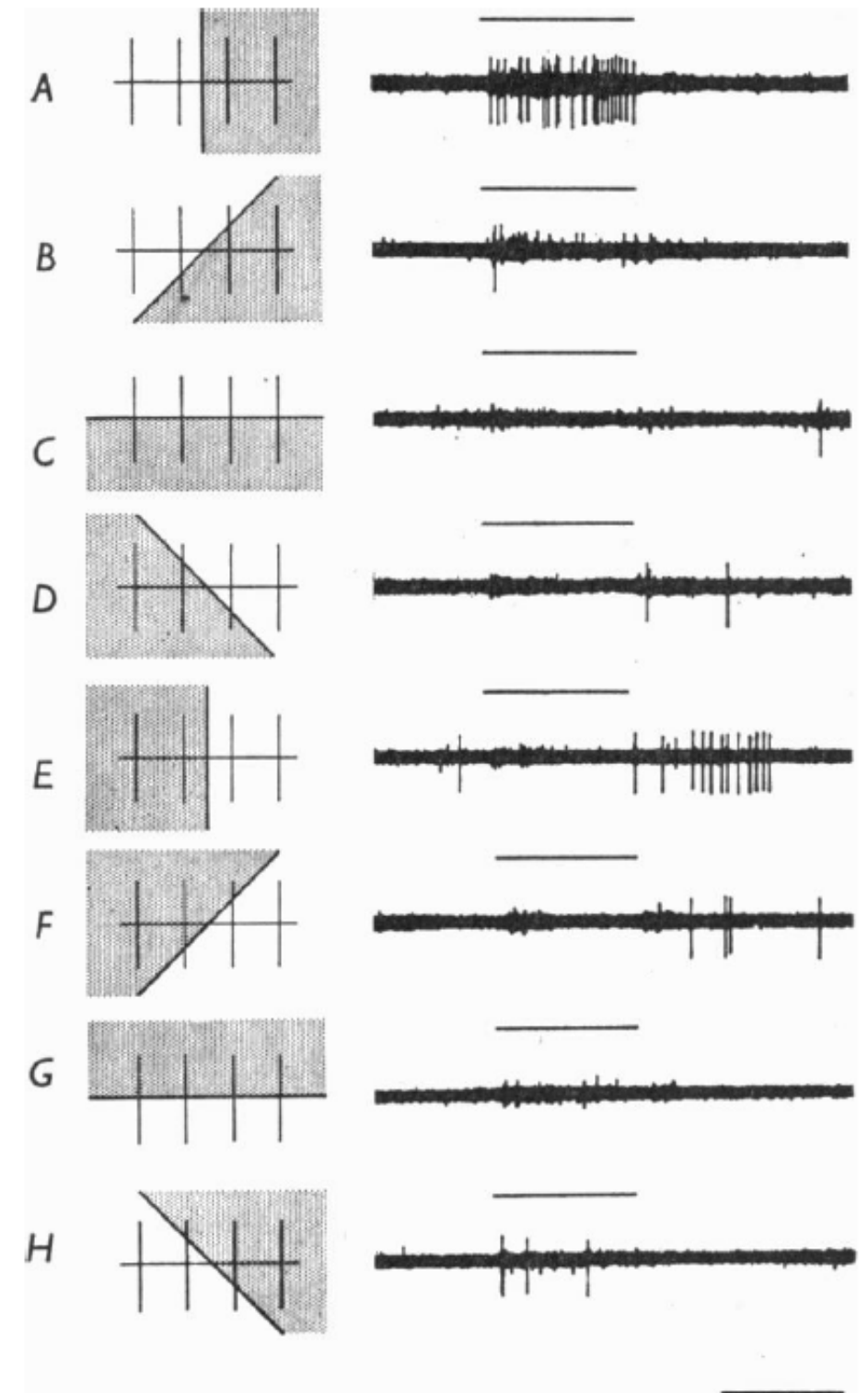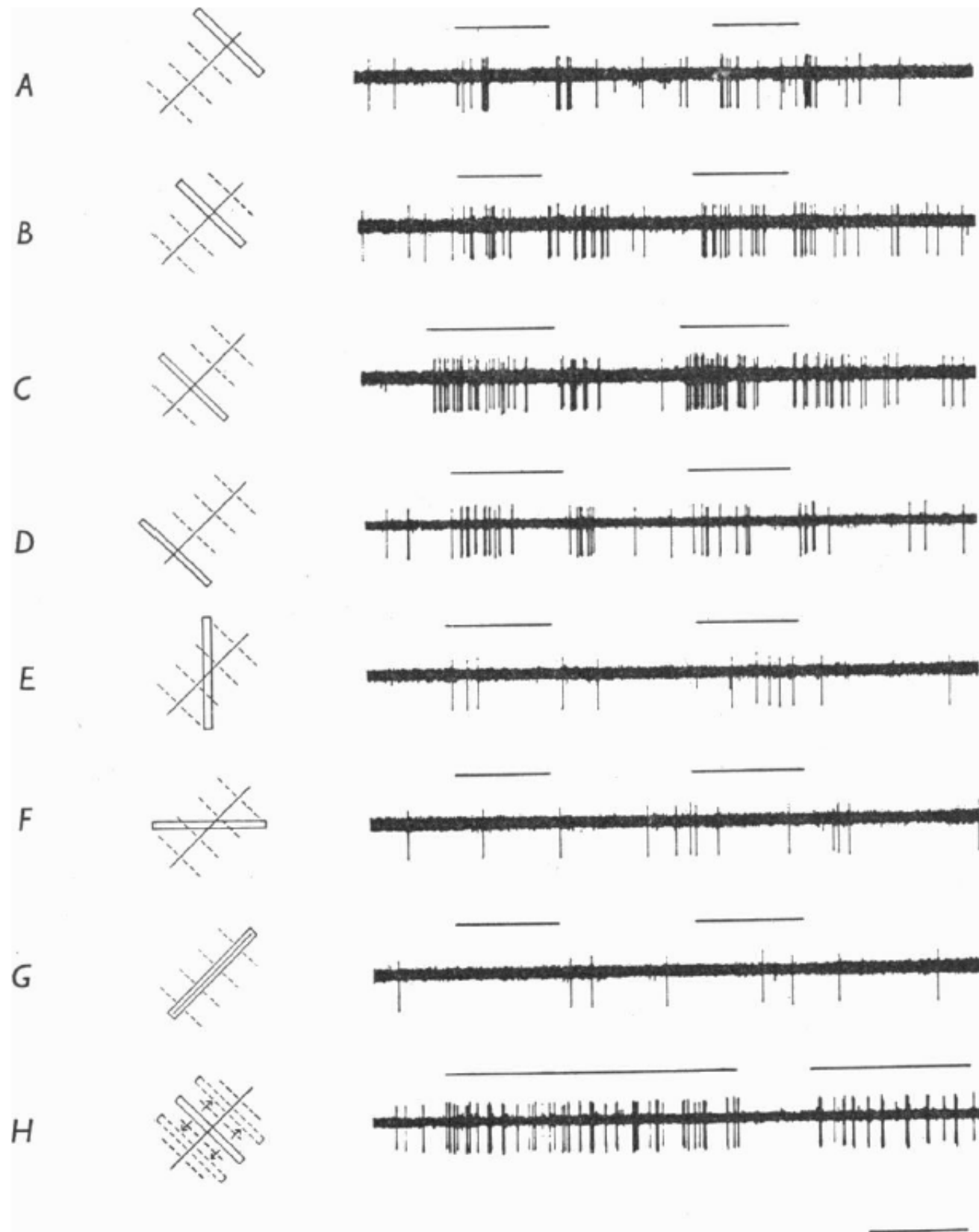
## PART I

## ORGANIZATION OF RECEPTIVE FIELDS IN CAT'S VISUAL CORTEX: PROPERTIES OF 'SIMPLE' AND 'COMPLEX' FIELDS

Stimulus:  on  off

Light stimulus

OFF
ON
OFF
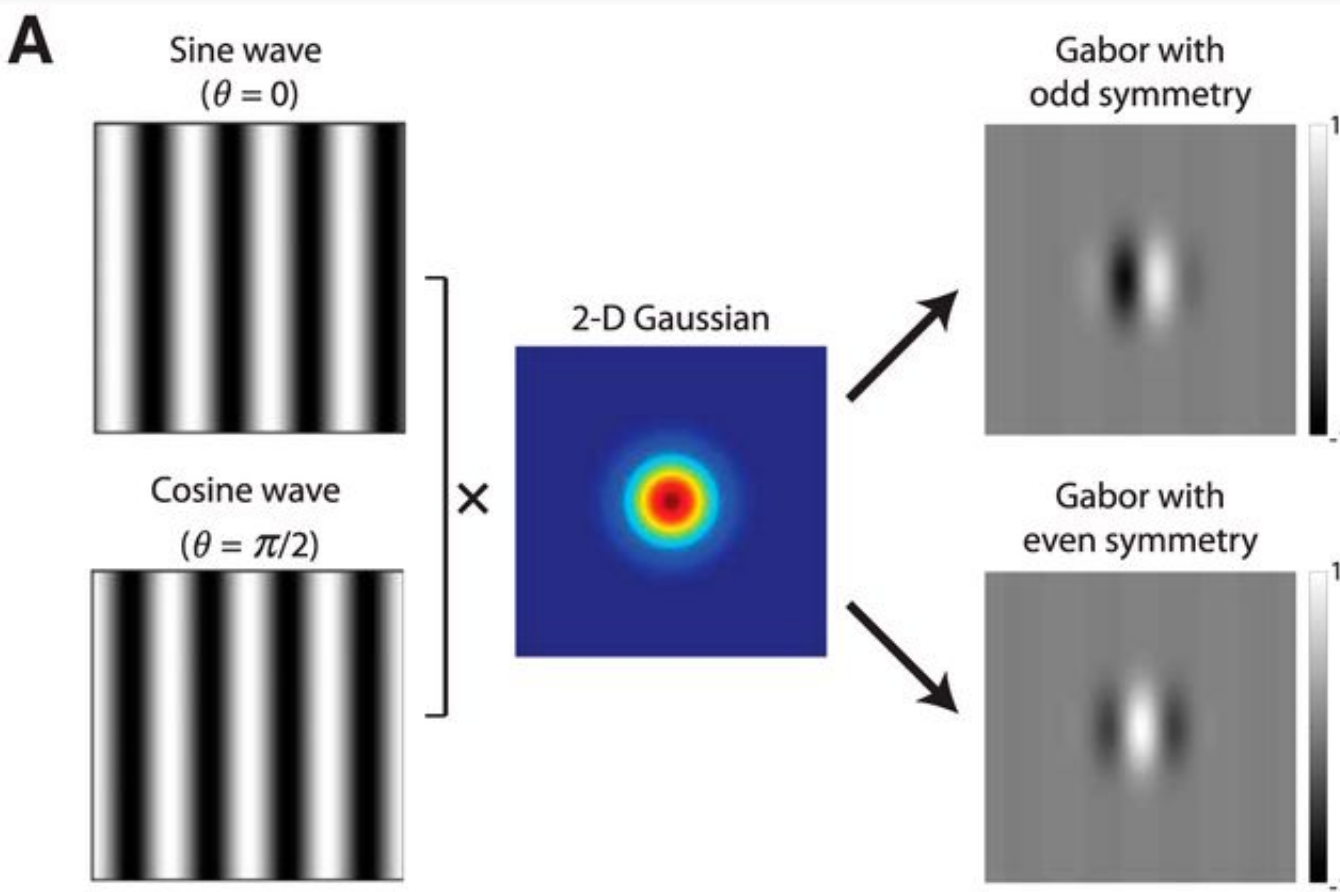ON
OFF
ON
OFF

Receptive field
of LGN cell

2  1
3

OFF ON OFF

Receptive field
of a simple cell in
primary visual cortex

1 - Active

2 - Inhibited

3 - No change

Light on

0    1    2    3
Time (sec)

from Ayzenshtat et al (2016)

*Lauritzen et al 2001*

Orientation tuning

*Lauritzen et al 2001*

$$\textbf{Circular Variance} = 1 - \frac{\sum_k r_k \mathbf{e}^{2i\theta_k}}{\sum_k r_k}$$

$r_k$ = neuron **r**'s response to stimulus with pure orientation **k**

## Orientation tuning

*Lauritzen et al 2001*

—— 40% base
– – 40% base + 40% mask
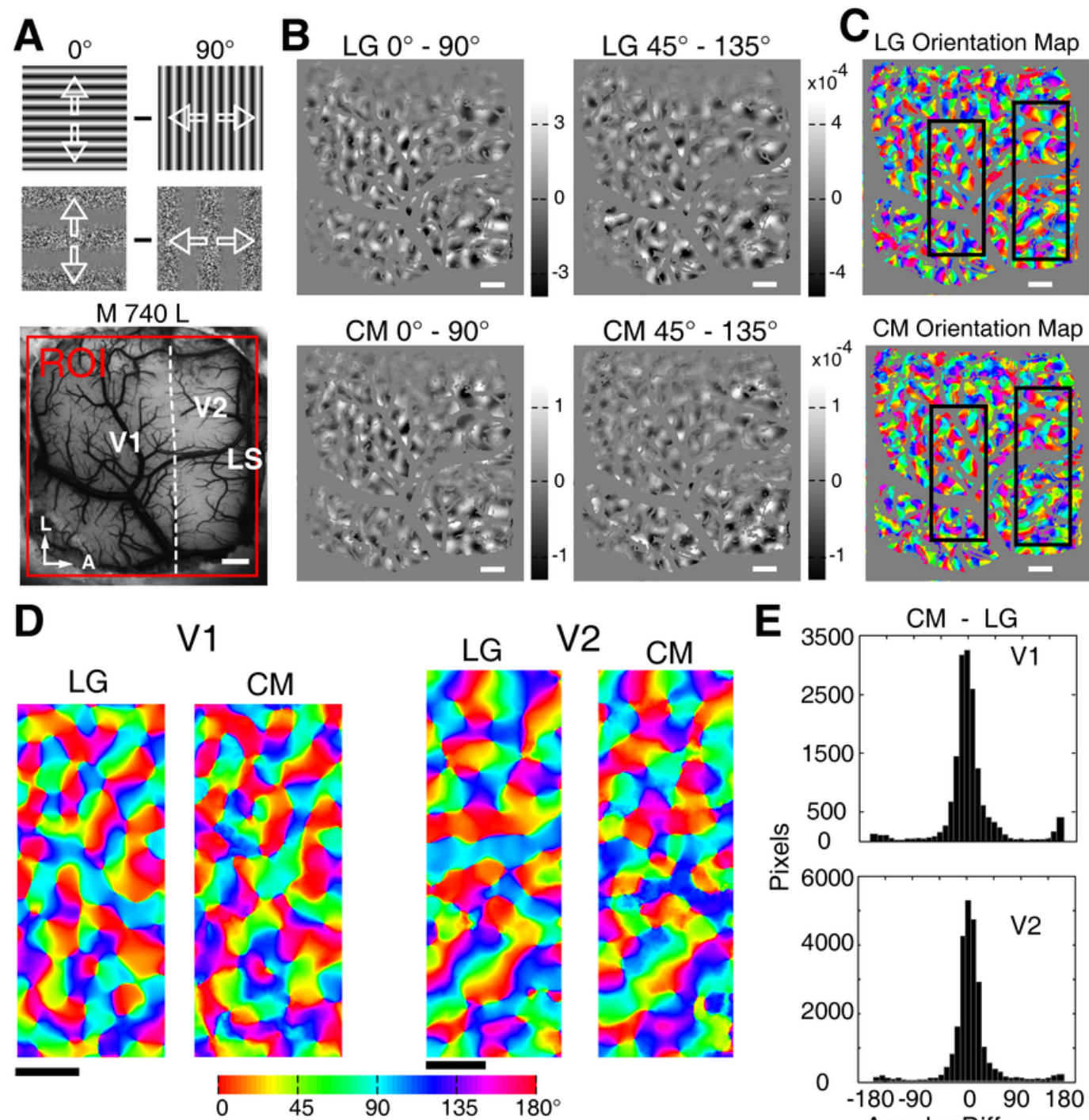· · · · 40% mask + 40% mask normalized

to make it "circular"

$$\textbf{Circular Variance} = 1 - \frac{\sum_k r_k \boxed{\mathbf{e}^{2i\theta_k}}}{\boxed{\sum_k r_k}}$$
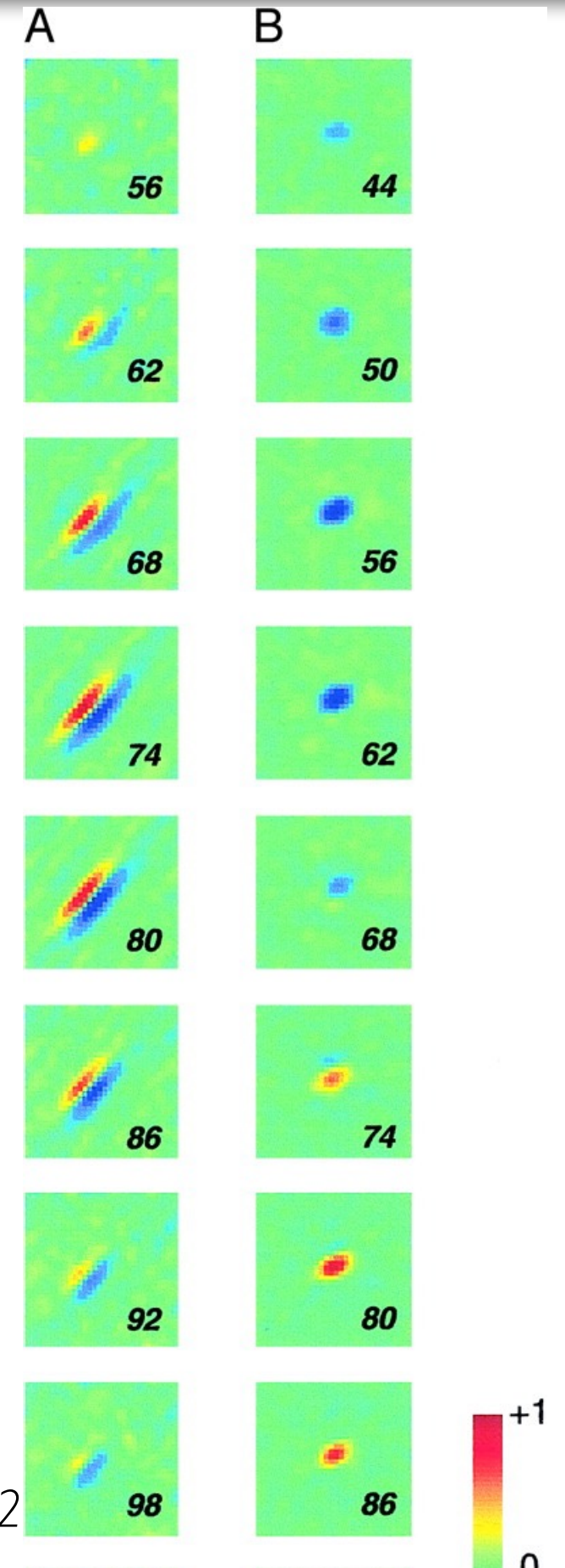
total response

$r_k$ = neuron **r**'s response to stimulus with pure orientation **k**

*An et al 2015*

*Ringach 2002*

# Simple V1 cells  Daugman, 1985



2D Receptive Field

2D Gabor Function

Difference
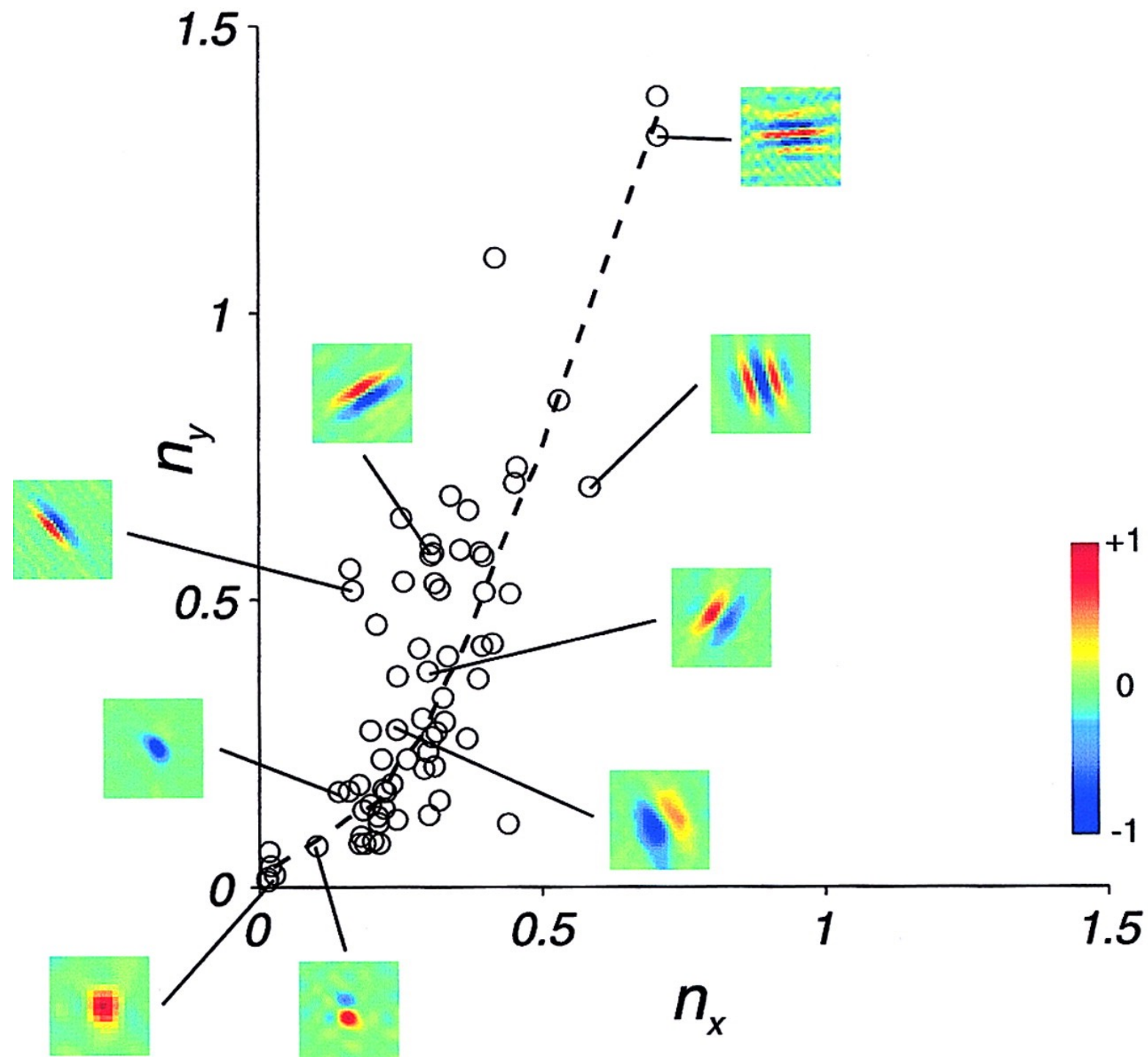
Receptive fields in primary visual cortex (Jones and Palmer, 1987)

**Gabor wavelets:** localized sine and cosine waves

$$G(x) \propto \exp\{-\frac{1}{2}[\frac{x_1^2}{\sigma_1^2} + \frac{x_2^2}{\sigma_2^2}]\}e^{ix_1}$$

Transation, rotation, dilation of the above function

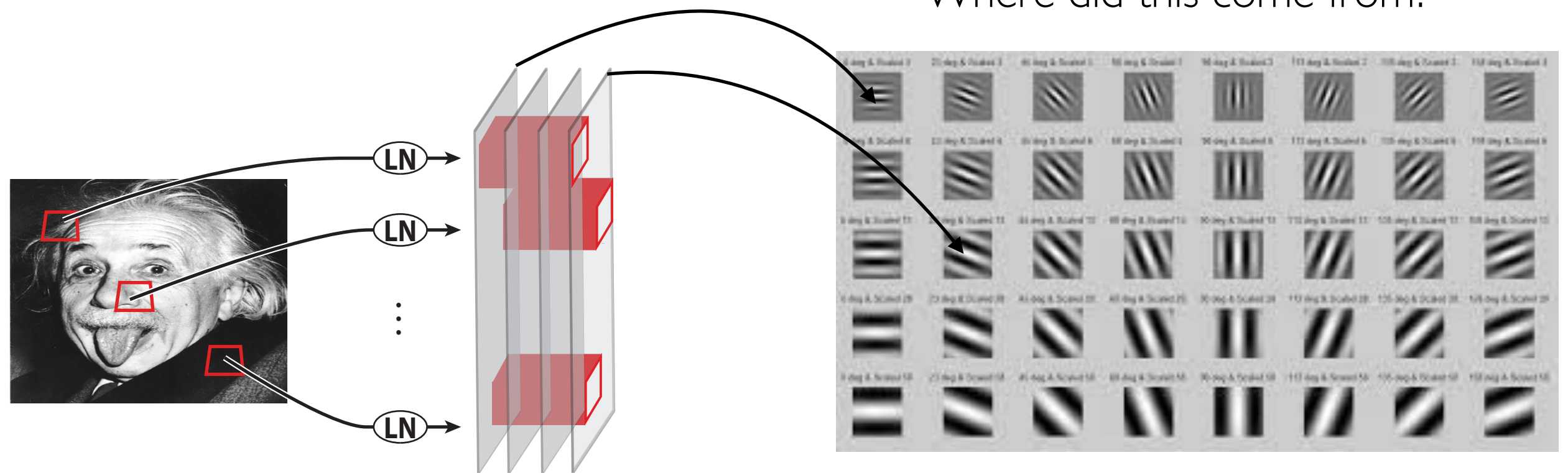There is a frequency-orientation relationship:



*from Ringach 2002*

Where did this come from?

**Two strategies to find the correct parameters.**
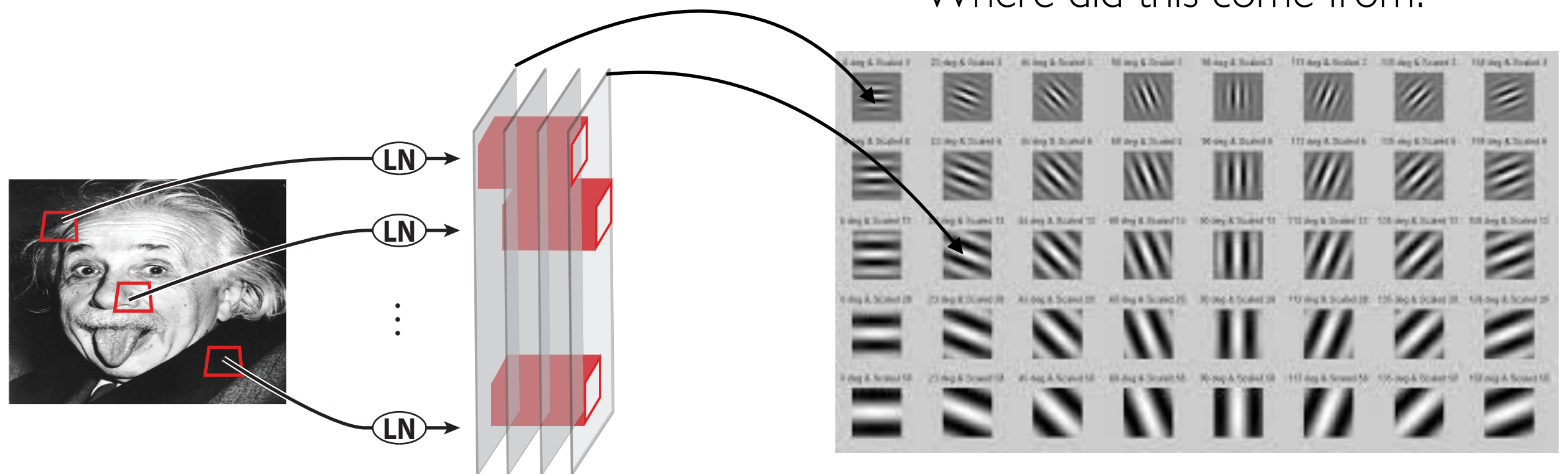
*less normative theory*

*more normative theory*

1. Fit neural data

2. Solve a high-level ecological task

…

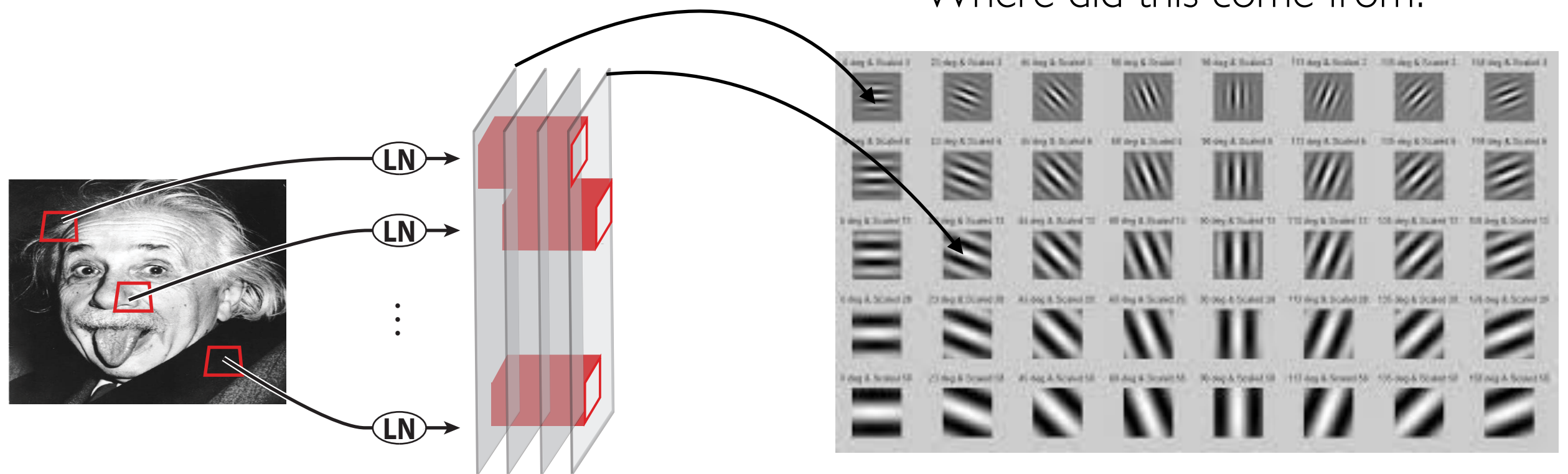compare to neural data *and* Turing Test

Where did this come from?



(1) "Hubel and Wiesel's Intuition"
    ~1970s and formalized later

→ e.g. there is a "fixed basis set" that just "makes sense" if we're smart enough

Where did this come from?



(1) "Hubel and Wiesel's Intuition"
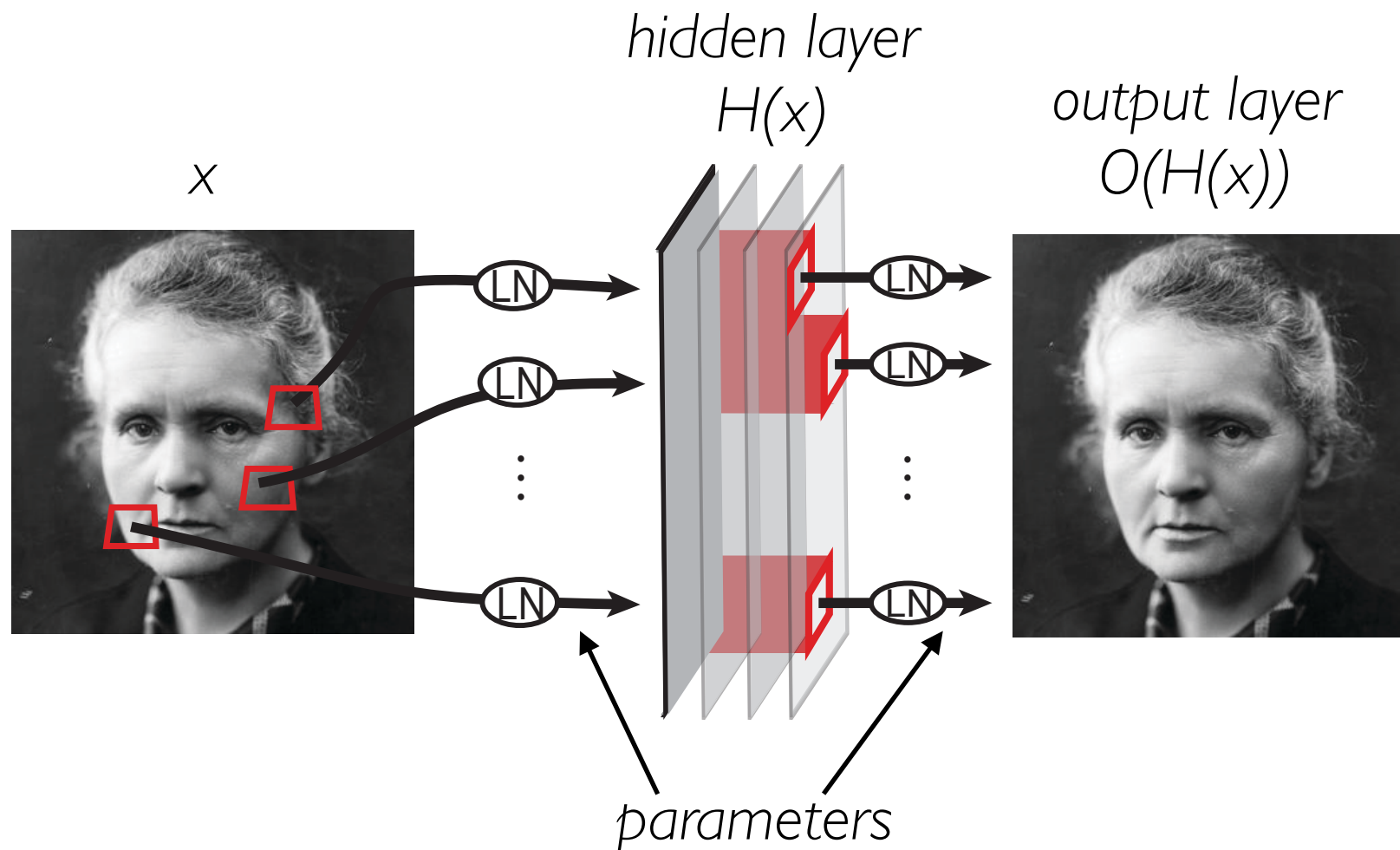    ~1970s and formalized later

→ e.g. there is a "fixed basis set" that just "makes sense" if we're smart enough

(2) Sparse Coding Foldiak, Olshausen, mid 1990s

→ neurons have to represent their environment, as efficiently as possible

*x*

hidden layer
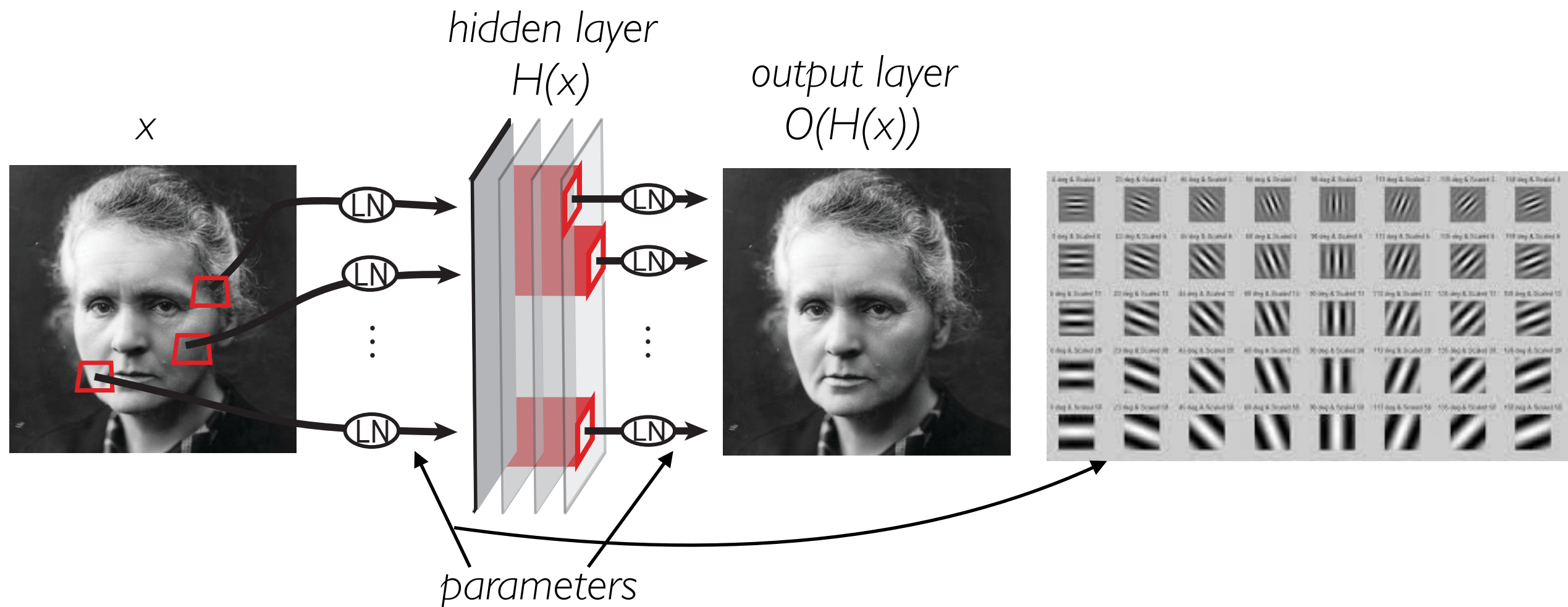*H(x)*

output layer
*O(H(x))*

parameters

$$L(x) = |x - O(H(x))|^2 + \lambda \cdot |H(x)|$$

(2) Sparse Coding Foldiak, Olshausen,
mid 1990s

→ neurons have to represent their
environment, as efficiently as possible

# Models of V1



*hidden layer*
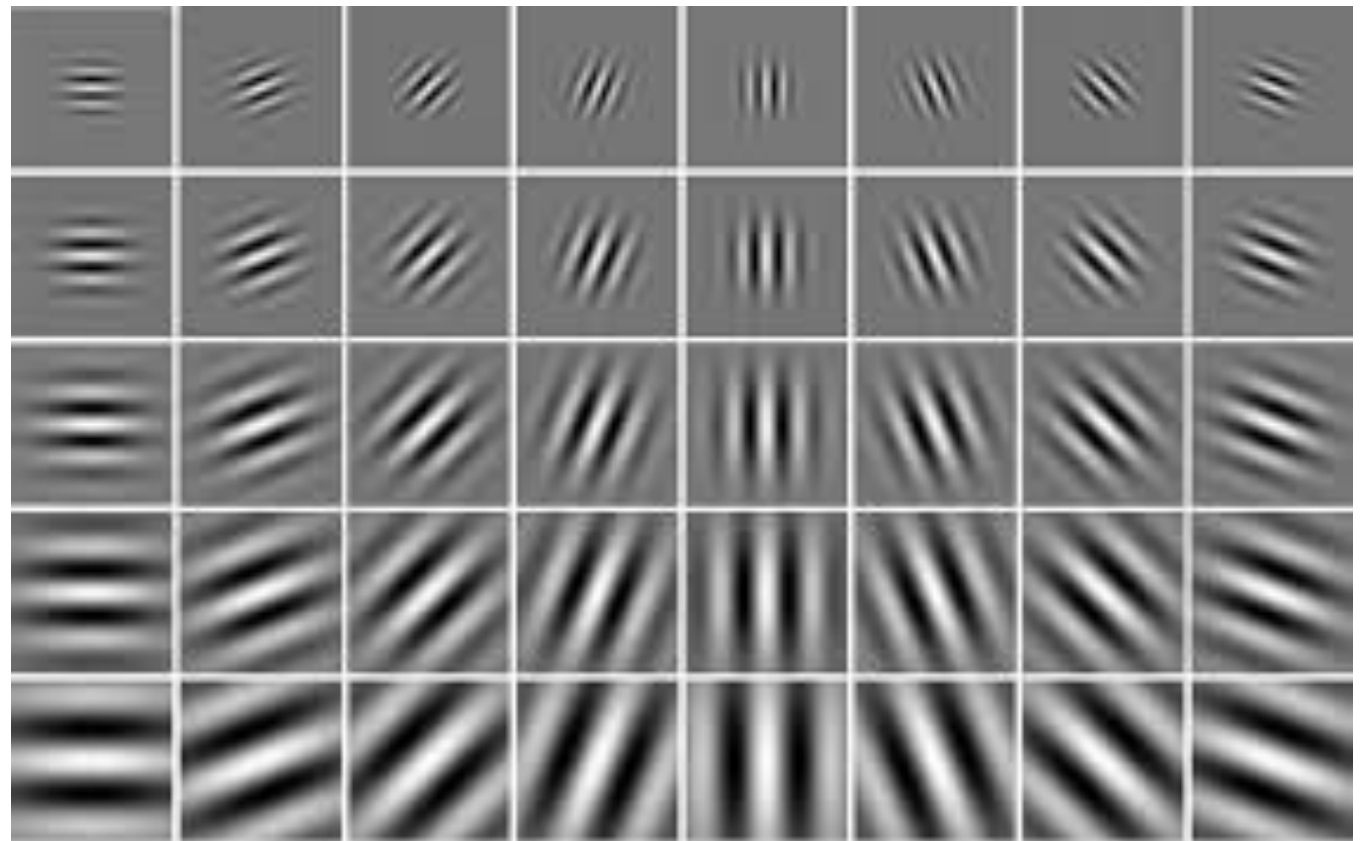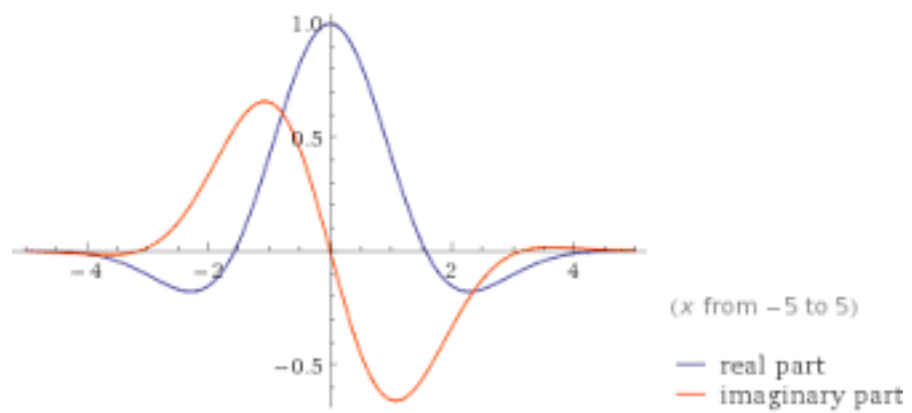*H(x)*

*output layer*
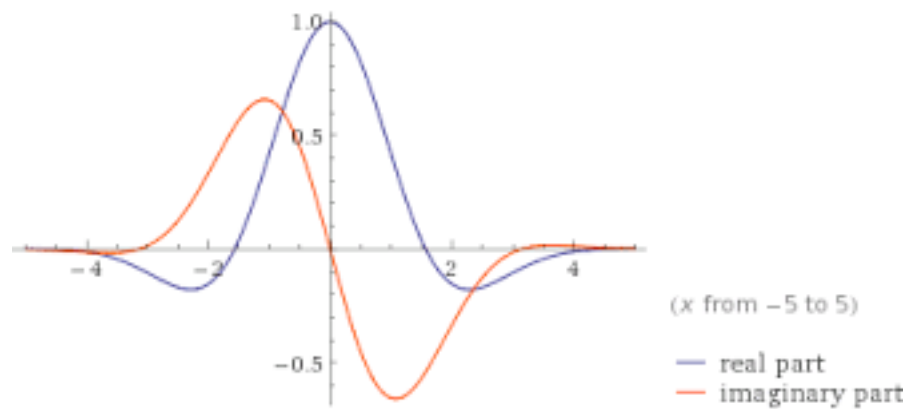*O(H(x))*

*x*

*parameters*

$$L(x) = |x - O(H(x))|^2 + \lambda \cdot |H(x)|$$

(2) Sparse Coding Foldiak, Olshausen,
   mid 1990s

→neurons have to represent their
environment, as efficiently as possible

(x from −5 to 5)

— real part
— imaginary part
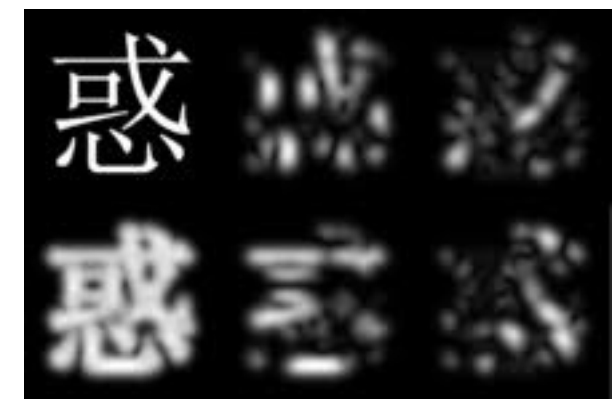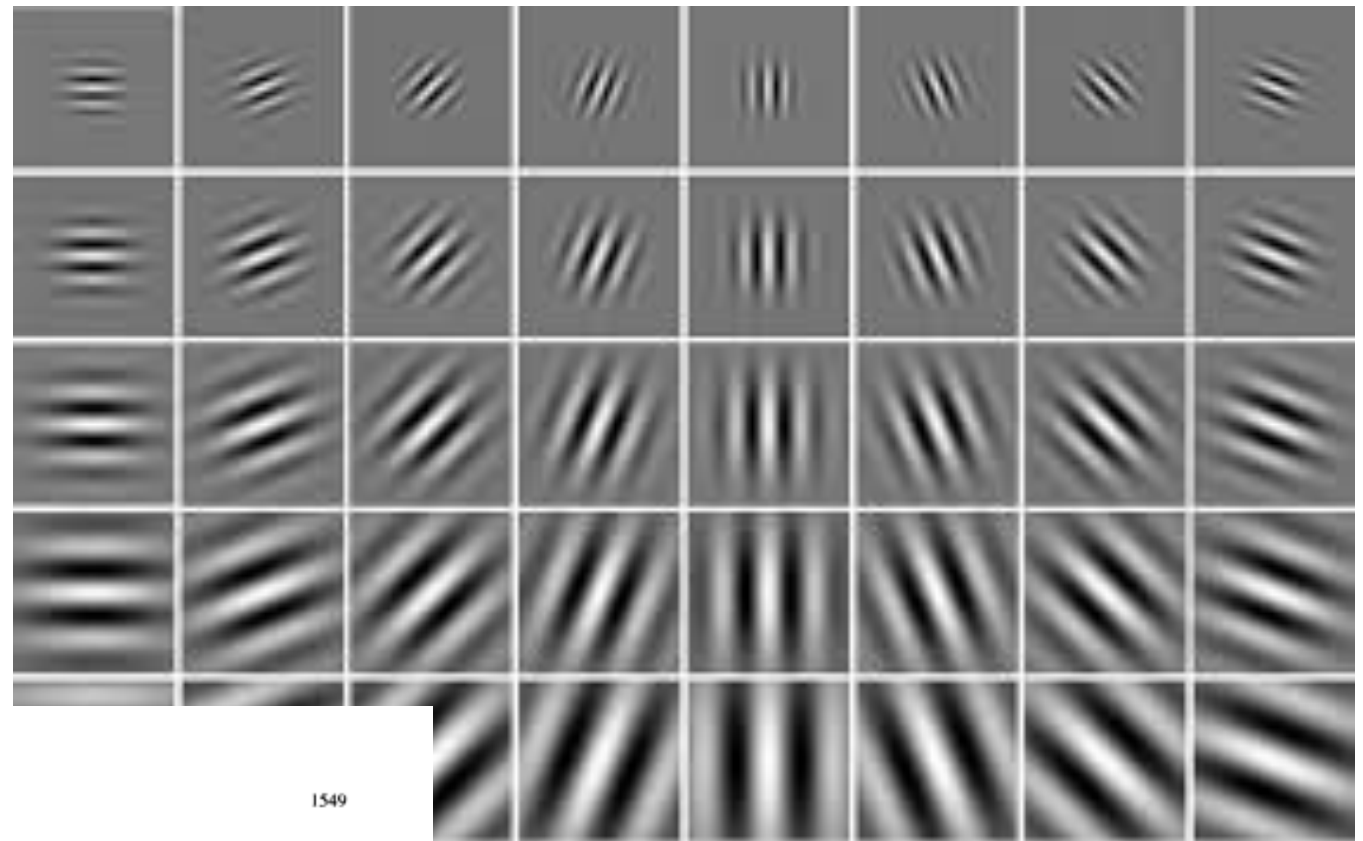
(x from −5 to 5)

— real part
— imaginary part

1549

# Texture Classification and Segmentation Using Wavelet Frames

Michael Unser, *Senior Member, IEEE*

*Abstract*—This paper describes a new approach to the characterization of texture properties at multiple scales using the wavelet transform. The analysis uses an overcomplete wavelet decomposition, which yields a description that is translation invariant. It is shown that this representation constitutes a tight frame of $l_2$ and that it has a fast iterative algorithm. A texture is characterized by a set of channel variances estimated at the output of the corresponding filter bank. Classification experiments with 12 Brodatz textures indicate that the discrete wavelet frame (DWF) approach is superior to a standard (critically sampled) wavelet transform feature extraction. These results also suggest that this approach should perform better than most traditional single resolution techniques (co-occurrences, local linear transform, and the like). A detailed comparison of the classification performance of various orthogonal and biorthogonal wavelet transforms is also provided. Finally, the DWF feature extraction technique is incorporated into a simple multicomponent texture segmentation algorithm, and some illustrative examples are presented.

reversible, which limits their applicability for texture synthesis. Most of these problems can be avoided if one uses the wavelet transform, which provides a precise and unifying framework for the analysis and characterization of a signal at different scales [16]–[19]. The use of a pyramid-structured wavelet transform for texture analysis was first suggested in the pioneering work of Mallat [19]. This initial proposal has been followed by several studies on texture classification with a particular attention to the use of wavelet packets [20], [21], which constitute a multiband extension of the pyramid-structured wavelet transform.
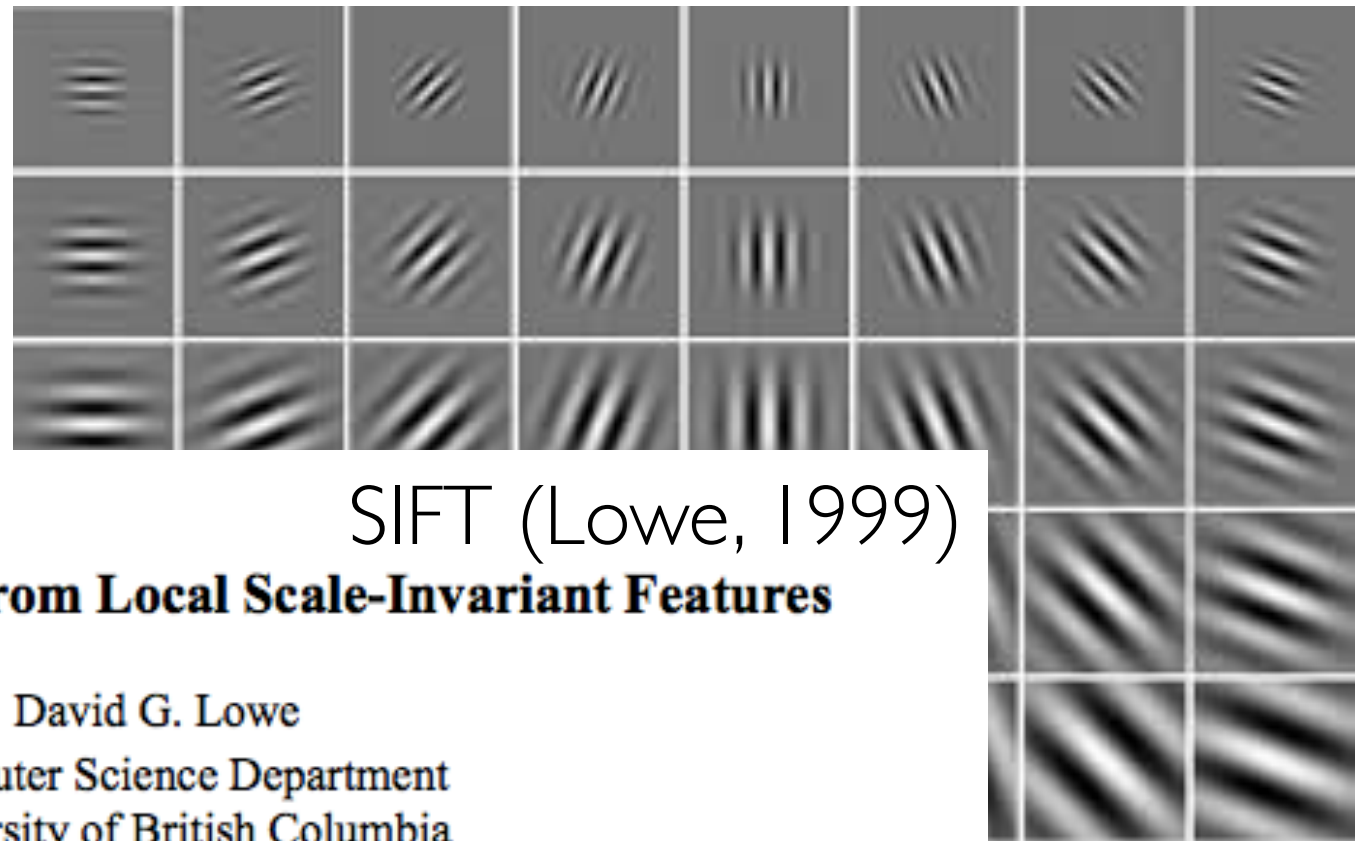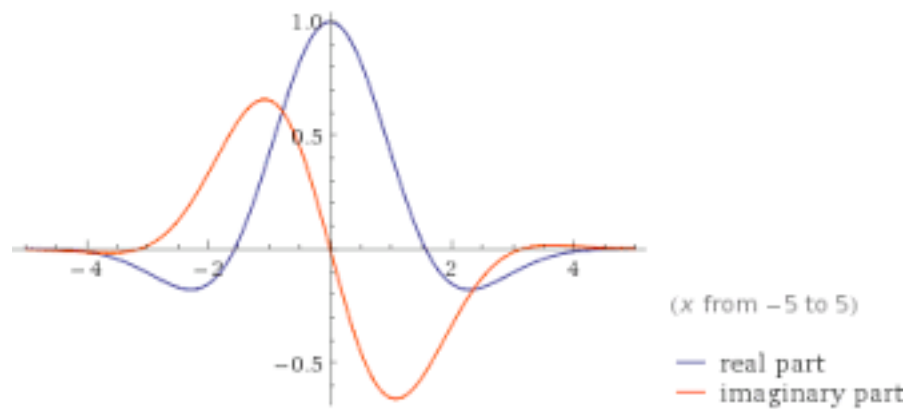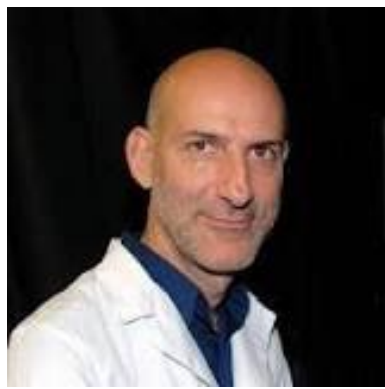
In this paper, a variation of the discrete wavelet transform is introduced for characterizing texture properties. This technique is applied to the problems of texture classification and segmentation. The present analysis method, which is described in Section II, uses an overcomplete wavelet decomposition (the discrete wavelet frame (DWF)) in which the output of



Many CV careers made on wavelets.

(x from −5 to 5)

— real part
— imaginary part



SIFT (Lowe, 1999)

## Object Recognition from Local Scale-Invariant Features

David G. Lowe

Computer Science Department
University of British Columbia
Vancouver, B.C., V6T 1Z4, Canada
lowe@cs.ubc.ca

IEEE TRANSACTIONS ON IMAGE PROCE

Segme

Michael Unser, *Senior Member, IEEE*

*Abstract*—This paper describes a new approach to the characterization of texture properties at multiple scales using the wavelet transform. The analysis uses an overcomplete wavelet decomposition, which yields a description that is translation invariant. It is shown that this representation constitutes a tight frame of $l_2$ and that it has a fast iterative algorithm. A texture is characterized by a set of channel variances estimated at the output of the corresponding filter bank. Classification experiments with 12 Brodatz textures indicate that the discrete wavelet frame (DWF) approach is superior to a standard (critically sampled) wavelet transform feature extraction. These results also suggest that this approach should perform better than most traditional single resolution techniques (co-occurrences, local linear transform, and the like). A detailed comparison of the classification performance of various orthogonal and biorthogonal wavelet transforms is also provided. Finally, the DWF feature extraction technique is incorporated into a simple multicomponent texture segmentation algorithm, and some illustrative examples are presented.

reversible, which limits their applicability for texture synthesis. Most of these problems can be avoided if one uses the wavelet transform, which provides a precise and unifying framework for the analysis and characterization of a signal at different scales [16]–[19]. The use of a pyramid-structured wavelet transform for texture analysis was first suggested in the pioneering work of Mallat [19]. This initial proposal has been followed by several studies on texture classification with a particular attention to the use of wavelet packets [20], [21], which constitute a multiband extension of the pyramid-structured wavelet transform.
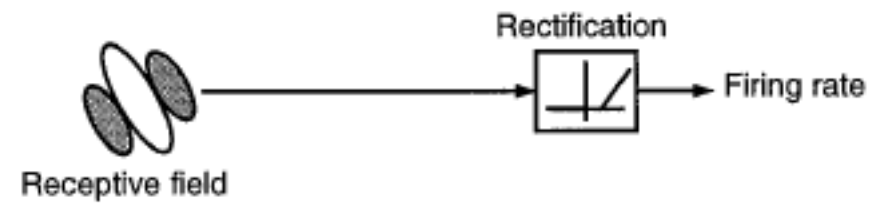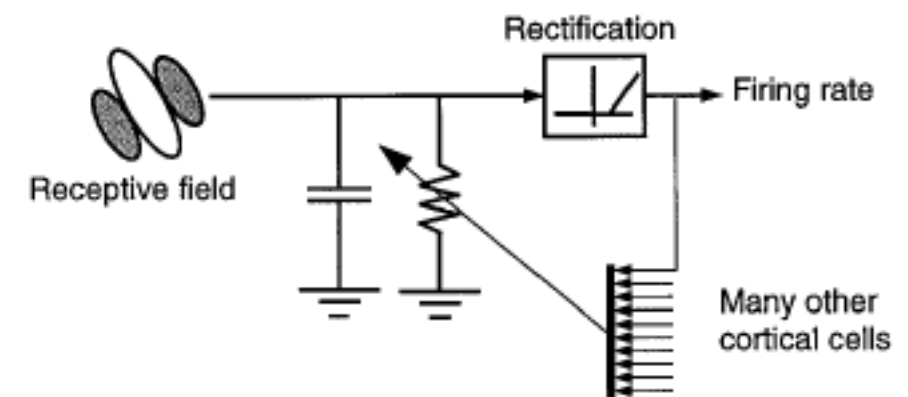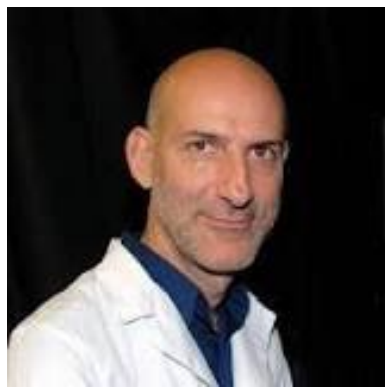
In this paper, a variation of the discrete wavelet transform is introduced for characterizing texture properties. This technique is applied to the problems of texture classification and segmentation. The present analysis method, which is described in Section II, uses an overcomplete wavelet decomposition (the discrete wavelet frame (DWF)) in which th



Many CV careers made on wavelets.

**A** Linear model

Rectification

Firing rate

Receptive field

**B** Normalization model

Rectification

Firing rate

Receptive field

Many other
cortical cells

**A** Linear model

Rectification

Firing rate

Receptive field

**B** Normalization model
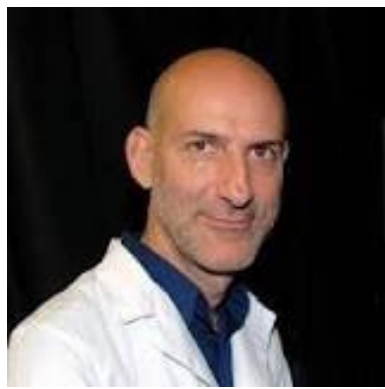
Rectification

Firing rate

Receptive field

Many other cortical cells

$$C\frac{dV}{dt} + gV = I$$

$$g = \frac{g_0}{\sqrt{1 - k \cdot \sum_{r \in R_x} r}}$$

$$R = max(0, V)$$

*Carandini, Heeger and Movshon (1997)*

**A** Linear model

**B** Normalization model

R

$$C\frac{dV}{dt} + gV = I$$

$$g = \frac{g_0}{\sqrt{1 - k \cdot \sum_{r \in R_x} r}}$$

$$R = max(0, V)$$

measure R from neural data

solve diff eq for equilibrium, estimate free parameters:    C, k, g0

**A** Linear model

Rectification

Firing rate

Receptive field

**B** Normalization model

Rectification

Firing rate

Receptive field

Many other cortical cells

$$y = R\,[\,W * x\,]$$

$$y = R\,[\,norm(W * x)\,]$$

OR: derive this expression: ——>
(basically)

$$norm(x) \sim \frac{x}{\left(\gamma + \alpha \cdot \displaystyle\sum_{r \in R_x} x_r^2\right)^{\beta}}$$

**A** Linear model

Rectification

Firing rate

Receptive field

**B** Normalization model

Rectification

Firing rate

Receptive field

Many other cortical cells

$$y = R\,[\,W * x\,]$$

$$y = R\,[\,norm(W * x\,)]$$

OR: derive this expression: ——> (basically)

$$norm(x) \sim \frac{x}{\left(\gamma + \alpha \cdot \sum_{r \in R_x} x_r^2\right)^{\beta}}$$

NB: (1) derivation involves "reasonable" assumption that "the normalization pool to contain quadruples of cells with the same amplitude response but with phases 90° apart." (2) **The above is how we now *define* local response normalization**

Interrelations and effects of the principal variables in the normalization model.



Relation between membrane potential and firing rate.

Relation between pool activity and membrane conductance.

Effects of conductance on the size and time course of the membrane potential.

The **I** in this equation:

$$C\frac{dV}{dt} + gV = I$$

is a sinusoid

neural response ~ A * sin (w t + d)

A = amplitude of cell          w = frequency of the drift          d = phase of cell

A, d are fit to the data

Responses to drifting sine gratings of different contrasts

What functions are A and d of a stimulus parameter — contrast?

the parameters **C, k, A, d** (basically) are fit to the neuron over a bunch of stimulus conditions

histogram of responses for different contrasts



response vs phase

each point is response to a different sinusoid

Now as a function of grating orientation

response vs phase

gray = -15deg, white = -45deg

*Carandini, Heeger and Movshon (1997)*



gray = 1.4 cyc/deg, white = 1.1 cyc/deg

Now as a function of temporal frequency

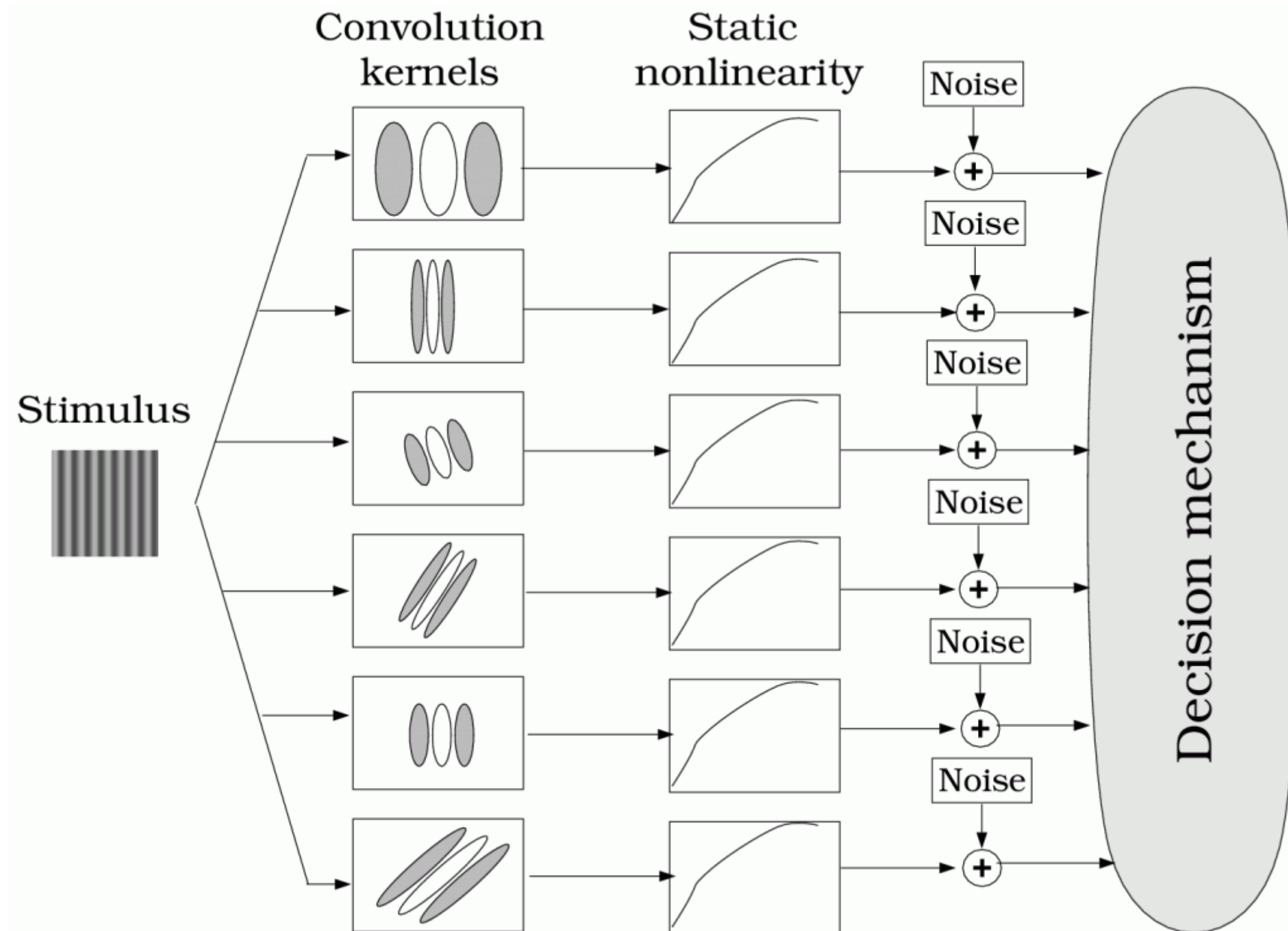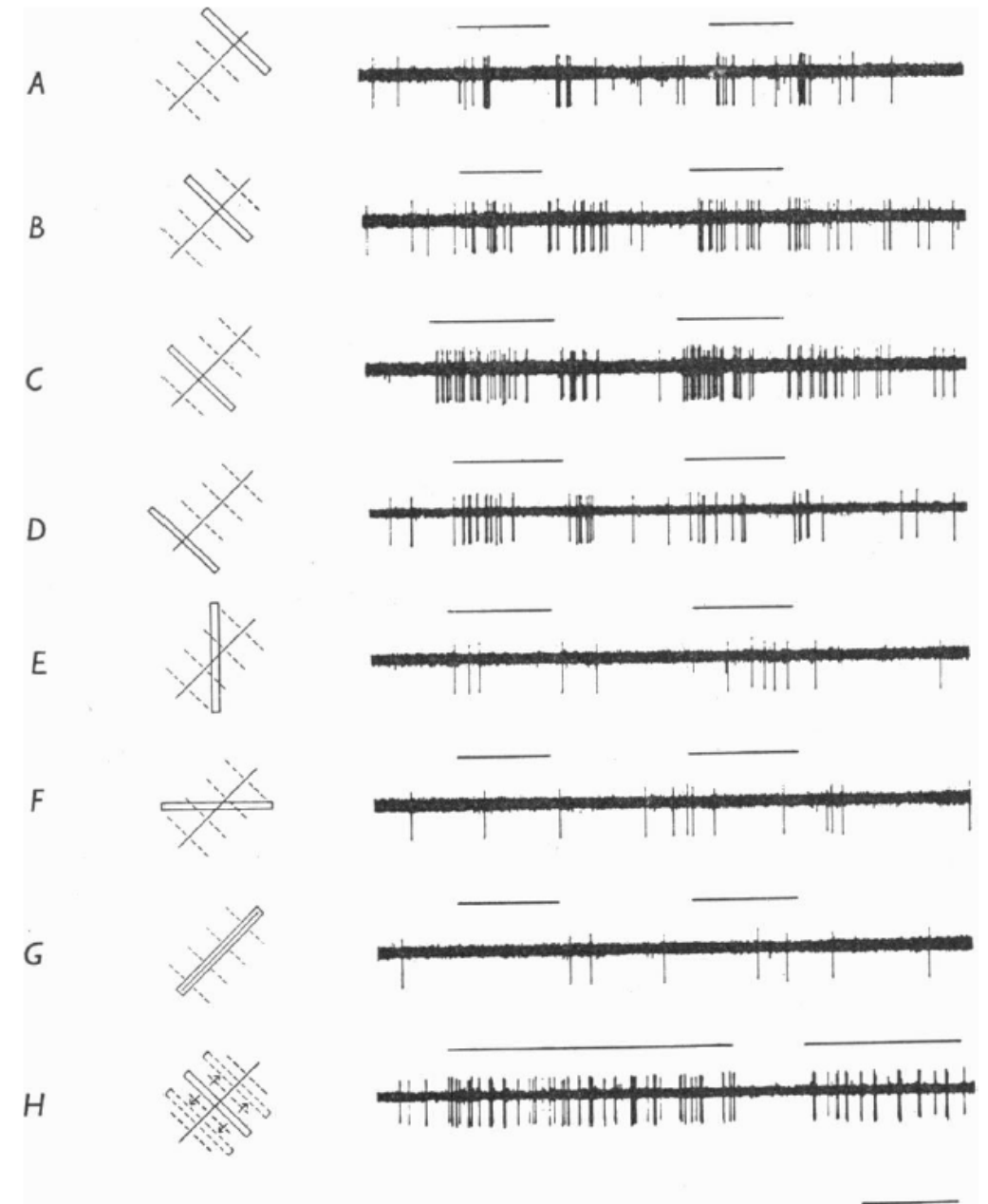*Carandini, Heeger and Movshon (1997)*

Colors
as in
panel **A**

*from Wandell 1996*

More generally, it was realized in computer vision that **pooling** was a good idea.

    recall Hubel & Wiesel's complex cell >>

$$y = \left( \frac{1}{|N_r|} \sum_{i \in N_r} x_i^p \right)^{1/p}$$
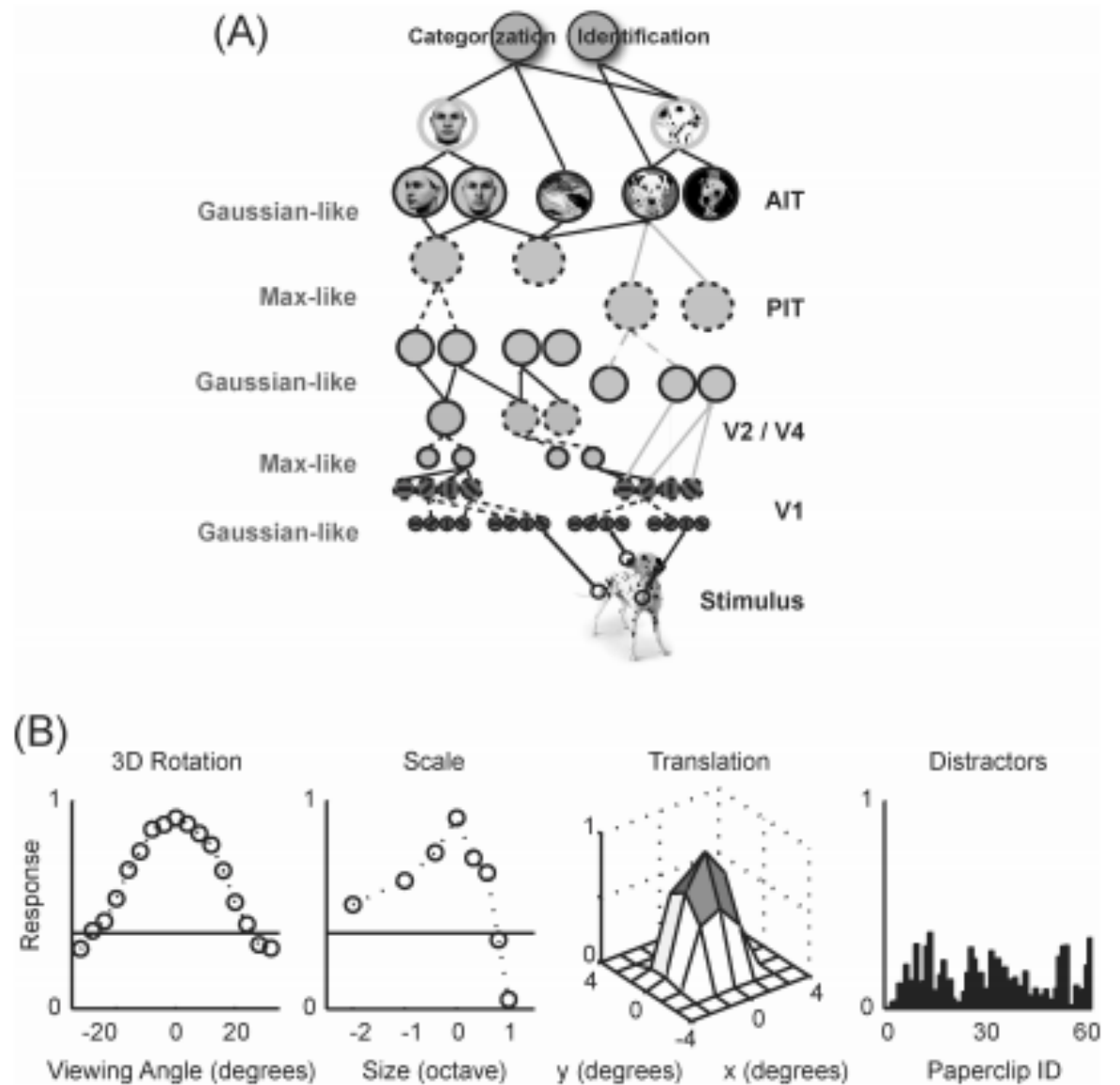


*(Actually, if you're running a CNN, you basically \*have\* to do pooling + downsampling, for memory reasons.)*

More generally, it was realized in computer vision that **pooling** was a good idea.

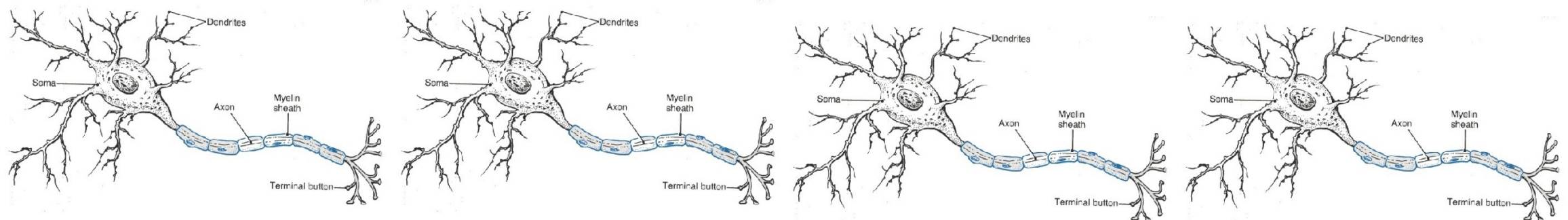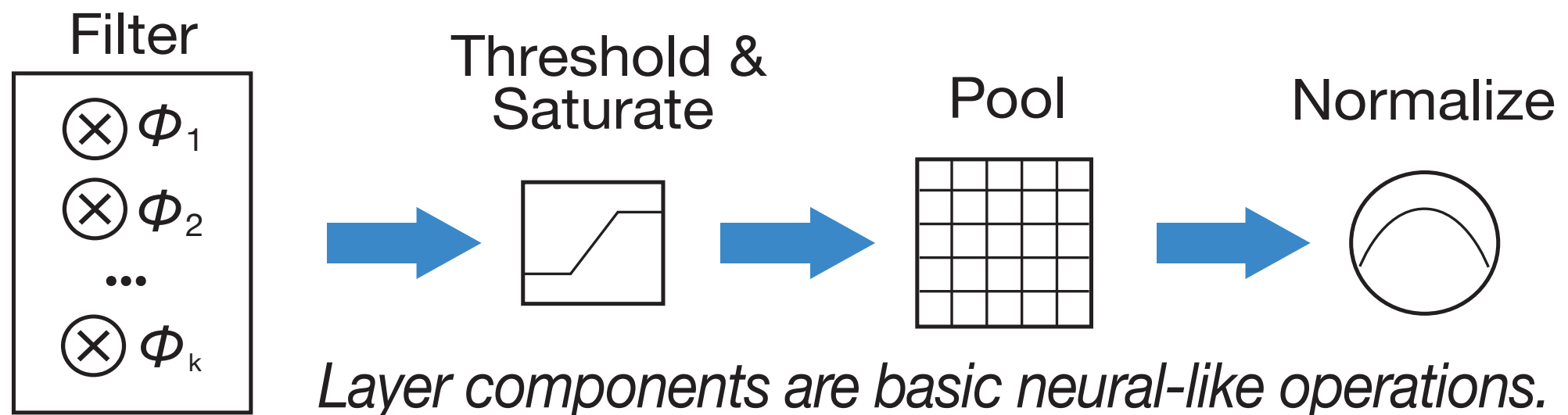$$y = \left( \frac{1}{|N_r|} \sum_{i \in N_r} x_i^p \right)^{1/p}$$



*from Kouh and Poggio (2008)*

*(Actually, if you're running a CNN, you basically \*have\* to do pooling + downsampling, for memory reasons.)*
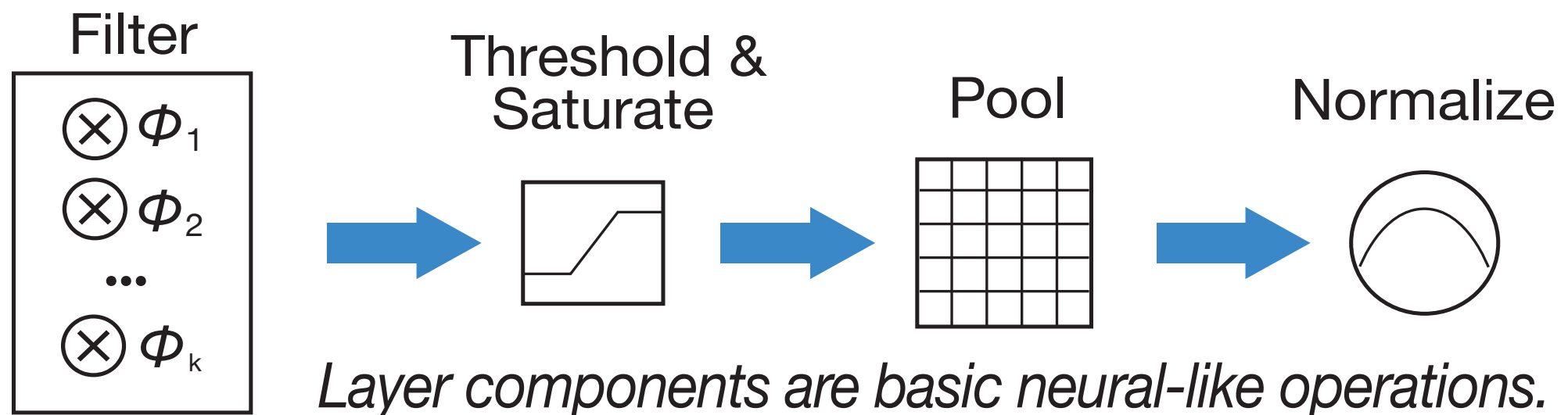
▶ Linear-Nonlinear neurally-plausible **basic operations** within layer



Filter

$\otimes \Phi_1$

$\otimes \Phi_2$

...

$\otimes \Phi_k$

Threshold & Saturate

Pool

Normalize

*Layer components are basic neural-like operations.*

▶ Linear-Nonlinear neurally-plausible **basic operations** within layer



Filter

$\bigotimes \phi_1$

$\bigotimes \phi_2$

...

$\bigotimes \phi_k$

Threshold & Saturate

Pool

Normalize

*Layer components are basic neural-like operations.*

**neuro:** synaptic weights patterns

**data:** untangling through dimension expansion

Hubel and Wiesel (1965-1975) Lecun (2004), Carandini et. al (2005), Lennie & Movshon (2005), DiCarlo (2012)

▶ Linear-Nonlinear neurally-plausible **basic operations** within layer



| | Filter | Threshold & Saturate | Pool | Normalize |
|---|---|---|---|---|

*Layer components are basic neural-like operations.*

**neuro:**

synaptic weights patterns

single-unit activations

**data:**

untangling through dimension expansion

"AND" operation by limiting dynamic range

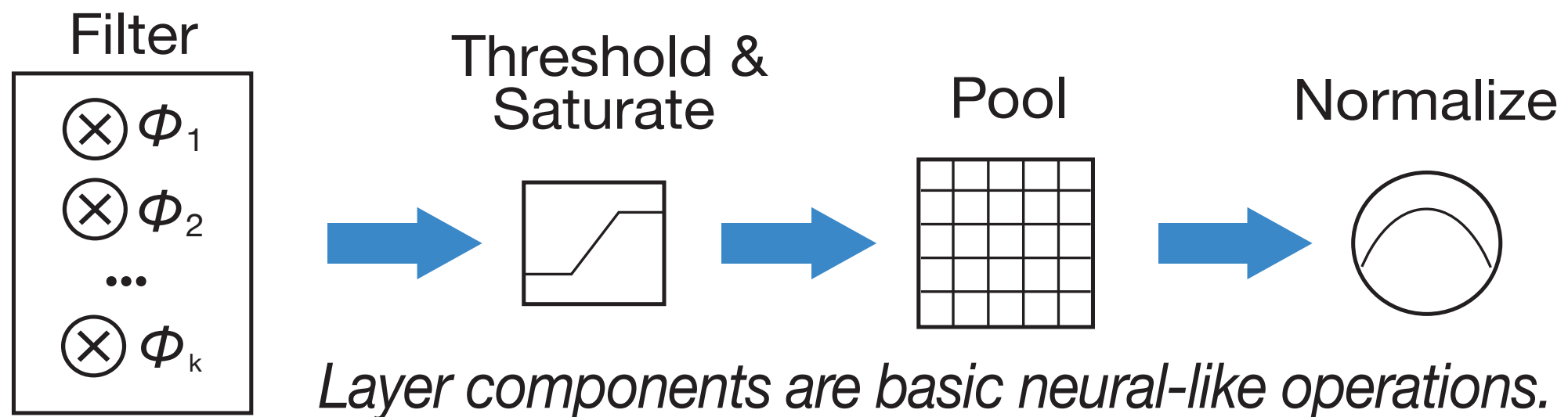Hubel and Wiesel (1965-1975) Lecun (2004), Carandini et. al (2005), Lennie & Movshon (2005) , DiCarlo (2012)

# Linear-Nonlinear Operations

▸ Linear-Nonlinear neurally-plausible **basic operations** within layer



| | Filter | Threshold & Saturate | Pool | Normalize |
|---|---|---|---|---|
| | $\otimes \Phi_1$ $\otimes \Phi_2$ ... $\otimes \Phi_k$ | | | |

*Layer components are basic neural-like operations.*

| | | | |
|---|---|---|---|
| **neuro:** | synaptic weights patterns | single-unit activations | complex cells |
| **data:** | untangling through dimension expansion | "AND" operation by limiting dynamic range | adding robustness by dimension reduction |

Hubel and Wiesel (1965-1975) Lecun (2004), Carandini et. al (2005), Lennie & Movshon (2005), DiCarlo (2012)

▶ Linear-Nonlinear neurally-plausible **basic operations** within layer



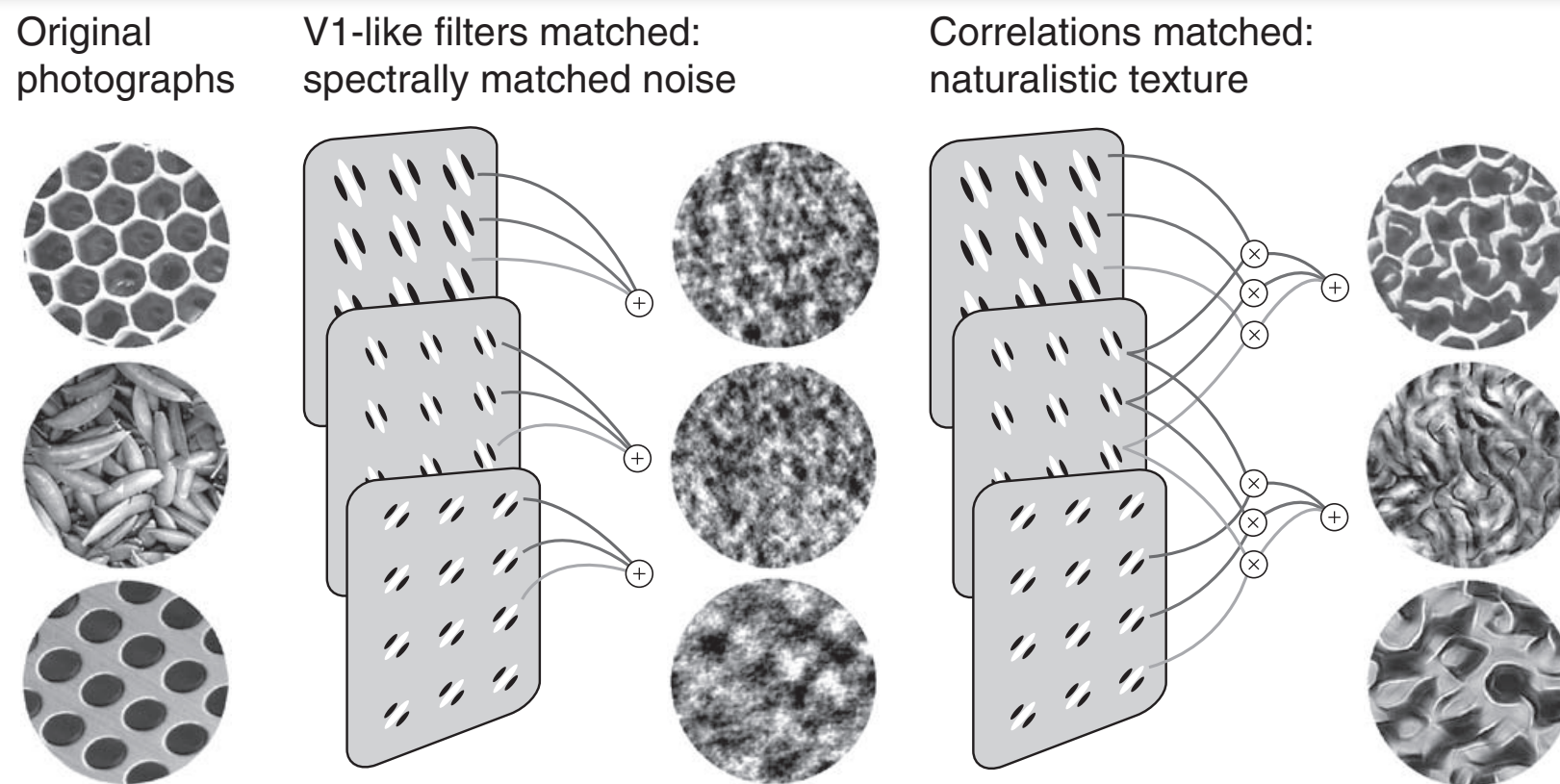| | Filter | Threshold & Saturate | Pool | Normalize |
|---|---|---|---|---|
| **neuro:** | synaptic weights patterns | single-unit activations | complex cells | competitive inhibition |
| **data:** | untangling through dimension expansion | "AND" operation by limiting dynamic range | adding robustness by dimension reduction | put results back into standard range |

*Layer components are basic neural-like operations.*

Hubel and Wiesel (1965-1975) Lecun (2004), Carandini et. al (2005), Lennie & Movshon (2005) , DiCarlo (2012)

▸ Linear-Nonlinear neurally-plausible **basic operations** within layer



**Linear**
Filter
$\otimes \Phi_1$
$\otimes \Phi_2$
...
$\otimes \Phi_k$

**Nonlinear**
Threshold & Saturate
Pool
Normalize

*Layer components are basic neural-like operations.*

| | | | |
|---|---|---|---|
| **neuro:** | synaptic weights patterns | single-unit activations | complex cells | competitive inhibition |
| **data:** | untangling through dimension expansion | "AND" operation by limiting dynamic range | adding robustness by dimension reduction | put results back into standard range |

Hubel and Wiesel (1965-1975) Lecun (2004), Carandini et. al (2005), Lennie & Movshon (2005) , DiCarlo (2012)

Linear-Nonlinear neurally-plausible **basic operations** within layer

**B**

| | | Latency |
|---|---|---|
| | ~10 M (IT representation) | |
| STPa | **AIT** ~16 M | ~100 ms |
| STPp | **CIT** ~17 M | ~90 ms |
| | **PIT** ~36 M | ~80 ms |
| 7a | VOT | |
| LIP MST FST | ~15 M (V4 representation) | |
| DP | **V4** ~68 M | ~70 ms |
| MIP PO MT | | |
| PIP V3A | | |
| V3 | ~29 M (V2 representation) | |
| **You are here.** | **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) | | |

*Adapted from DiCarlo et al. 2012*

Original photographs

V1-like filters matched: spectrally matched noise

Correlations matched: naturalistic texture

*Eero Simoncelli*

*Tony Movshon*

*Jeremy Freeman*

**b**

V1
n = 102

V2
n = 103

Naturalistic

Noise

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

**So, maybe a <u>hierarchically-built</u> sparse auto-encoding in a 2-layer model with max pooling??**

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

Adapted from DiCarlo et al. 2012

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

Jack Gallant

Anitha Pasupathy    Scott Brincat    Ed 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

A structural (parts-based) shape-coding scheme based on contour fragments. *A*, The example shape, a bold numeral 2, can be decomposed into contour fragments (*a-g*) with different curvatures, orientations, and positions. *B*, The curvature and orientation of each contour fragment is plotted on a 2-D domain. *C*, The positions of the contour fragments (relative to the object center) are plotted on a 2-D domain. Together, plots *B* and *C* represent a 4-D domain for describing contour fragments.

*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



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

# What shape features drive V4 response?



*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 for these types of stimuli

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

*Adapted from C.E. Connor*

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

<u>Hypothesis:</u>
V4 neurons are tuned in this basis

<u>Experimental result:</u>
Hypothesis explains ~50% of the explainable
response variance for these types of stimuli



**Problem:**
*No predictions for any other images.*
*i.e.*
*is not an "image-computable" model*

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

Adapted from DiCarlo et al. 2012

# IT statistics (rhesus monkey)

~ 7.7 cm$^2$

~ 8% of neocortex  (~ 15% of visual cortex)

~ 90 million neurons

Subregions: (PIT, CIT, AIT)     (TEO, TE)

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



*Increasing ability to drive this IT neuron -->*

*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).*

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



| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 2 | 3 | 3 | 4 | 4 | 5 | 6 |

*Increasing ability to drive this IT neuron -->*

*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).

Joyce Carol Oates!

Charlie Gross

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



Desimone et al. (1984)

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

That selectivity is tolerant to changes in position and size



Logothetis et al. (1995)

*Tanaka et al.*

Tanaka et al.

1 mm    Tsunoda et al.

## Face Patches in IT

*Winrich Freiwald*

*Doris Tsao*

ML

*Nancy Kanwisher*

**fMRI**

Faces vs Objects



Tsao, Freiwald, and Livingstone used fMRI to reveal a set of face selective regions in macaque IT (aka "face patches")

Most of the single neurons in these regions showed a preference for frontal faces

Tsao et al., *Science* 2006

Multi-array electrophysiology in macaque V4 and IT.



V4

10mm

IT

☐ = Array

About 300 total sites

Ha Hong

Jim DiCarlo

# Multi-array Electrophysiology Experiment

5760 images

64 objects

8 categories

uncorrelated photo backgrounds

Low variation

... *640 images*

Medium variation

... *2560 images*

High variation

... *2560 images*

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

Pose, position, scale, and background variation

complex, uncorrelated backgrounds **prevent low-level cheating**

part of what we mean by "complex task"

**Ellie.** *C. Shay & K. Kar (Winter 2019)*

complex, uncorrelated backgrounds **prevent low-level cheating**

part of what we mean by "complex task"

# Multi-array Electrophysiology Experiment



10mm

V4

IT

□ = Array

About 300 total sites

*Output = Binned spike counts 70ms-170ms post stimulus presentation averaged over 25-50 reps of each image.*

Img 1     Img 2     Img **5760**

...

Neuron 1
Neuron 2
Neuron 3
⋮
Neuron **296**

...

# Neural-Behavior Decoding

# Neural-Behavior Decoding



linear combination of units

# Neural-Behavior Decoding

V4 loses out at higher variation:



Basic
categorization



Performance
*(% correct)*

V4 NEURONS

Low Variation          Medium Variation          High Variation
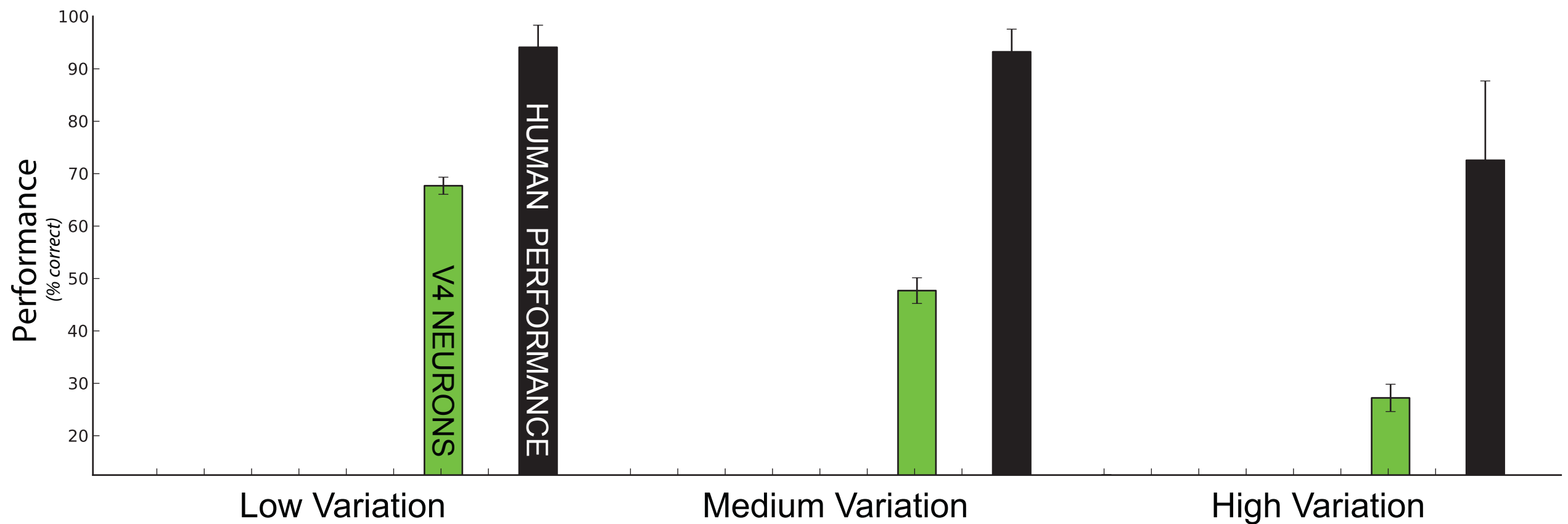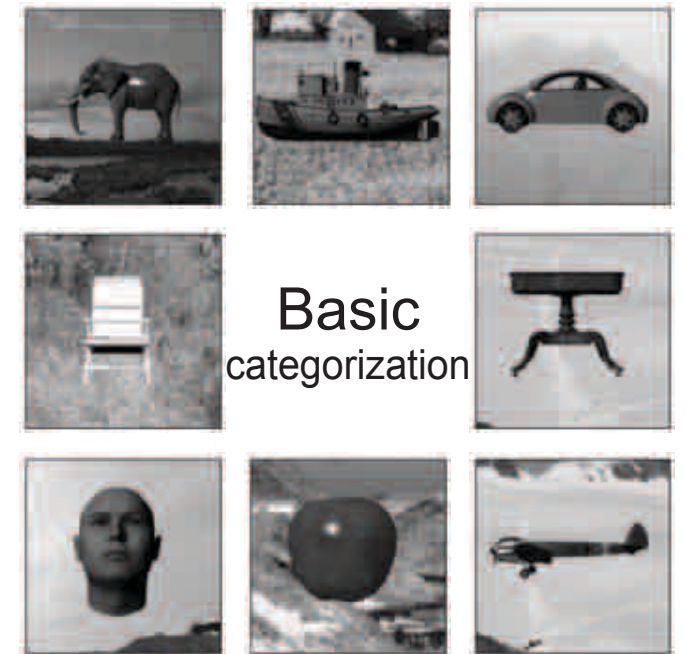
Variation Level

Low    Medium    High

at
ceiling …

… at
chance

Decoding Behaviorally Output from Neural Populations

V4 loses out at higher variation:

… but humans are much less affected.

Basic categorization

Performance (% correct)

V4 NEURONS

HUMAN PERFORMANCE

Low Variation

Medium Variation

High Variation

Yamins* and Hong* et. al. **PNAS** (2014)

V4 loses out at higher variation:

… but humans are much less affected.

… as is the IT neural population.



Basic
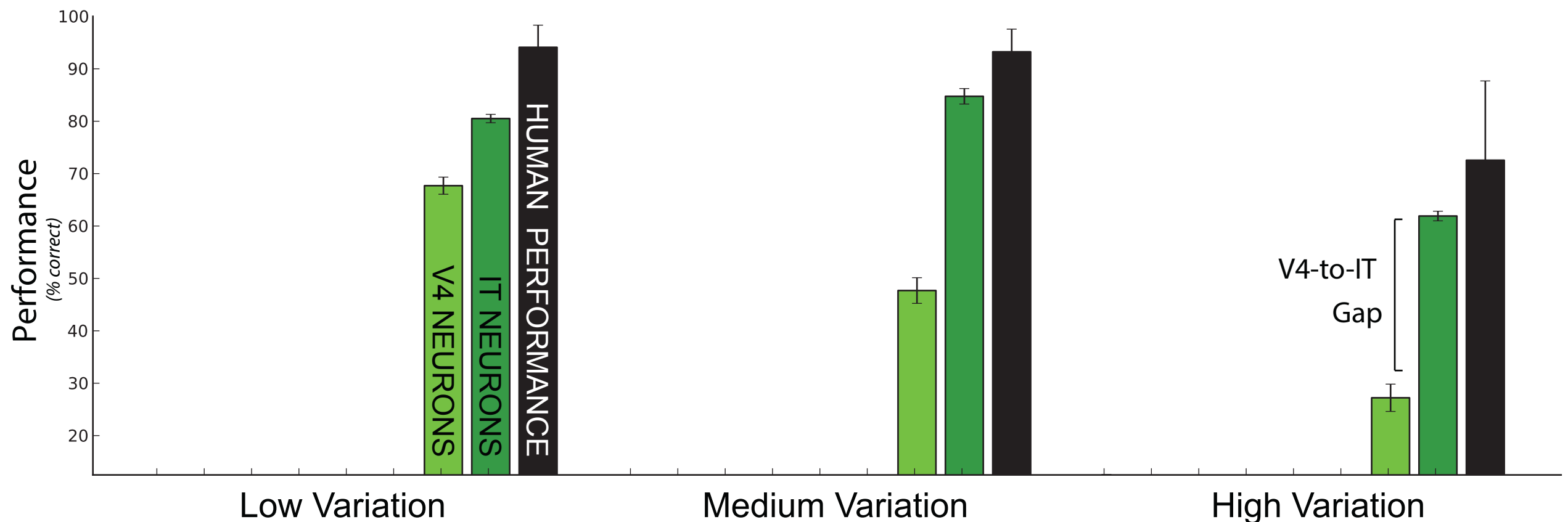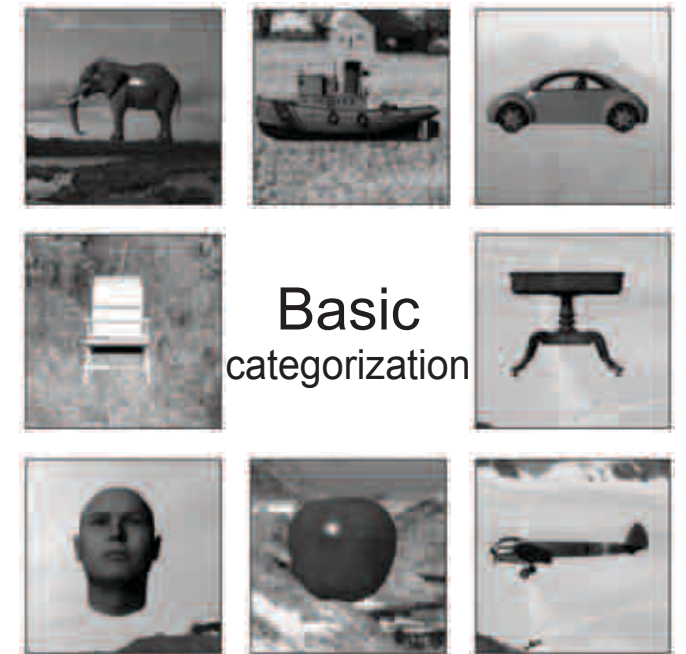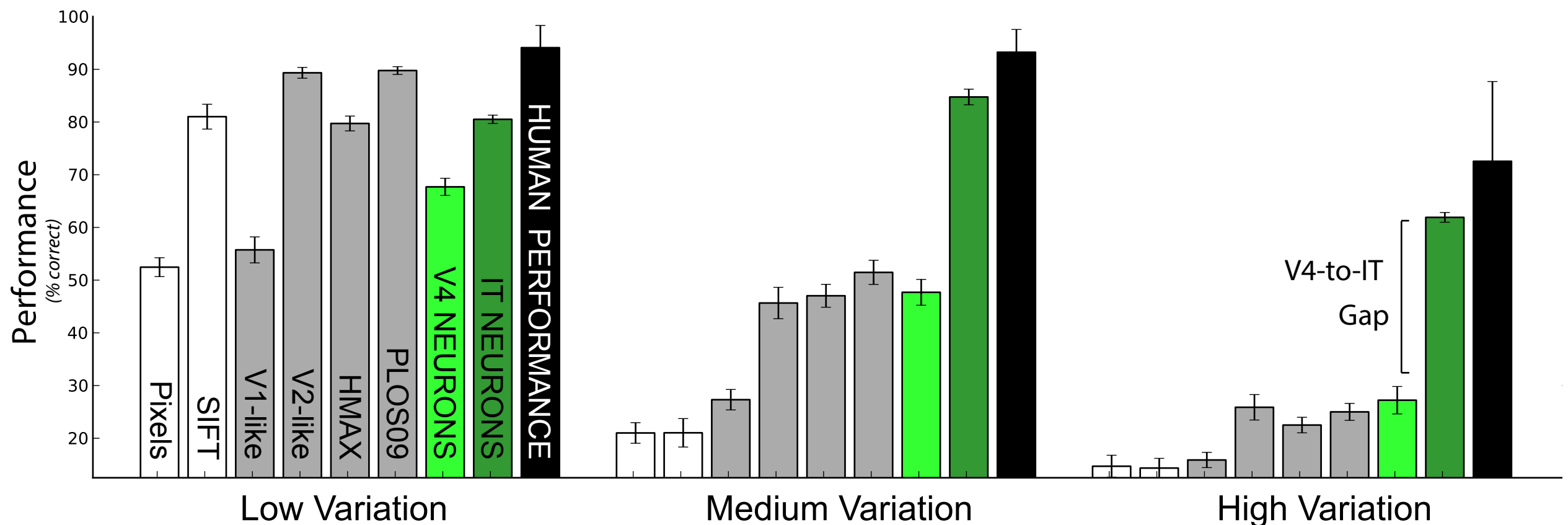categorization
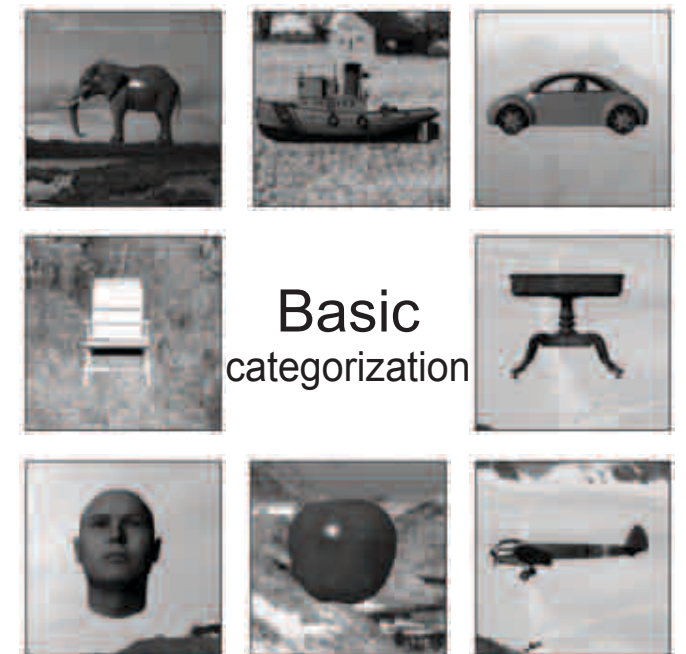


*Yamins\* and Hong\* et. al.* **PNAS** *(2014)*

# IT Neurons Track Human Performance

V4 loses out at higher variation:

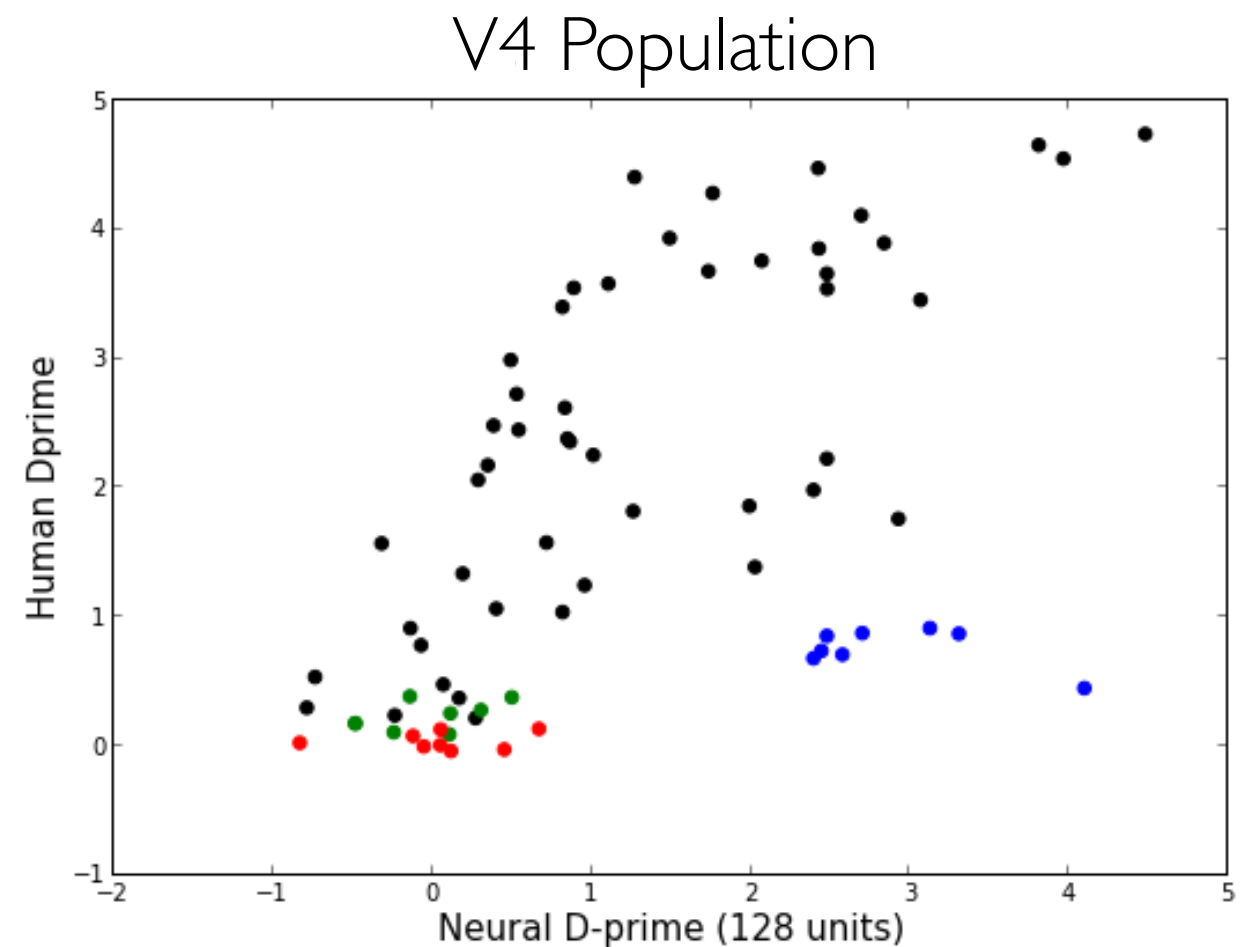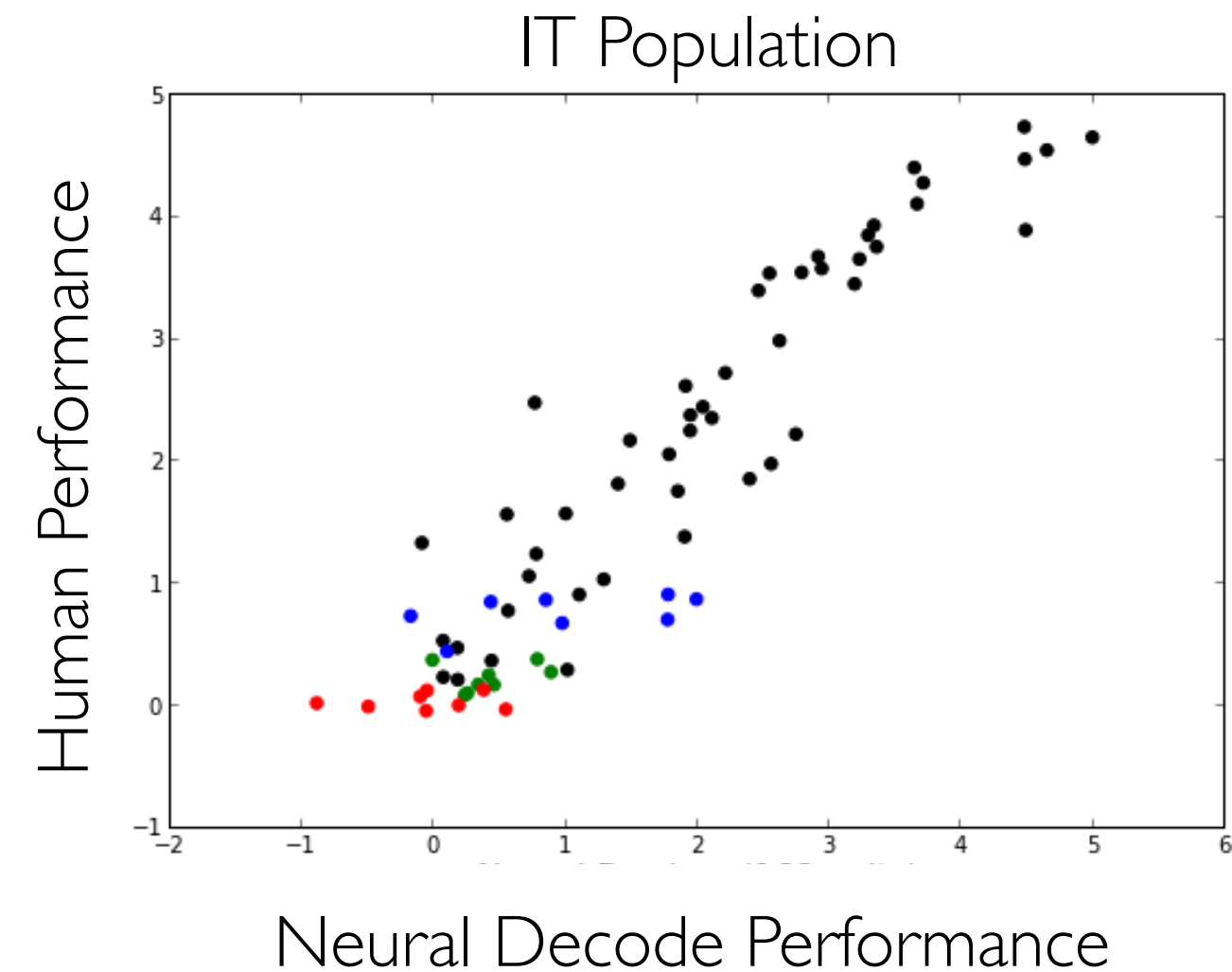… but humans are much less affected.

… as is the IT neural population.


Basic categorization



At <u>high variation levels</u>, IT much better than V4 and existing models.

*Yamins\* and Hong\* et. al.* **PNAS** *(2014)*

# IT Neurons Track Human Performance

IT matches human error patterns as well as raw performance.

IT Population



V4 Population

Human Performance

Neural Decode Performance

Human Dprime

Neural D-prime (128 units)

● Low-Variation Face subordinate tasks.

**Human**

**Rhesus monkey**

"camel" confused with "dog"

"tank" confused with "truck"

**Upshot:  human and non-human primate basic level core object percept (sp. identification) are indistinguishable**

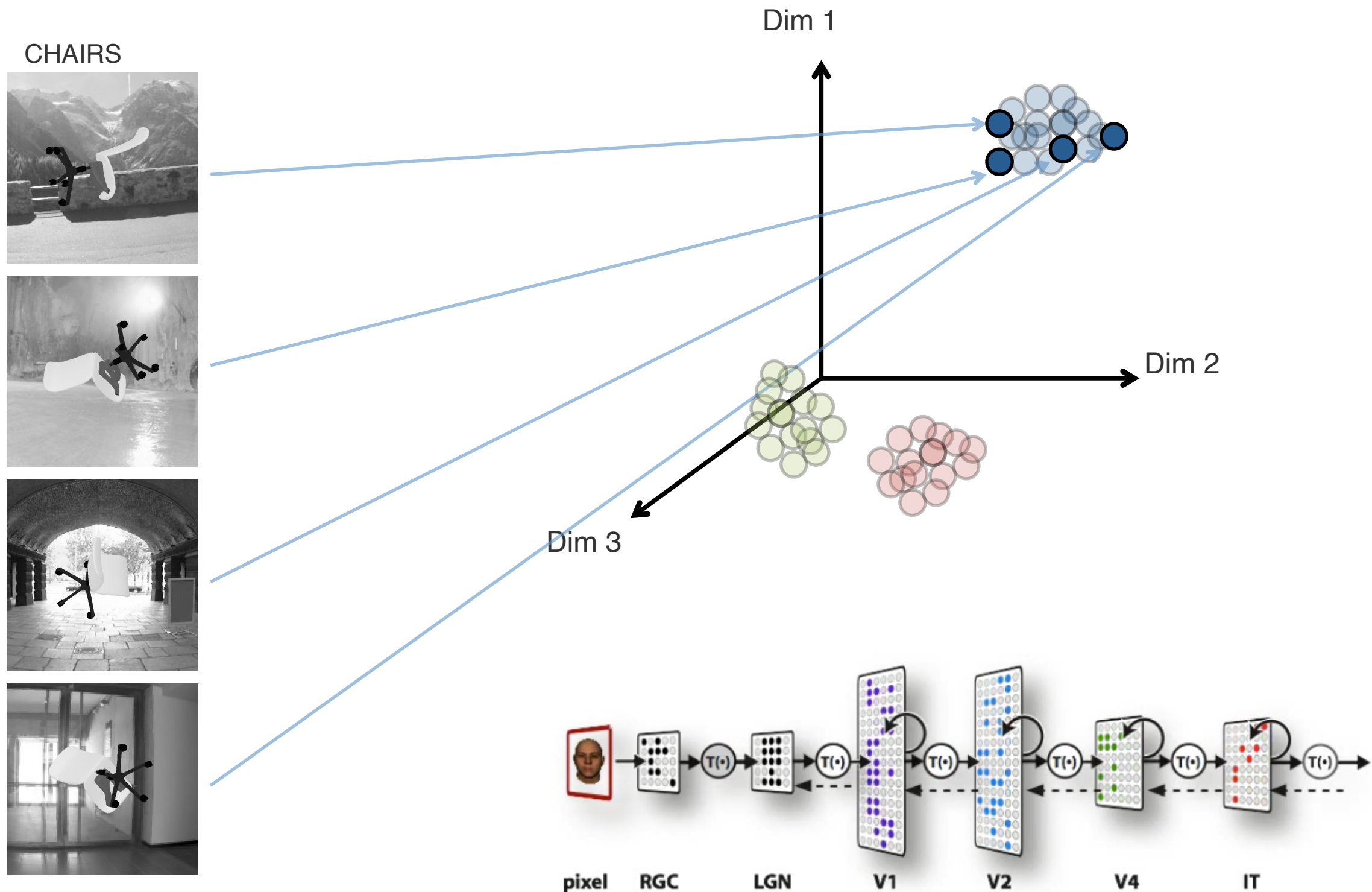**Does not depend on reporting effector (touch vs. eye movement)**

Adapted from Motter and Mountcastle 1981

Pixel space: $R^{\sim 1000000}$

Feature space: $R^{4000(?)}$

Dim 1

Dim 2

Dim 3

CHAIRS



pixel   RGC   LGN   V1   V2   V4   IT
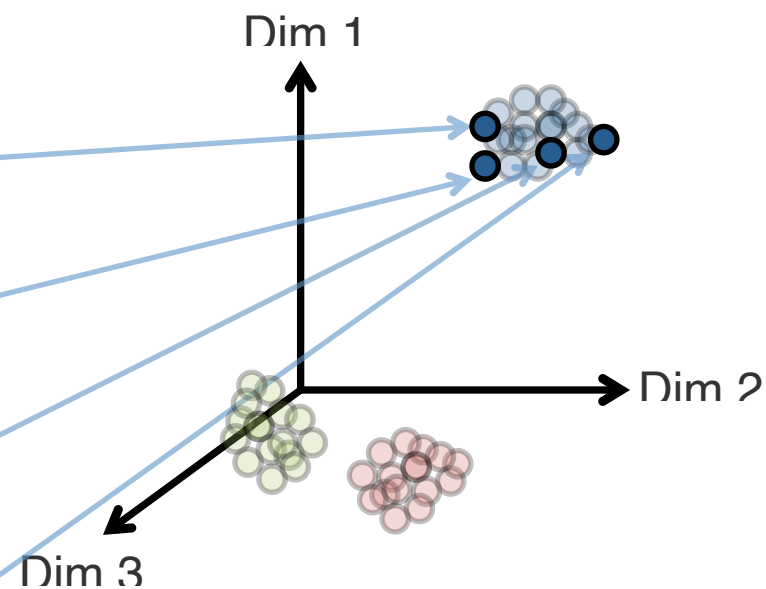
# Feature Space as Encoding

Behavior = Feature space + Simple decision rule
= encoding + decoding

Pixel space: $R^{1000000}$
Output

Feature space: $R^{4000(?)}$

Behavioral

CHAIR

Dim 1

Dim 2

Dim 3

Linear Classifier

Category Judgement

Linear Regressor

Localization

Distance Function

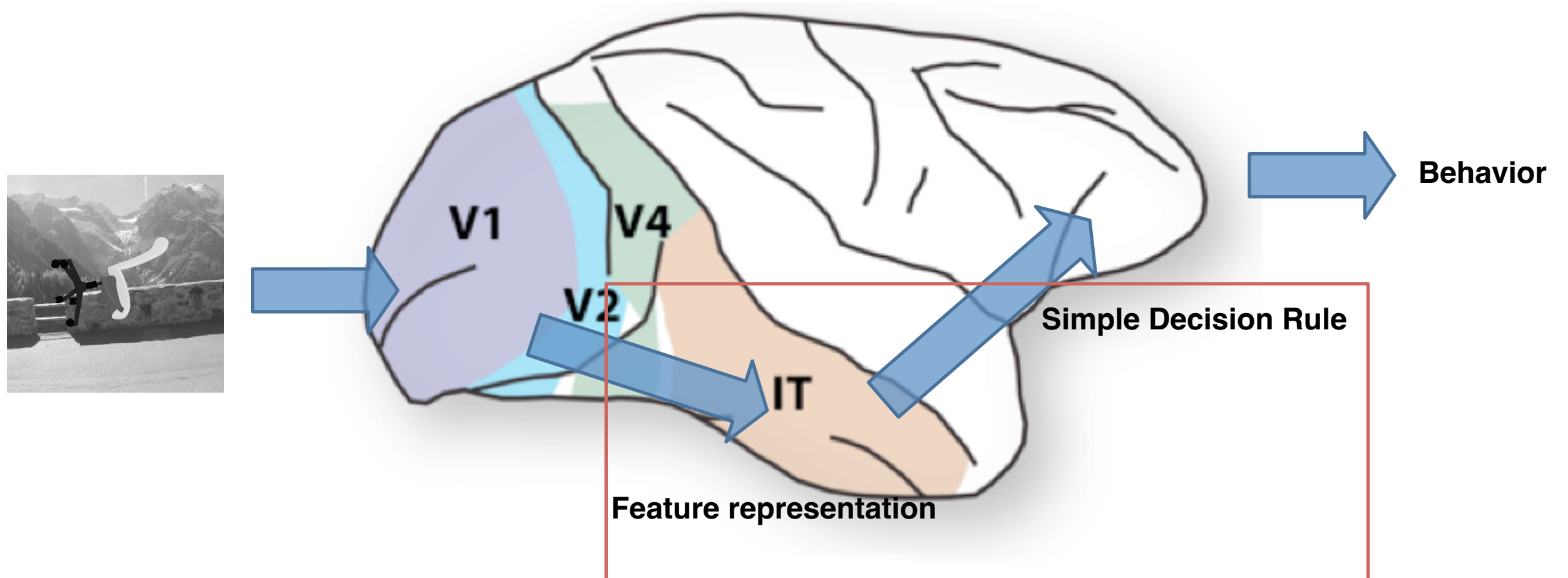"Subjective" Similarity judgement
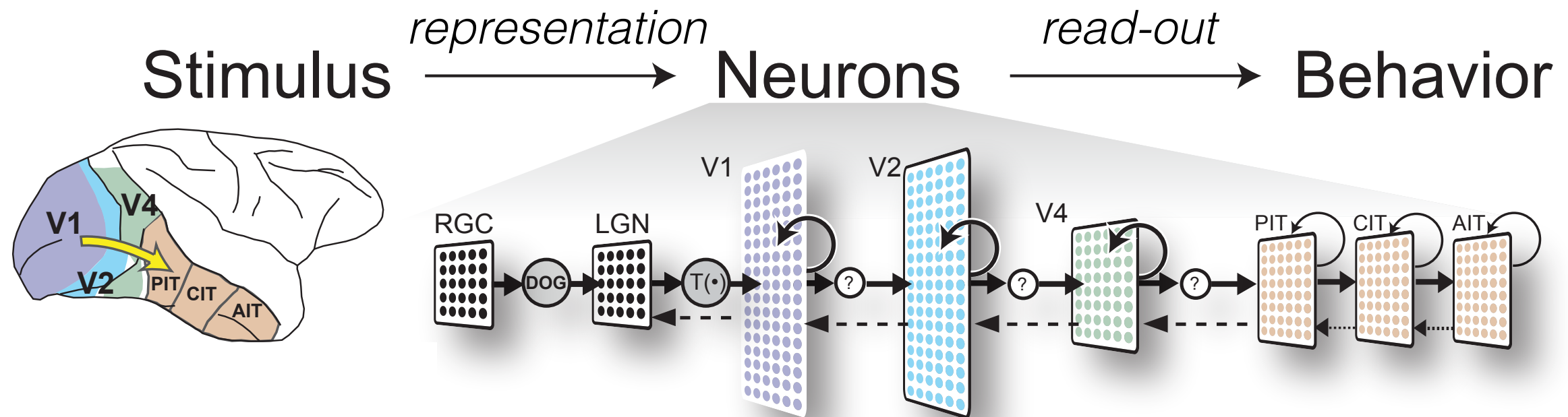
Behavior = Feature space + Simple decision rule
= encoding + decoding

# Encoding & Decoding

Stimulus → *representation* → Neurons → *read-out* → Behavior

visual representation

Category
Location
Size
Pose
Depth relationships

*very nonlinear**

**which is presumably why so much brainmeat needs to devoted to it.*
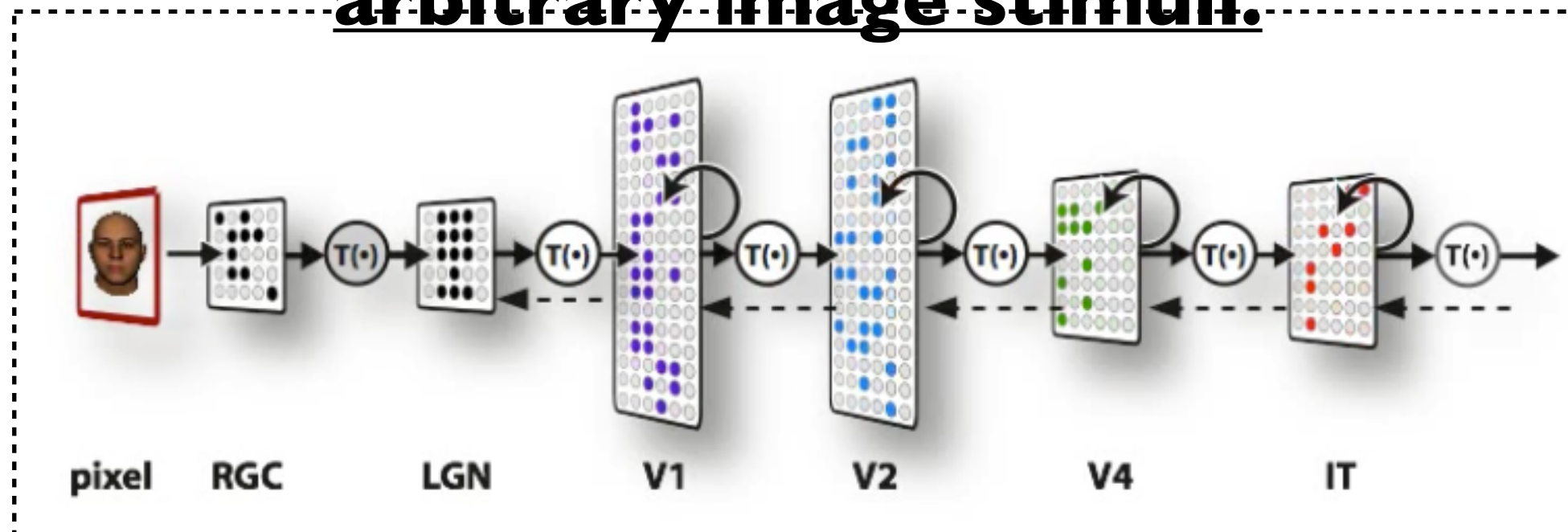
# GOAL: Predictive model of single-neuron responses throughout the ventral stream to arbitrary image stimuli.



Key questions:
    (a) how many layers?

# GOAL:  Predictive model of single-neuron responses throughout the ventral stream to arbitrary image stimuli.



pixel   RGC        LGN        V1         V2         V4         IT

Key questions:
   (a)  how many layers?
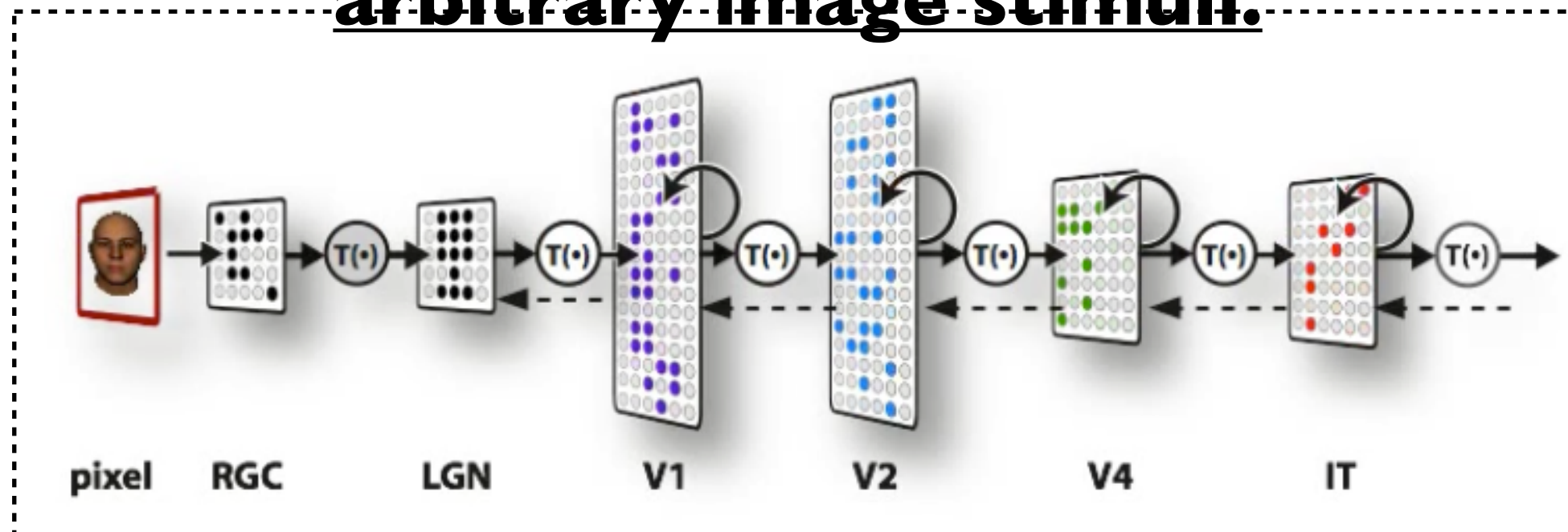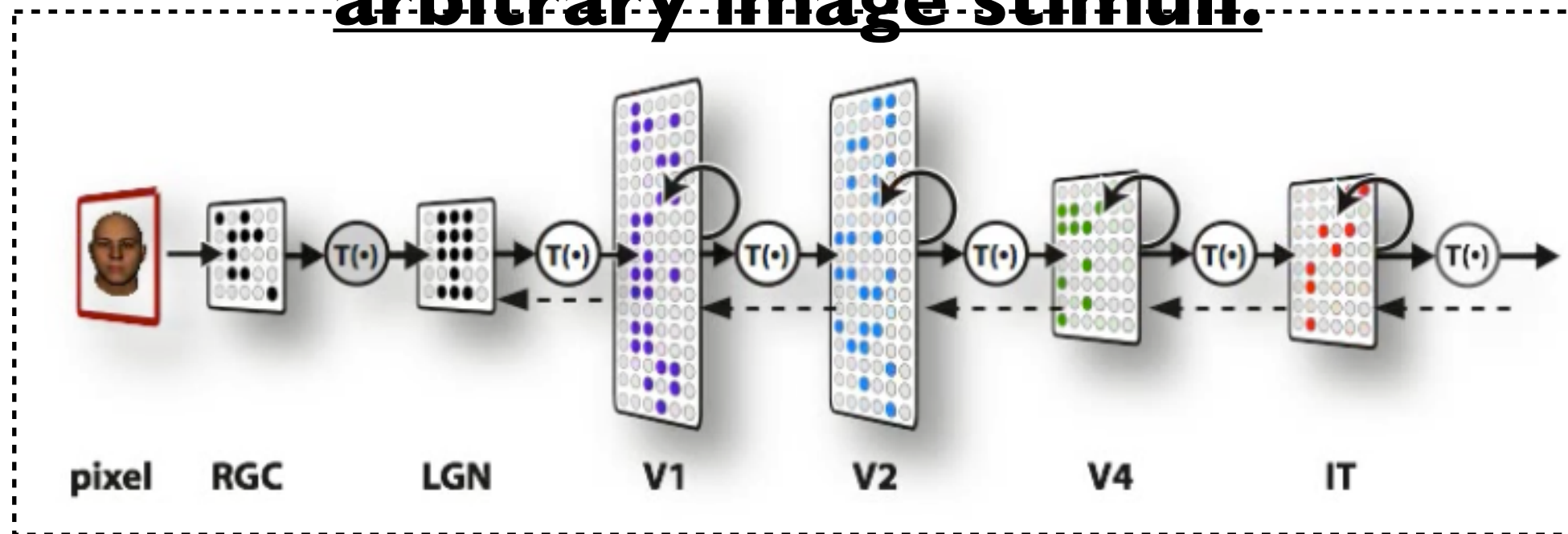
   (b)  what's in each layer, specifically?

# GOAL: Predictive model of single-neuron responses throughout the ventral stream to arbitrary image stimuli.



Key questions:
   (a) how many layers?

   (b) what's in each layer, specifically?

   (c) what behavioral goals and biophysical facts constrain it to be as it is?

**How are we supposed to use all this hard-won (Retina-IT) neuroscience knowledge to make an actual model?**