# Informatics 1 Cognitive Science

Lecture 15: Vision Part 4

---

Matthias Hennig

School of Informatics
University of Edinburgh
mhennig@inf.ed.ac.uk
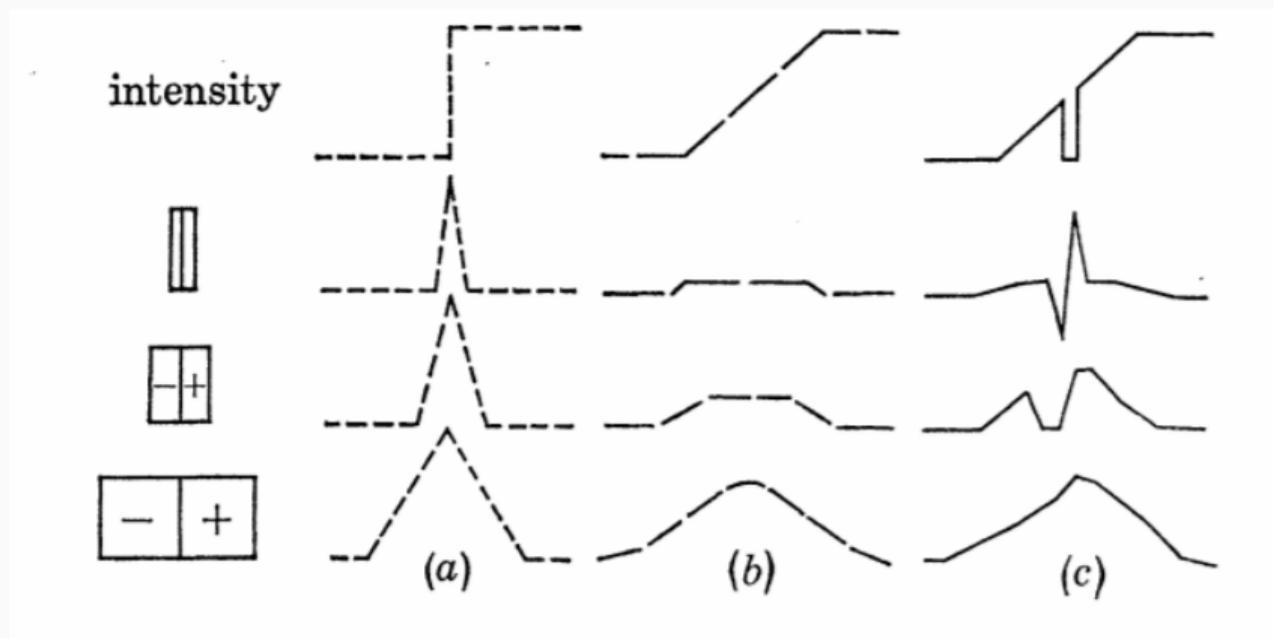
# Marr and Poggio's Vision Theory
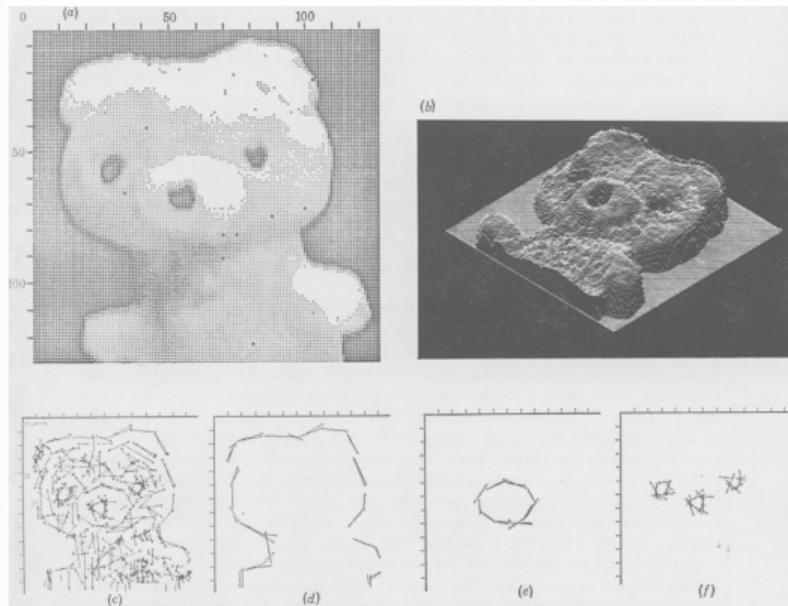
## Understanding Vision: Marr & Poggio

1. Primal sketch: local features including edges, regions, etc.
2. 2.5D sketch: surfaces with depth/orientation — shape as seen by the viewer
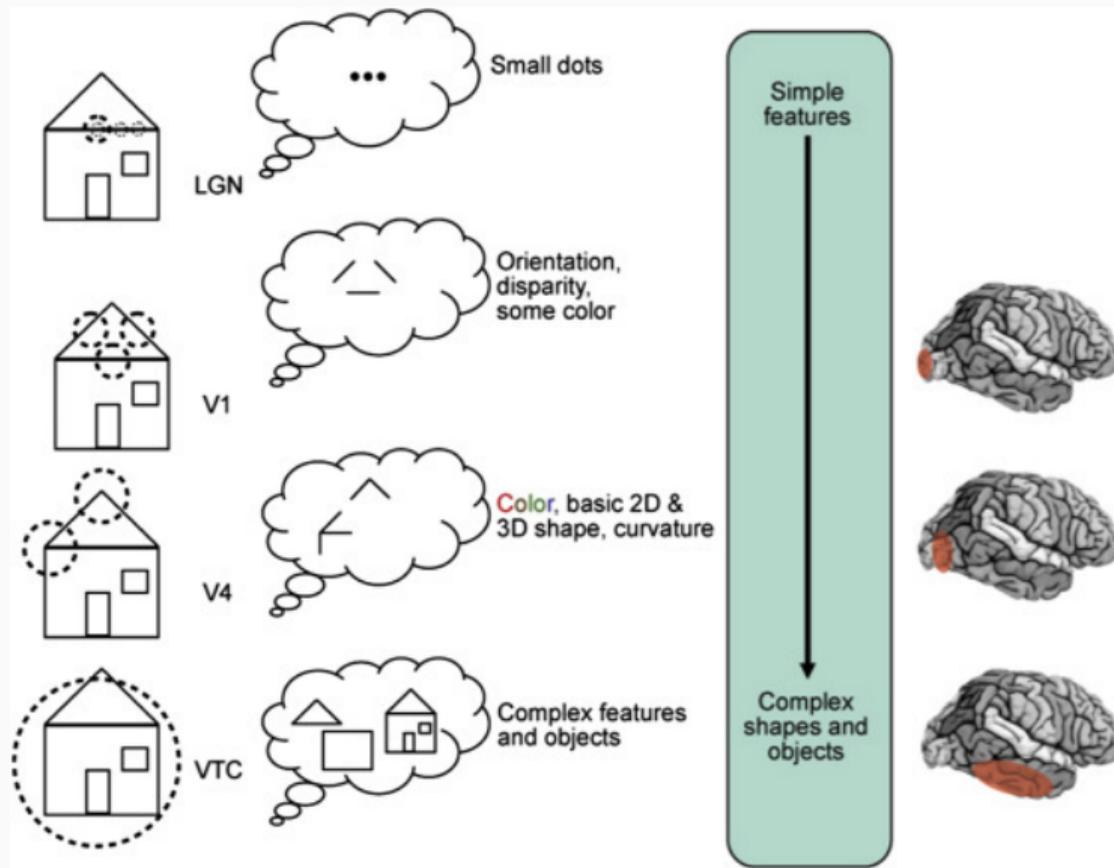3. 3 D model: represents objects in terms of 3D geometric primitives

## The Primal Sketch



Filters are required to capture contrast changes at different scales. These resemble V1 simple cells RFs.

4

The primal sketch (c), and three principal forms extracted from the primal sketch
(d-f). Marr's idea was that these primitives can now be combined to 3D object
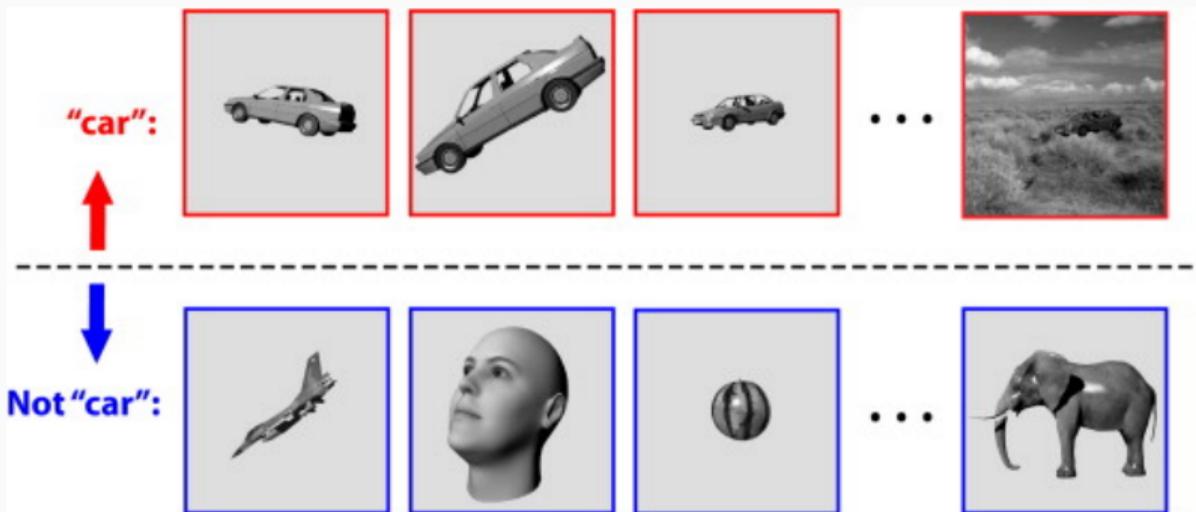descriptions. This is a hard problem that deep learning can now solve.
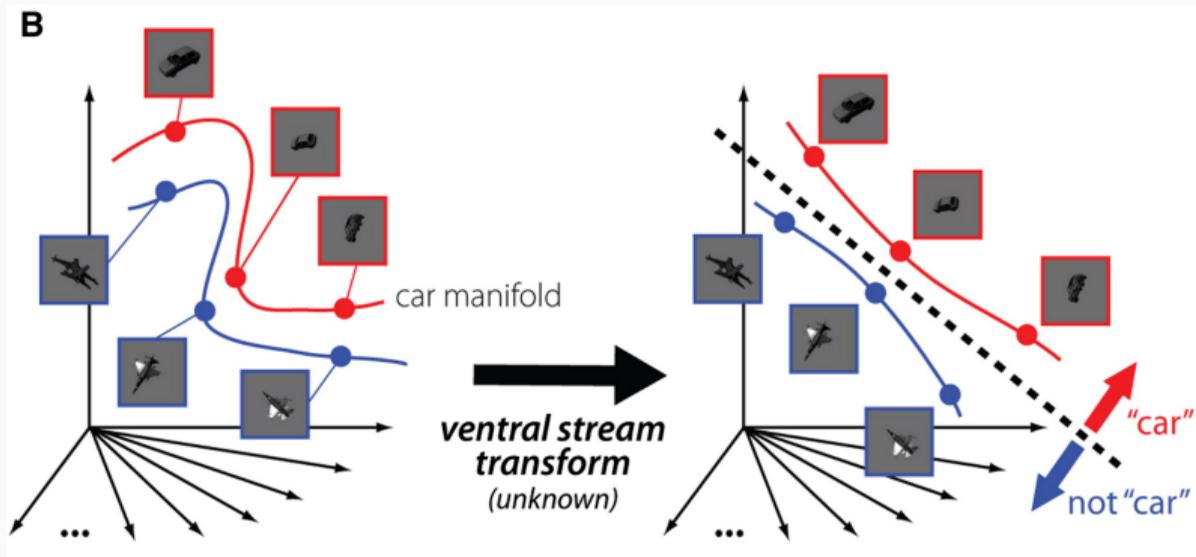
# Object Recognition in the Ventral Pathway
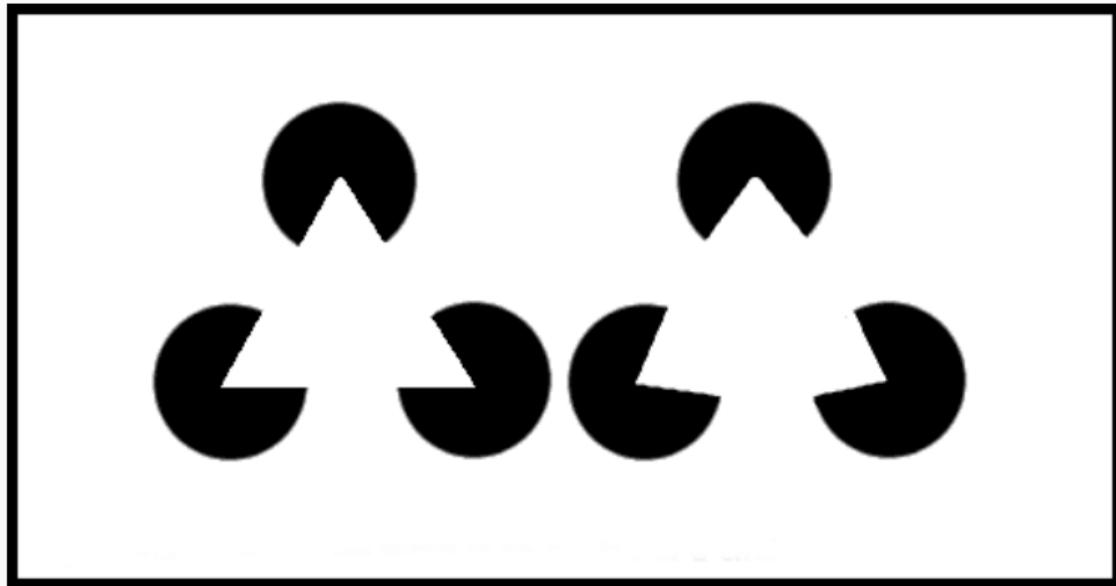
## Object Recognition



Object recognition is the ability to rapidly (200 ms viewing duration) discriminate a given visual object (e.g., a car, top row) from all other possible visual objects (e.g., bottom row) without any object-specific or location-specific pre-cuing.

**B**

car manifold

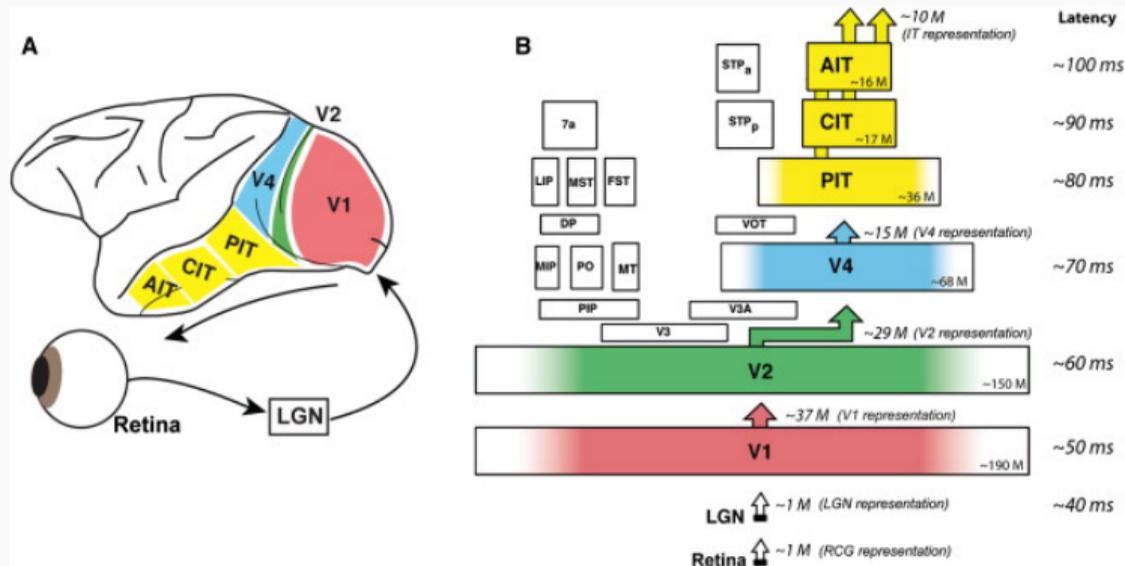*ventral stream transform (unknown)*

"car"

not "car"

In images and responses in the early visual system, object identity is hidden in curves and tangled "manifolds". Solution: a series of successive re-representations along the ventral stream to area IT that allows easy separation of the object manifold.

8

## Illusory Contours have a Neural Correlate



Responses corresponding to the non-existing lines in these images are recorded in area V2. This suggests the cortex actively interprets images according to common ecological properties.
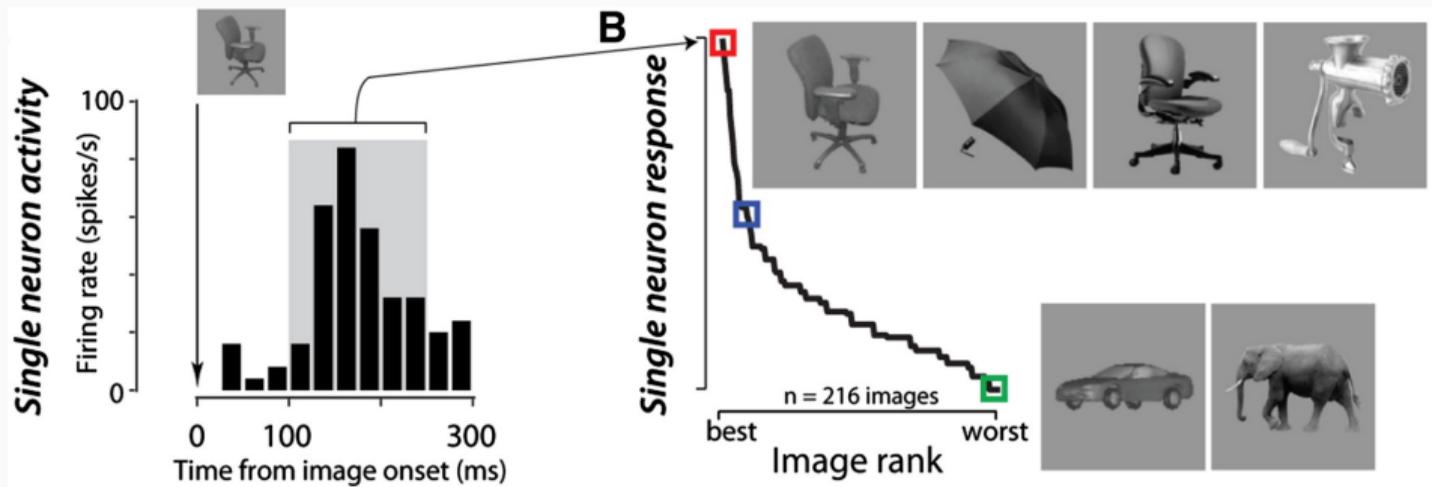
V2: Like V1 and orientation of illusory contours and figure/ground separation
V3: intermediate complexity object features, simple geometric shapes (2.5D-like), but tuning difficult to measure
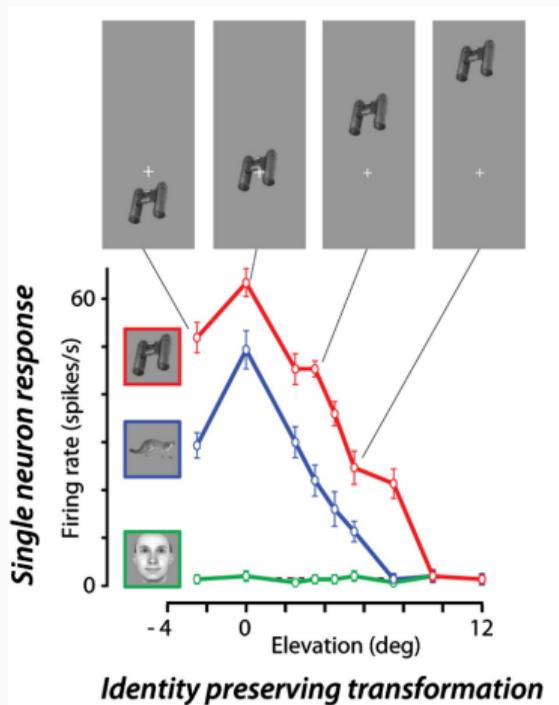Inferotemporal cortex (IT): complex shapes, objects, and faces

10

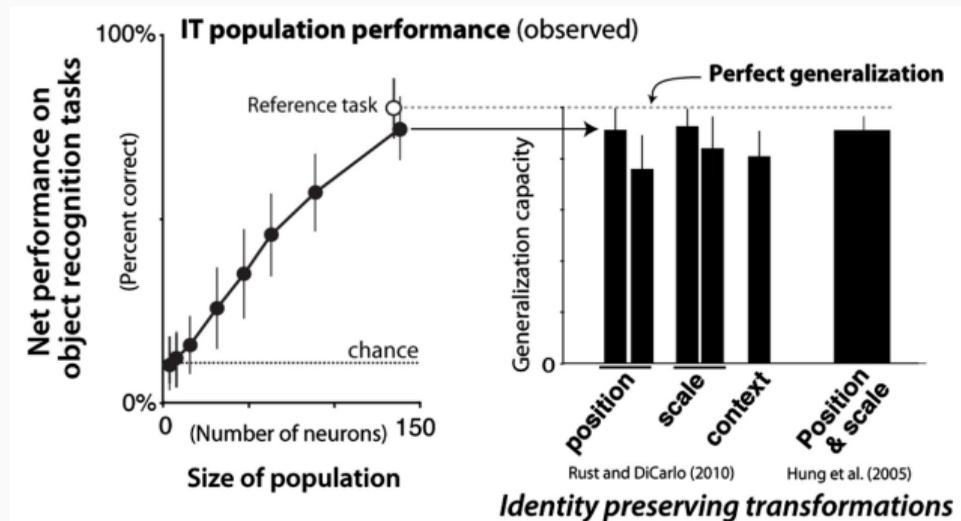IT neurons respond to pictures of objects with relatively high selectivity. (pictures from DiCarlo et al., Neuron, 2012)

Single neuron response

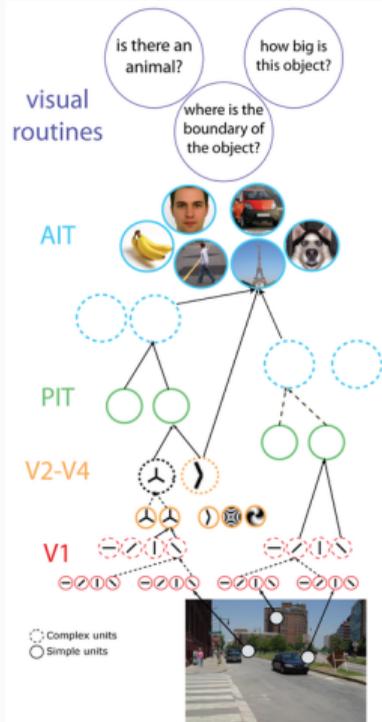*Identity preserving transformation*

Object preference is preserved over a wide range of elevations.

## Decoding Object Identity from IT Neurons



IT population performance (observed)

Rust and DiCarlo (2010)   Hung et al. (2005)

*Identity preserving transformations*

Object classification is near perfect using about 100 IT neurons, and generalisation across position and scale is robust. (reference is a based on SVM classifier on full population)
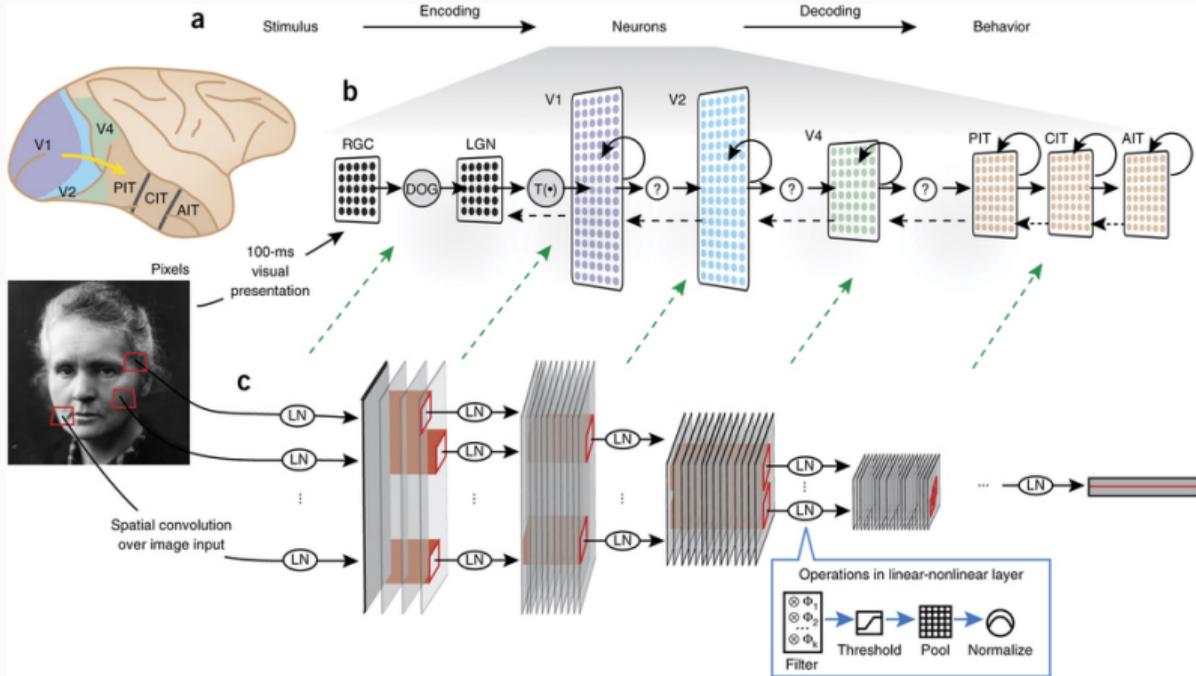
- hierarchical, local layer-wise pooling on multiple scales
- increasing size of RFs
- max pooling in higher layers
- includes learning at the top layer (and intermediate layers in newer version)
- performance ranges 50%-90% in 10 class image data sets
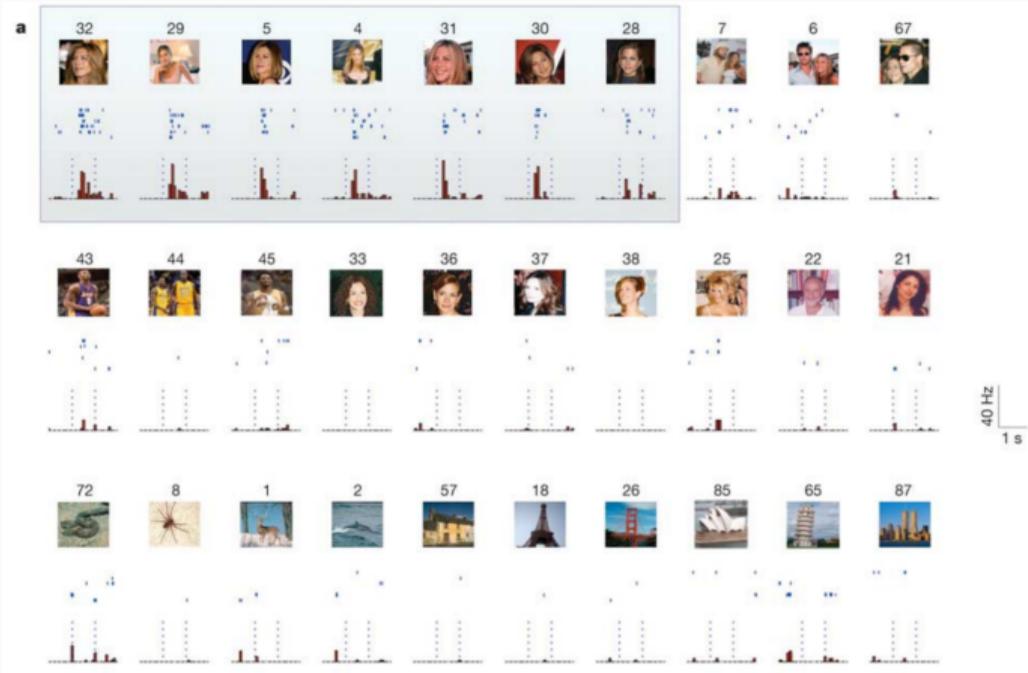
# Deep Neural Networks and the Ventral Pathway

Activations in a deep net trained to classify images mirror recorded activity in the ventral stream, and its hierarchical organisation (Yamis, DiCarlo 2012, 2016).

# Concept Cells

# Jennifer Aniston or Grandmother Cells



A single unit in the hippocampus that responds selectively to images ($+$ e.g. written or spoken name) of Jennifer Aniston (Quiroga et al., 2005).

# CLIP models also have concept cells



CLIP model: trained jointly on text and images
Paper: https://distill.pub/2021/multimodal-neurons/

Stroop effect:

green, blue, red

## Summary

- Marr & Poggio: primal sketch $\rightarrow$ 2.5D sketch $\rightarrow$ 3D model.

- The ventral stream supports object recognition via hierarchical increase in selectivity + invariance (V1/V2 $\rightarrow$ IT).

- Higher areas (e.g., V2) reflect interpretation, not just stimuli (e.g. illusory contours).

- Models: HMAX (pooling / increasing RF size) and deep networks.

- "Concept cells" illustrate highly selective coding; such units are also found in multimodal models (e.g. CLIP).