

## A very brief introduction of what is known about vision experimentally

### Neurons, neural circuits, and brain regions along the visual pathway

Multiple brain areas are involved in the processing of visual information. Visual inputs from the retina are transmitted to the lateral geniculate nucleus (LGN) of the thalamus and then form projections to the primary visual cortex (V1, striate cortex). Afterwards, information is transmitted towards the frontal areas to extrastriate cortex (V2, V3, and V4, middle temporal area MT/V5) and other brain regions such as inferotemporal cortex (IT), lateral intraparietal area (LIP), and frontal eye field (FEF). These areas can be divided into the ventral visual pathway (“what”: recognition; V2, V4, IT) and the dorsal visual pathway (“where”: motion, speed; V3, MT, LIP).

Each neuron typically responds to a limited region of the visual space, known as the receptive field. The size of the receptive fields (measured in visual angle) increase along the visual pathway. Visual inputs that excite retinal neurons are dots in contrasting backgrounds. Neurons in V1, on the other hand, are excited by static or moving bars and edges of varying orientations. Further along the visual pathway, receptive fields become more complex and depend on the animal’s conscious state and what the animal is paying attention to.

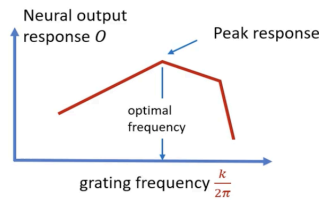
### Retina

The signals from the photoreceptors are processed by horizontal, bipolar, and amacrine cells. They are then transmitted to the retinal ganglion cells that project to the central brain through the optic nerve. The response of a typical ganglion cell is approximated as a weighted summation of photoreceptor signals. The weights are determined by a filter  $K(x)$  that is non-zero within the receptive field of the ganglion cell. This filter is a center-surround receptive field and is modeled as a difference between two Gaussians:

$$K(x) = \frac{\omega_c}{\sigma_c} \exp\left(-\frac{x^2}{2\sigma_c^2}\right) - \frac{\omega_s}{\sigma_s} \exp\left(-\frac{x^2}{2\sigma_s^2}\right)$$

Receptive field mapping is done by shining light at different locations of the receptive field (S. W. Kuffler, 1953), by using disks of light centered at the receptive field (D. H. Hubel, T. N. Wiesel, 1959), or by using grating patterns with varying frequency. A grating can be described as a cosine wave with a spatial frequency  $k$ , grating amplitude  $S_k$ , and spatial phase  $\phi$ :

$$S(x) = S_k \cos(kx + \phi) + \text{constant}$$



**Fig. 1. Contrast sensitivity function of a retinal ganglion cell.**

The neural response as a function of the grating frequency  $k$  is called the contrast sensitivity function  $g(k)$  (Fig. 1). It peaks at the frequency where the wavelength is similar to the size of the receptive field.

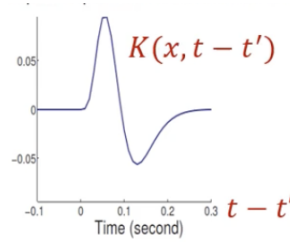
The receptive field can be approximated as a weighted summation Fourier waves. The amplitude of the waves scales with the contrast sensitivity function  $g_c(k)$ . Then, the spatial shape of the receptive field  $K(x)$  can be constructed as a weighted summation of cosine waves by measuring the contrast sensitivity function  $g_c(k)$  through neuron's responses to perfectly aligned gratings:

$$O = \int dx K(x) S_k \cos(kx + \phi) \propto g_c(k) S_k \cos(\phi)$$

The spatial input  $S(x)$  and spatial filter  $K(x)$  can be generalized to spatiotemporal input  $S(x, t)$  and spatiotemporal filter  $K(x, t - t')$  so that the neural response evolves in time:

$$O(t) = \int dt' dx K(x, t - t') S(x, t') + \text{spontaneous firing rate}$$

The neuronal response right after stimulus onset (small  $t > 0$ ) is called the transient response; otherwise, it is sustained response ( $t \rightarrow \infty$ ). The impulse response function describes a neuron's reaction to a brief input signal (Fig. 2).



**Fig. 2. Impulse response function of a model neuron to a brief input signal.** The neuron initially becomes excited; the excitation is followed by a period of inhibition and recovery to the baseline.

The center-surround receptive field model can be expanded to include time as well:

$$K(x, t) = \frac{K_c(t) \omega_c}{\sigma_c^2} \exp\left(-\frac{|x|^2}{2\sigma_c^2}\right) - \frac{K_s(t) \omega_s}{\sigma_s^2} \exp\left(-\frac{|x|^2}{2\sigma_s^2}\right)$$

The filters  $K_s(t)$  and  $K_c(t)$  have the same shape as the impulse response function (Fig. 4). However, the surround filter  $K_s(t)$  has a longer temporal width than  $K_c(t)$ . Such filters are sensitive to both spatial and temporal contrasts in visual inputs.

The retina has three types of cones: red (long wavelength, L cones), green (medium wavelength, M cones), and blue (short wavelength, S cones) cones. The response of a retinal ganglion cell is the summation of the contributions of different cones. Cones are more densely packed around the fovea and their density reduces with eccentricity. The sizes of the receptive fields increase with eccentricity. Another

way to represent the 3-dimensional color inputs is to use luminance, red-green, and blue-yellow channels. The signals in these channels are decorrelated from each other.

The two main types of cells in the retina are the parvocellular (P) cells and the magnocellular (M) cells. The P cells are more common, and have better spatial resolution but worse temporal resolution. M cells display non-linear responses to visual inputs, are more sensitive to luminance inputs and not sensitive to color differences.

### Primary visual cortex (V1)

Visual inputs from the retinal ganglion cells travel to V1 through LGN. V1 of different hemispheres receives inputs from half of the visual field. A larger area of V1 is dedicated for processing visual information from the fovea, and the devoted V1 area decreases drastically as the eccentricity increases (cortical magnification factor). The receptive fields in V1 resemble multiple retinal ganglion receptive fields added together, and are tuned to (tilted) bars and edges. They are modeled using Gabor filters:

$$K(x, y) \propto \exp\left(-\frac{x^2}{2\sigma_x^2} - \frac{y^2}{2\sigma_y^2}\right) \cos(\hat{k}x + \phi)$$

The preferred frequencies are in the Gaussian envelope around the center frequency  $\hat{k}$ . The range of the preferred frequencies is inversely related to the  $\sigma_x$ . Cells tuned to different frequency ranges are needed to cover the whole frequency range at a single spatial location (multiscale coding). Sensitivity to chromatic inputs is focused on lower spatial frequencies than sensitivity to luminance. Thus, many color-selective neurons are not tuned to orientation and have larger receptive fields. Spatiotemporal receptive fields are sensitive to both contrast and time, and are able to detect motion. These filters can be space-time separable, and space-time non-separable. A space-time non-separable can be created by combining two or more space-time separable filters and are selective to motion direction.

Binocular disparity refers to the difference in image location of an object seen by the left and right eyes. It signals visual inputs at different depths in a 3-dimensional visual world. Monocular neurons prefer inputs from one eye over another eye (ocular dominance) and binocular neurons are equally sensitive to inputs from both eyes.

V1 neurons belong to one of the two classes: either simple cells, or complex cells. Simple cells have static nonlinearity and can be approximated as linear. They prefer features such as orientation, scale, color, motion direction, disparity, and spatial location of visual inputs. Simple cells can be on-cells and off-cells and detect opposite features. These on- and off-cells are combined to approximate a linear filter. A complex cell receives inputs from two linear simple cells in a quadrature pair (energy model):

$$L_1 = \int dx K_1(x)S(x), L_2 = \int dx K_2(x)S(x) \rightarrow E \equiv L_1^2 + L_2^2$$

A complex cell is sensitive to relevant visual features (e.g., orientation). However, it is insensitive to small changes in spatial location and the contrast polarity (e.g., not sensitive to the phase of the grating).

Suppression to neural responses by visual inputs outside the receptive field is called surround (contextual) suppression (extra-classical receptive fields). This suppression is the strongest when the contextual input has the same vertical orientation as the visual input within the classical receptive field. Some contextual inputs can increase the response (contextual facilitation). Contextual inferences are likely caused by interactions between neighboring V1 neurons and feedback from higher visual areas.

### **Higher visual cortical areas**

V2 follows V1 along the visual pathway and has a similar size. Neural properties are similar to those in V1, however, the receptive fields are slightly larger. V2 neurons are sensitive to real and illusory contours, border ownership, and depth order between surfaces.

MT / V5 is along the dorsal visual pathway. Most neurons are tuned to motion direction and binocular disparity and their receptive field diameter is about 10 times larger than those of V1 neurons. MT responds to plaids (pattern cells) or component motions (component cells), and groups of moving dots.

V4 is in the ventral visual stream. Its neurons are tuned to similar features as V1. The diameter of receptive fields in V4 is 4-7 times larger as compared to V1. V4 lesions mildly impair simple visual perception, however, they severely impair shape discrimination.

IT in the ventral stream. Its neurons have large receptive fields that include the central fovea. They are driven by moderately complex shapes (face, bush, hand, etc.) rather unselectively. However, some patches include cells exclusively selective to face-like features or body parts.

Central vision loss is caused by macular degeneration through loss of retinal ganglion cells and results in difficulties recognizing faces, reading, etc. Tunnel vision is caused by the vision loss in the peripheral visual field (glaucoma, damage to optic nerve) and causes difficulties in navigation. Damage to V1 makes the person clinically blind. However, there are some residual visual abilities called blindsight (deny seeing an object but can navigate around it). Hemineglect is a deficit of attentional awareness to a half of the visual field and is caused by damage to one (right) hemisphere of the brain.

### **Eye movements and visual attention**

Visual acuity is poor away from the center of gaze. As a result, the gaze shifts around 3 times per second. Gaze trajectories are affected by the internal goal of the observer.

Two types of eye movements include saccades (jump between objects) and smooth pursuits (following an object). Eye movements are related to attention. For example, just before an impending saccade, object discrimination is best at the destination of the saccade as the attention is directed to that location. It is possible to keep the gaze fixed while directing attention elsewhere; however, shifting gaze without shifting attention is impossible.

Multiple brain regions are involved in eye movements. SC receives inputs from the retina and visual cortical areas and sends commands to control saccadic eye movements. The frontal eye field can bypass SC to directly command eye movements. The upper layers of SC (sensory) have a retinotopic map of the visual field with small receptive fields. In the deeper layers of SC (motor) the neurons are active before or

during saccades to their movement fields (the movement fields coincide with the receptive fields in the upper layers).

Most LIP and FEF neurons have visual receptive fields and many respond before or during saccades to their movement fields. Electrical stimulation of V1, V2, superior colliculus, LIP, and FEF can produce saccadic eye movements towards the receptive field or the movement field of the stimulated neurons. FEF stimulation still evokes saccadic eye movements if the superior colliculus is lesioned (has direct command lines to the brain stem). Lesioning SC reduces the number of spontaneous saccades and short latency saccades (express saccades) are eliminated. FEF lesions eliminate memory-guided or non-visually guided saccades (unable to move eyes in response to motor commands). V1 lesions abolish all visually-guided saccades until 2 months after the lesion which means that direct visual inputs from the retina to SC are not enough to drive saccades in monkeys.

Neurons within LIP, FEF, and SC increase their responses to visual input within their receptive fields if this input is the target of an upcoming saccade. Moreover, they shift their visual receptive fields so that they respond to visual inputs that are about to be brought into the classical receptive fields by the impending saccades, a process known as receptive field remapping.

Paying attention to a location is similar to stimulating the corresponding SC neuron. Visual discrimination is better at the receptive field of the more active neuron in LIP. Subthreshold stimulation of FEF or SC improves task performance inside the movement fields of the stimulated neurons.

Attention is measured by its effect to increase the performance or speed of object recognition at the attended location versus a non-attended location or feature. Attention can be controlled endogenously (voluntarily) or exogenously (involuntarily). Usually, the effects of endogenous attention are measured. They are observed in V2, V4, IT, MT; the effects in V1 are very weak. Attention increases sensitivity and effective input strength, scales feature tuning curve and changes feature tuning width (biased competition: visual inputs compete for being represented in the brain).

Exogenous attentional shift is caused by external visual inputs, so it is difficult to distinguish between (1) changes of external visual inputs, and (2) changes in exogenous attention as the reason for changes in the neural responses. Although top-down attentional influences on V1 responses are weak or absent, contextual influences in V1 could depend on what task a monkey is performing and on the monkey's experience.

## **Visual behavior**

A typical visual behavioral experiment involves an observer performing a visual task and their reaction is observed by pressing buttons, tracking gaze position or pupil size, measuring their brain waves. Categories of visual behavior can be viewed in terms of encoding (visual input sampling and image processing), selection (visual attention and segmentation), and decoding (recognition).

Examples of remarkable visual behavioral phenomena include: inattention blindness, random dot stereograms, visual crowding, Adelson's checker-shadow illusion, Kanizsa's triangle, hyperacuity, reverse Phi motion illusion, scintillating grid illusion, occlusion and recognition, ambiguous perception, binocular rivalry, gaze captured by a non-distinctive and task-irrelevant object.