

Understanding early visual coding from information theory

By Li Zhaoping

EU-India-China summer school, ISI, June 2007

Contact: z.li@ucl.ac.uk

Facts: neurons in early visual stages: retina, V1, have particular receptive fields. E.g., retinal ganglion cells have center surround structure, V1 cells are orientation selective, color sensitive cells have, e.g., red-center-green-surround receptive fields, some V1 cells are binocular and others monocular, etc.


Question: Can one understand, or derive, these receptive field structures from some first principles, e.g., information theory?

Example: visual input, 1000x1000 pixels, 20 images per second --- many megabytes of raw data per second.

Information bottle neck at optic nerve.

Solution (**Infomax**): recode data into a new format such that data rate is reduced without losing much information.

Redundancy between pixels.

1 byte per pixel at receptors  0.1 byte per pixel at retinal ganglion cells?

Consider redundancy and encoding of stereo signals

S^L



S^R



Redundancy is seen at correlation matrix (between two eyes)

$$R^S \equiv \begin{pmatrix} \langle S_L^2 \rangle & \langle S_L S_R \rangle \\ \langle S_R S_L \rangle & \langle S_R^2 \rangle \end{pmatrix} = \langle S_L^2 \rangle \begin{pmatrix} 1 & r \\ r & 1 \end{pmatrix}$$

$$0 \leq r \leq 1.$$

Assume signal (S_L, S_R) is gaussian, it then has probability distribution:

$$P(\mathbf{S}) \propto \exp\left(-\sum_{ij} S_i \tilde{S}_j (R^S)^{-1}_{ij} / 2\right)$$

An encoding:

$$O_+ = S_+ \equiv (S_L + S_R)/\sqrt{2}, \quad O_- = S_- \equiv (S_L - S_R)/\sqrt{2}$$

Gives zero correlation $\langle O_+ O_- \rangle$ in output signal (O_+, O_-) , leaving output Probability

$$P(O_+, O_-) = P(O_+) P(O_-) \quad \text{factorized.}$$

The transform S to O is linear.

O_+ is binocular, O_- is more monocular-like.

Note: S_+ and S_- are eigenvectors or principal components of the correlation matrix R^S , with eigenvalues $\langle S_{\pm}^2 \rangle = (1 \pm r) \langle S_L^2 \rangle$

In reality, there is input noise $N_{L,R}$ and output noise $N_{o,\pm}$, hence:

$$O_{\pm} = [(S_L + N_L) \pm (S_R + N_R)]/\sqrt{2} + N_{o,\pm},$$

Effective output noise:

$$N_{\pm} = (N_L \pm N_R)/\sqrt{2} + N_{o,\pm}.$$

Let:

$$\langle N^2 \rangle \equiv \langle N_L^2 \rangle = \langle N_R^2 \rangle, \text{ and } \langle \bar{N}_o^2 \rangle \equiv \langle N_{o,+}^2 \rangle = \langle N_{o,-}^2 \rangle.$$

Input $S_{L,R} + N_{L,R}$ has

$$I_{L,R} = \frac{1}{2} \log_2 \frac{\langle S_{L,R}^2 \rangle + \langle N^2 \rangle}{\langle N^2 \rangle}$$

Bits of information about signal $S_{L,R}$

Input $S_{L,R} + N_{L,R}$ has

$$I_{L,R} = \frac{1}{2} \log_2 \frac{\langle S_{L,R}^2 \rangle + \langle N^2 \rangle}{\langle N^2 \rangle}$$

bits of information about signal $S_{L,R}$

Whereas outputs $O_{+,-}$ has

$$I_{\pm} = \frac{1}{2} \log_2 \frac{\langle O_{\pm}^2 \rangle}{\langle N_{\pm}^2 \rangle} = \frac{1}{2} \log_2 \frac{\langle S_{\pm}^2 \rangle + \langle N^2 \rangle + \langle N_o^2 \rangle}{\langle N^2 \rangle + \langle N_o^2 \rangle}$$

bits of information about signal $S_{L,R}$

Note: redundancy between S_L and S_R cause higher and lower signal powers $\langle O_+^2 \rangle$ and $\langle O_-^2 \rangle$ in O_+ and O_- respectively, leading to higher and lower information rate I_+ and I_- .

If cost $\sim \langle O_{\pm}^2 \rangle$

$$I_{\pm} = \frac{1}{2} \log_2(\langle O_{\pm}^2 \rangle) + \text{constant} = \frac{1}{2} \log_2(\text{cost}) + \text{constant},$$

Gain in information per unit cost $(\Delta I / \Delta \text{cost})$

smaller in O_+ than in O_- channel.

If cost $\sim \langle O_{\pm}^2 \rangle$

$$I_{\pm} = \frac{1}{2} \log_2(\langle O_{\pm}^2 \rangle) + \text{constant} = \frac{1}{2} \log_2(\text{cost}) + \text{constant},$$

Gain in information per unit cost $(\Delta I / \Delta \text{cost})$

smaller in O_+ than in O_- channel.

Hence, gain control on O_{\pm} is motivated.

$$O_{\pm} \rightarrow g_{\pm} O_{\pm}$$

To balance the cost and information extraction, optimize by finding the gain g_+ such that

$$E(V_{\pm}) \equiv \sum_a (\langle O_a^2 \rangle) - \lambda \sum_a (I_a) = \text{cost} - \lambda \cdot \text{Information}$$

Is minimized. This gives, for $k = +$ or $-$

$$g_k^2 \propto \text{Max} \left\{ \left[\frac{1}{2} \frac{\langle S_k^2 \rangle}{\langle S_k^2 \rangle + \langle N^2 \rangle} \left(1 + \sqrt{1 + \frac{4\lambda}{(\ln 2) \langle N_o^2 \rangle} \frac{\langle N^2 \rangle}{\langle S_k^2 \rangle}} \right) - 1 \right], 0 \right\}$$

$$g_k^2 \propto \text{Max} \left\{ \left[\frac{1}{2} \frac{\langle S_k^2 \rangle}{\langle S_k^2 \rangle + \langle N^2 \rangle} \left(1 + \sqrt{1 + \frac{4\lambda}{(\ln 2) \langle N_o^2 \rangle} \frac{\langle N^2 \rangle}{\langle S_k^2 \rangle}} \right) - 1 \right], 0 \right\}$$

In the zero noise limit when $\frac{\langle S_{\pm}^2 \rangle}{\langle N^2 \rangle} \gg 1$,

$$g^2 \propto \langle S^2 \rangle^{-1}$$

This equalizes the output power $\langle O_+^2 \rangle \approx \langle O_-^2 \rangle$ --- whitening

When output noise N_o is negligible, output O and input $S+N$ convey similar amount of information about signal S , but uses much less output power with small gain g_{\pm}

$\langle O_+^2 \rangle \sim \langle O_-^2 \rangle$ --- whitening also means that output correlation matrix

$$R_{ab}^o = \langle O_a O_b \rangle$$

Is proportional to identity matrix, (since $\langle O_+ O_- \rangle = 0$).

Any rotation (unitary or ortho-normal transform):

$$\begin{pmatrix} O_1 \\ O_2 \end{pmatrix} = \begin{pmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{pmatrix} \begin{pmatrix} O_+ \\ O_- \end{pmatrix} = \begin{pmatrix} \cos(\theta)O_+ + \sin(\theta)O_- \\ -\sin(\theta)O_+ + \cos(\theta)O_- \end{pmatrix}.$$

Preserves de-correlation $\langle O_1 O_2 \rangle = 0$

Leaves output cost $\text{Tr}(R^o)$ unchanged

Leaves amount of information extracted $I = \frac{1}{2} \log \frac{\det R^o}{\det R^N}$, unchanged

Tr, det, denote trace and determinant of matrix.

Both encoding schemes:

$$S_{L,R} \rightarrow O_{\pm} \text{ and } S_{L,R} \rightarrow O_{1,2}$$

With former a special case of latter, are optimal in making output decorrelated (non-redundant), in extracting information from signal S, and in reducing cost.

In general, the two different outputs:

$$\begin{pmatrix} O_1 \\ O_2 \end{pmatrix} = \begin{pmatrix} S_L(\cos(\theta)V_+ + \sin(\theta)V_-) + S_R(\cos(\theta)V_+ - \sin(\theta)V_-) \\ S_L(-\sin(\theta)V_+ + \cos(\theta)V_-) + S_R(-\sin(\theta)V_+ - \cos(\theta)V_-) \end{pmatrix}$$

prefer different eyes. In particular, $\theta = 45^\circ$ gives

$$O_{1,2} \sim S_L(g_+ \mp g_-) + S_R(g_+ \pm g_-)$$

The visual cortex indeed has a whole spectrum of neural ocularity.

Summary of the coding steps:

Input: $S+N$, with signal correlation (input statistics) R^s

get eigenvectors (principal components) S' of R^s

$$S + N \longrightarrow S' + N' = K_o(S + N)$$

↙ rotation of coordinates

gain control V on each principal component

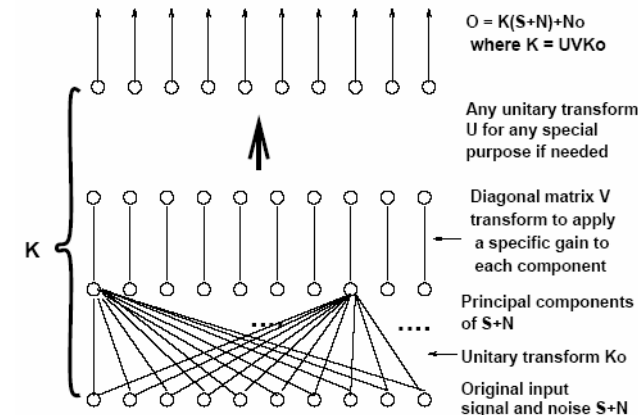
$$S' + N' \longrightarrow O = g(S' + N') + N_o$$

rotation U' (multiplexing) of O

$$O' \longrightarrow U'O = U'g K_o S + \text{noise}$$

$$\text{Neural output} = \underbrace{U'g K_o}_{\text{Receptive field, encoding kernel}} \text{ sensory input} + \text{noise}$$

Receptive field, encoding kernel



Variations in optimal coding:

Factorial codes

Minimum entropy, or minimum description length codes

Independent components analysis

Redundancy reduction

Sparse coding

Maximum entropy code

Predictive codes

Minimum predictability codes, or least mutual information between output Channels.

They are all related!!!

Another example, visual space coding, i.e., spatial receptive fields

Signal at spatial location x is $S_x = S(x)$

Signal correlation is $R^S_{x,x'} = \langle S_x S_{x'} \rangle = R^S(x-x')$ --- translation invariant

Principal components S_k are Fourier transform of S_x

$$S_x \rightarrow \underline{S}_k \sim \sum_x K_o^{kx} S_x \sim \sum_x e^{-ikx} S_x$$

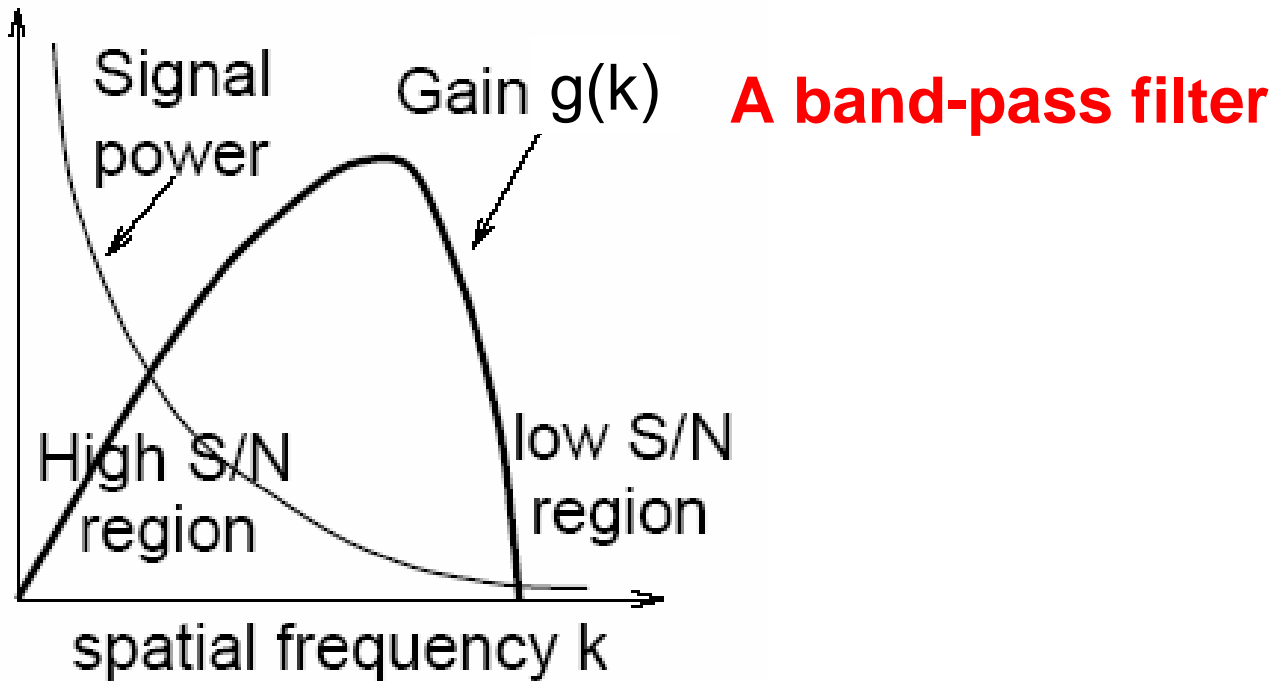
Eigenvalue spectrum (power spectrum): $\langle \underline{S}_k^2 \rangle \sim 1/k^2$

Assuming white noise power $\langle N_k^2 \rangle = \text{constant}$, high S/N region is at low frequency, i.e., small k , region.

Gain control, $V(k) \sim \langle S_k^2 \rangle^{-1/2} \sim k$, --- whitening in space

At high k , where S/N is small, $V(k)$ decays quickly with k to cut down noise according to

$$g_k^2 \propto \text{Max} \left\{ \left[\frac{1}{2} \frac{\langle S_k^2 \rangle}{\langle S_k^2 \rangle + \langle N^2 \rangle} \left(1 + \sqrt{1 + \frac{4\lambda}{(\ln 2) \langle N_o^2 \rangle} \frac{\langle N^2 \rangle}{\langle S_k^2 \rangle}} \right) - 1 \right], 0 \right\}$$

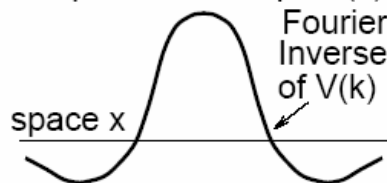


Let the multiplexing rotation be the inverse Fourier transform: $U^{x'k} \sim e^{ikx'}$

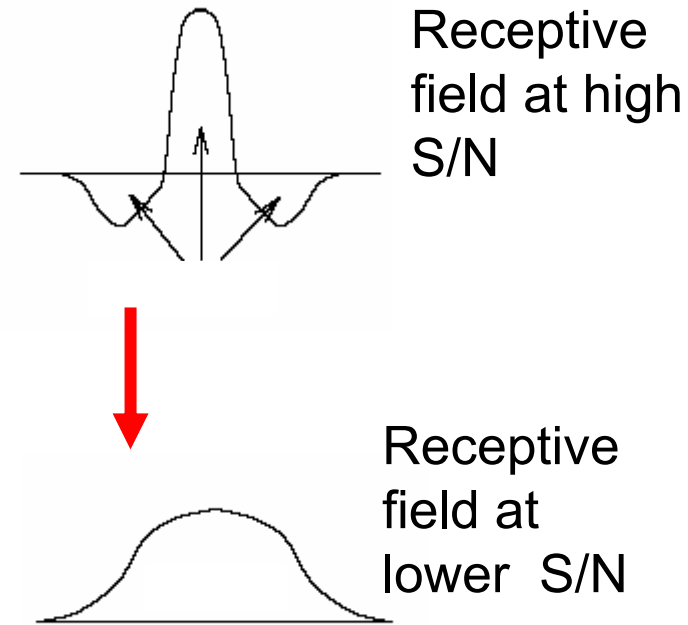
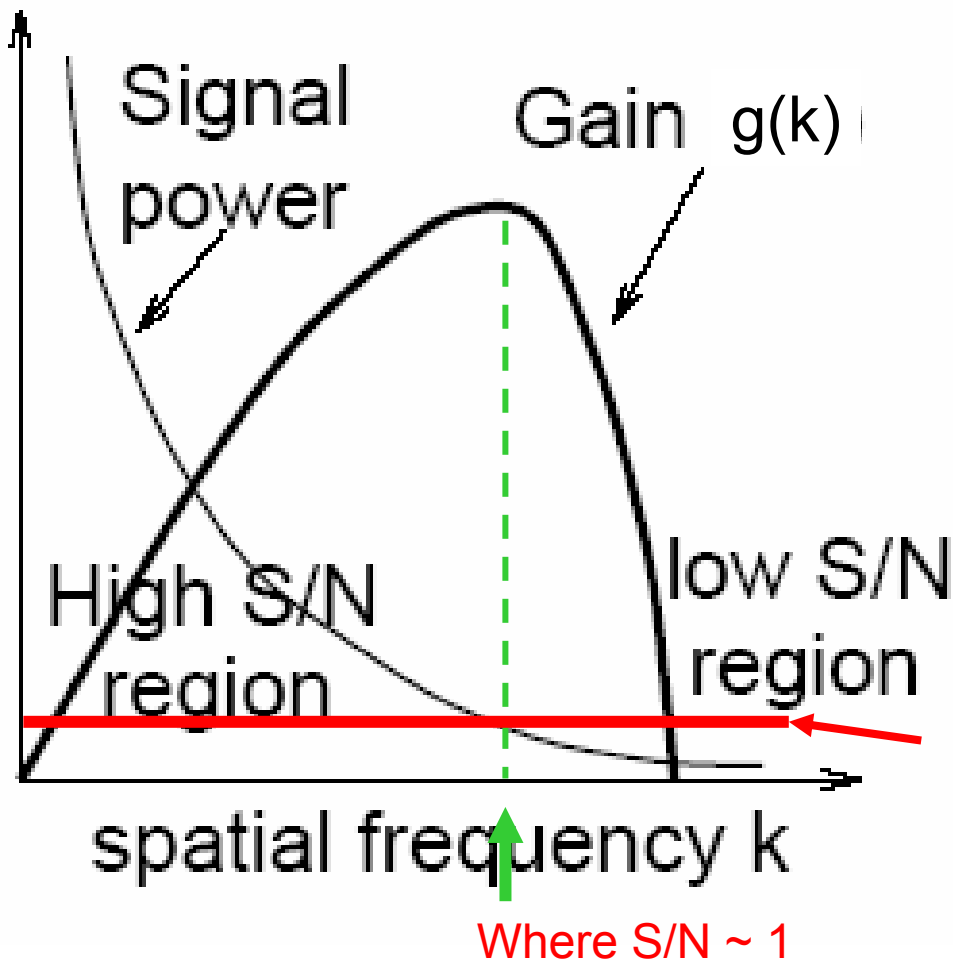
The full encoding transform is

$$O_{x'} = \sum_k U^{x'k} g(k) \sum_x e^{-kx} S_x = \sum_k g(k) \underbrace{\sum_x e^{-k(x'-x)} S_x}_{\text{Fourier Inverse of } V(k)} + \text{noise}$$

center-surround spatial filter
receptive field shape $K(x)$



Understanding adaptation by input strength



When overall input strength is lowered, the peak of $V(k)$ is lowered to lower spatial frequency k , a band-pass filter becomes a low pass (smoothing) filter.

Another example: optimal color coding

Analogous to stereo coding, but with 3 input channels, red, green, blue.

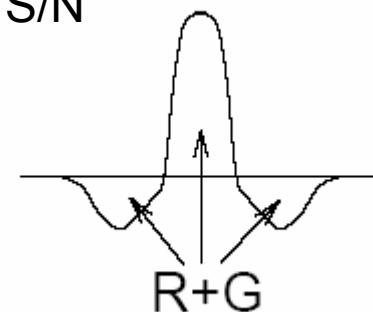
For simplicity, focus only on red and green

Input signal S_r, S_g Input correlation $R_{rg}^S > 0$

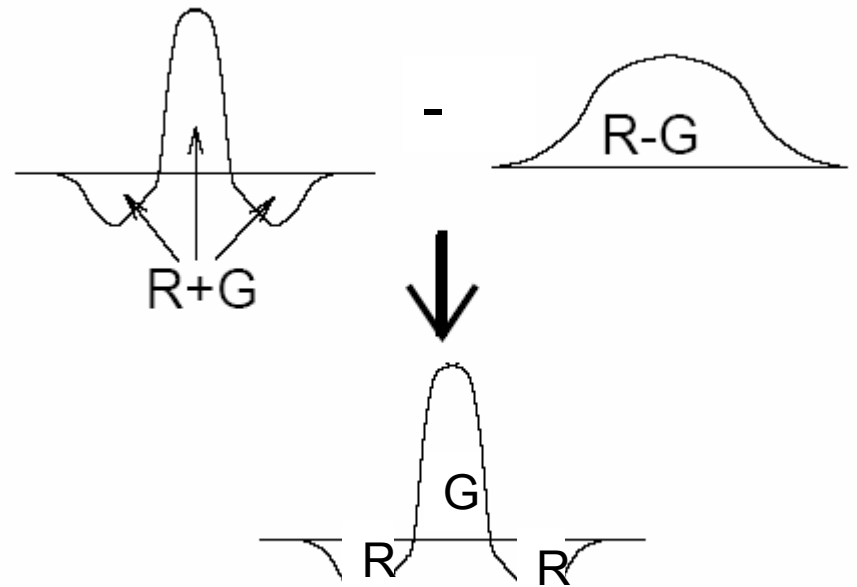
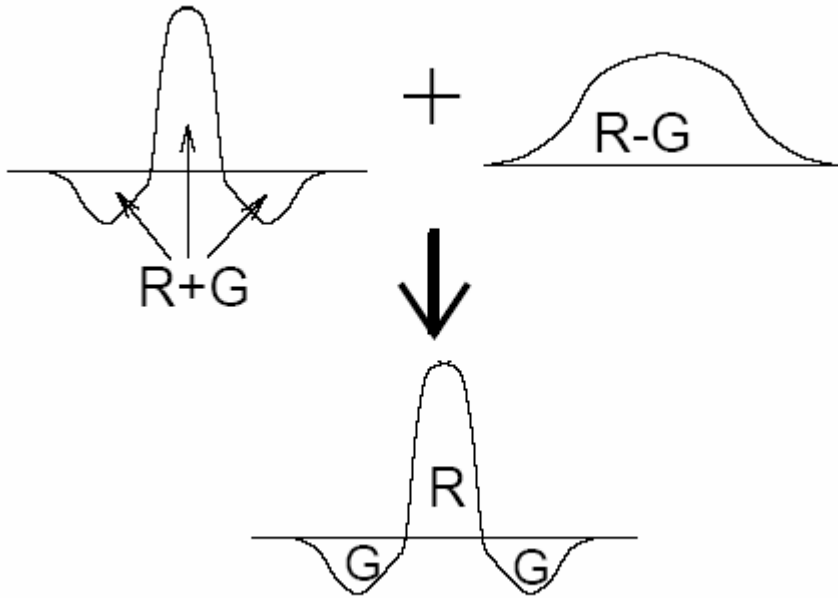
Eigenvectors: $S_r + S_g$ ← Luminance channel, higher S/N
 $S_r - S_g$ ← Chromatic channel, lower S/N

Gain control on $S_r + S_g$ --- lower gain until at higher spatial k

Gain control on $S_r - S_g$ --- higher gain then decay at higher spatial k



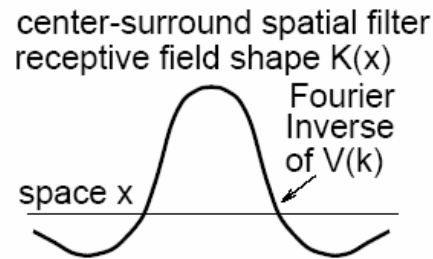
Multiplexing in the color space:



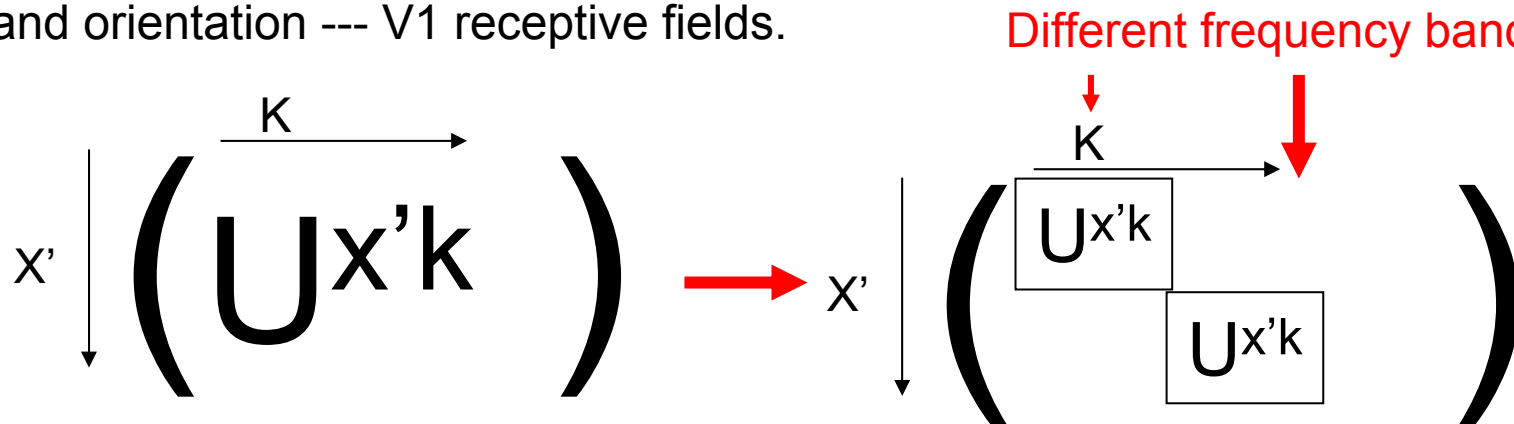
How can one understand the orientation selective receptive fields in V1?

Recall the retinal encoding transform:

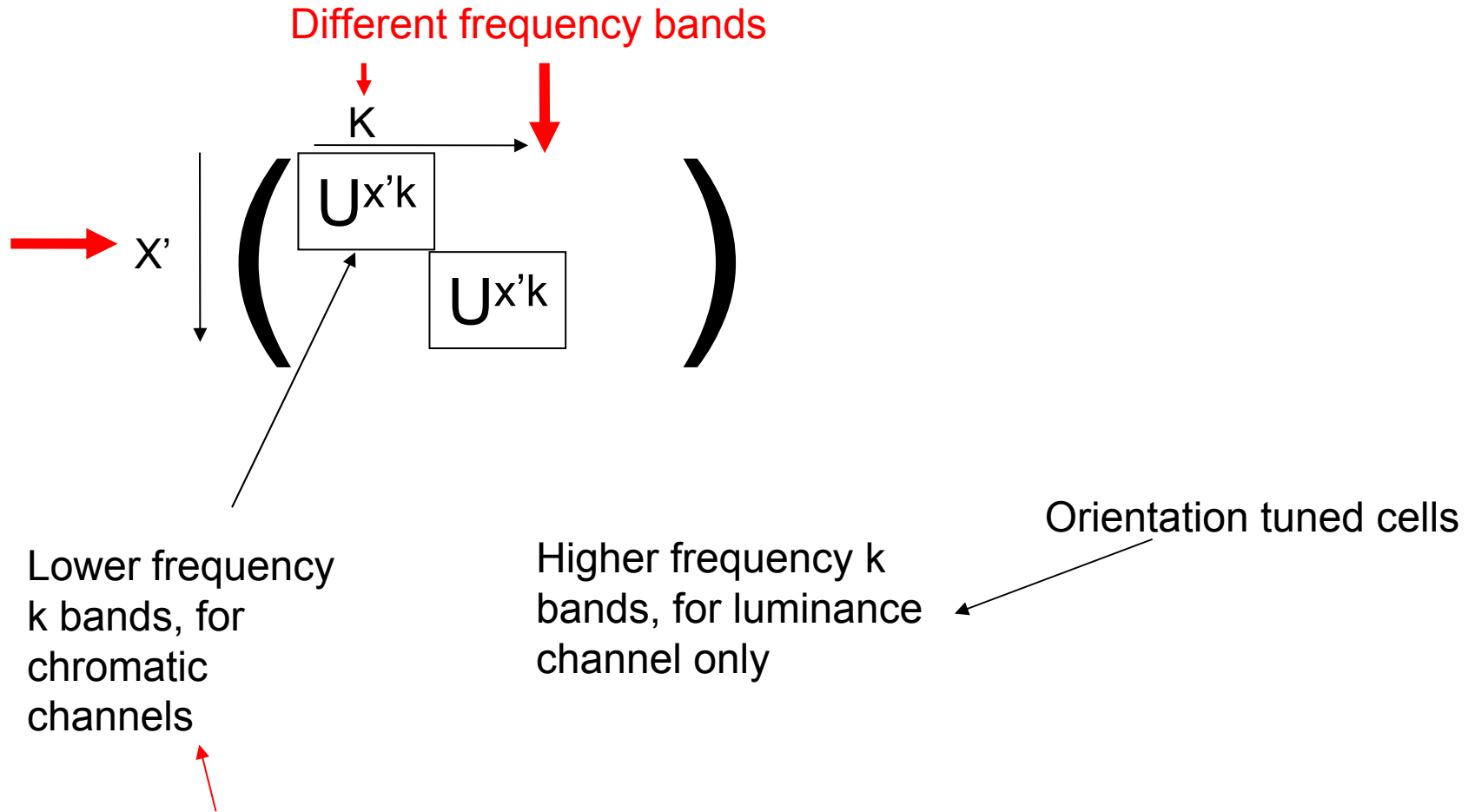
$$O_{x'} = \sum_k U^{x'k} g(k) \sum_x e^{-kx} S_x = \sum_k g(k) \underbrace{\sum_x e^{-k(x'-x)} S_x}_{\text{center-surround spatial filter}} + \text{noise}$$



If one changes the multiplexing filter $U^{x'k}$, such that it is block diagonal, and for each output cell x' , it is limited in frequency band in frequency magnitude and orientation --- V1 receptive fields.



V1 Cortical color coding



In V1, color tuned cells have larger receptive fields, have double opponency

Question: if retinal ganglion cells have already done a good job in optimal coding by the center-surround receptive fields, why do we need change of such coding to orientation selective? As we know such change of coding does not improve significantly the coding efficiency or sparseness.

Answer? Ref: (Olshausen, Field, Simoncelli, etc)

Why is there a large expansion in the number of cells in V1?

This leads to increase in redundancy, response in V1 from different cells are highly correlated.

What is the functional role of V1? It should be beyond encoding for information efficiency, some cognitive function beyond economy of information bits should be attributed to V1 to understand it.