

# Center for Brains, Minds & Machines

---

CBMM Memo No. 017

June 6, 2014

## Computational role of eccentricity dependent cortical magnification

by

**Tomaso Poggio<sup>\*,†</sup>, Jim Mutch<sup>\*</sup>, Leyla Isik<sup>\*</sup>**

<sup>\*</sup> Center for Brains, Minds and Machines, McGovern Institute, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>†</sup> Istituto Italiano di Tecnologia, Genova, Italy

**Abstract:** We develop a sampling extension of M-theory focused on invariance to scale and translation. Quite surprisingly, the theory predicts an architecture of early vision with increasing receptive field sizes and a high resolution fovea – in agreement with data about the cortical magnification factor, V1 and the retina. From the slope of the inverse of the magnification factor, M-theory predicts a cortical “fovea” in V1 in the order of 40 by 40 basic units at each receptive field size – corresponding to a foveola of size around 26 minutes of arc at the highest resolution,  $\approx 6$  degrees at the lowest resolution. It also predicts uniform scale invariance over a fixed range of scales independently of eccentricity, while translation invariance should depend linearly on spatial frequency. Bouma’s law of crowding follows in the theory as an effect of cortical area-by-cortical area pooling; the Bouma constant is the value expected if the signature responsible for recognition in the crowding experiments originates in V2. From a broader perspective, the emerging picture suggests that visual recognition under natural conditions takes place by composing information from a set of fixations, with each fixation providing recognition from a space-scale image fragment – that is an image patch represented at a set of increasing sizes and decreasing resolutions.



This work was supported by the Center for Brains, Minds and Machines (CBMM), funded by NSF STC award CCF - 1231216.

# Computational role of eccentricity dependent cortical magnification

June 9, 2014

Tomaso Poggio<sup>\*,†</sup>, Jim Mutch<sup>\*</sup>, Leyla Isik<sup>\*</sup>

<sup>\*</sup> Center for Brains, Minds and Machines, McGovern Institute, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>†</sup> Istituto Italiano di Tecnologia, Genova, Italy

## Summary

*We develop a sampling extension of M-theory focused on invariance to scale and translation. Quite surprisingly, the theory predicts an architecture of early vision with increasing receptive field sizes and a high resolution fovea – in agreement with data about the cortical magnification factor, V1 and the retina. From the slope of the inverse of the magnification factor, M-theory predicts a cortical “fovea” in V1 in the order of 40 by 40 basic units at each receptive field size – corresponding to a foveola of size around 26 minutes of arc at the highest resolution,  $\approx 6$  degrees at the lowest resolution. It also predicts uniform scale invariance over a fixed range of scales independently of eccentricity, while translation invariance should depend linearly on spatial frequency. Bouma’s law of crowding follows in the theory as an effect of cortical area-by-cortical area pooling; the Bouma constant is the value expected if the signature responsible for recognition in the crowding experiments originates in V2. From a broader perspective, the emerging picture suggests that visual recognition under natural conditions takes place by composing information from a set of fixations, with each fixation providing recognition from a space-scale image fragment – that is an image patch represented at a set of increasing sizes and decreasing resolutions.*

## Motivation and plan

A recent theory – dubbed “M-theory”[1] – conjectures that the main computational task of the ventral stream of visual cortex is to compute a representation of any new image in terms of a signature which is invariant to transformations that have been experienced in the past with other images and objects. The conjecture

is motivated by theorems showing that invariant representations can significantly reduce the sample complexity of a supervised learning classifier. The magic theory proposes a family of algorithms for unsupervised learning of an invariant representation. The computation of an invariant representation starts in the early visual areas with affine transformations, mainly 2D translations and scale changes (and some rotation). Here we show that an analysis of the sampling aspects of the theory offers an intriguing interpretation of the role of eccentricity-dependent size of receptive fields in V1, V2, V4 (and IT) – and of the linear dependency with eccentricity of the spatial sampling of the photoreceptors in the primate retina. It also yields several predictions.

This memo is divided into two parts: Part I analyzes eccentricity dependence in the first stage of vision (eg V1, eg the first layer in a hierarchical network), while Part II discusses plausible scenarios for the higher visual areas (eg higher layers in the hierarchical network). The main new results of Part I are:

- simultaneous invariance to translation and scale leads to an inverted, truncated pyramid as a model of the set of receptive fields in V1; this is a new alternative, as far as we know, to the usual smooth empirical fit of the cortical magnification factor data.
- the model above predicts, as a special case, results from Anstis about recognition of letters at different eccentricities [2].
- the model also predicts that shift invariance depends on resolution whereas scale invariance is uniform across eccentricities.
- the size of the flat, high acuity foveal region, which we identify with the *foveola* can be inferred from the slope of the eccentricity dependent acuity. Existing neural data from macaque V1 suggest a foveola with a diameter of around 20' of arc.
- scale invariance turns out to be more natural and complete than shift invariance. This is to be expected in organisms in which eye fixations can easily take care of shift transformations whereas more expensive motions of the whole body would be required to change the scale.

A knowledge of the main results of the magic theory is assumed. Much of the results and discussion below is in computational terms without reference to neuroscience. Biological implications and predictions are presented in the final section.

## 1 Part I: V1 simple cell receptive fields as stored transformations of templates

The theory shows how the first step of the invariance computation may consist of unsupervised learning of the set of transformations. This is done by

1. storing template  $t_k$  (which is an image patch and could be chosen at random),
2. storing all its observed transformations (bound together by continuity in time), and
3. repeating the above process for a set of  $K$  templates.

The simplest and most common image transformations are similitude transformations, that is shifts in position and uniform changes in scale. Though the templates can be arbitrary image patches, we will assume that V1 templates are Gabor-like functions (windowed Fourier transforms).

### 1.1 Geometry

The geometry of scaling is of course independent of eccentricity dependence in the retina or V1. Under scaling, a pattern exactly centered at 0 on the  $x$  axis (eg the center of the fovea) will increase in size without any translation while its boundaries will shift in  $x$ ; a pattern at say  $x_0$  will increase in size and its center will translate in the  $s, x$  plane, according to  $g_a I(x - x_0) = I(\frac{x-x_0}{a})$ , see Figures 1 and 2. In the  $s, x$  plane the slope of the trajectory of a pattern under scaling is a *straight line through the origin* with a slope that depends on the size of the pattern  $s$  and the associated position.

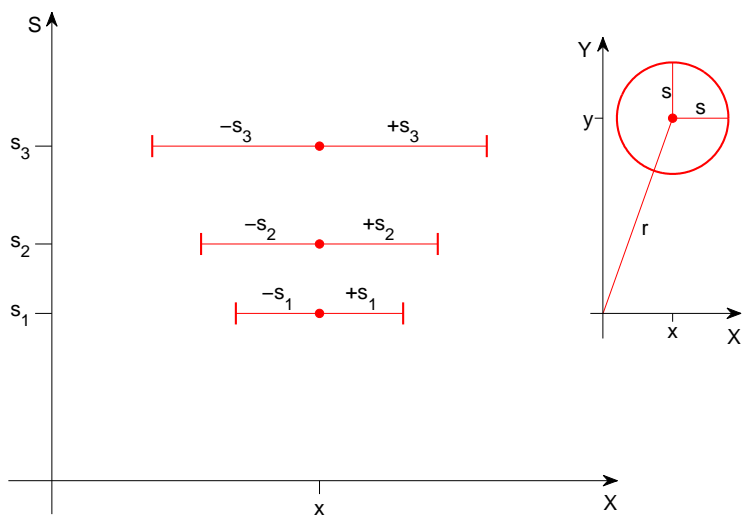


Figure 1: The  $s, x$  plane, showing example templates of radius  $s_1$ ,  $s_2$ , and  $s_3$ , all centered at position  $x$ . Note that while  $s$  and  $x$  are both measured in degrees of visual angle, in this plot the two axes are not shown to the same scale. Here, as in the rest of the paper, we show only one spatial coordinate ( $x$ ); everything we say can be directly extended to the  $x, y$  plane. Inset: an example template of radius  $s$  shown in the original  $x, y$  plane.  $x$  in the  $s, x$  plane corresponds to  $r$  in the  $x, y$  plane. We will assume later that the smallest template is the smallest simple cell in the fovea with a radius of around  $40''$  [11].

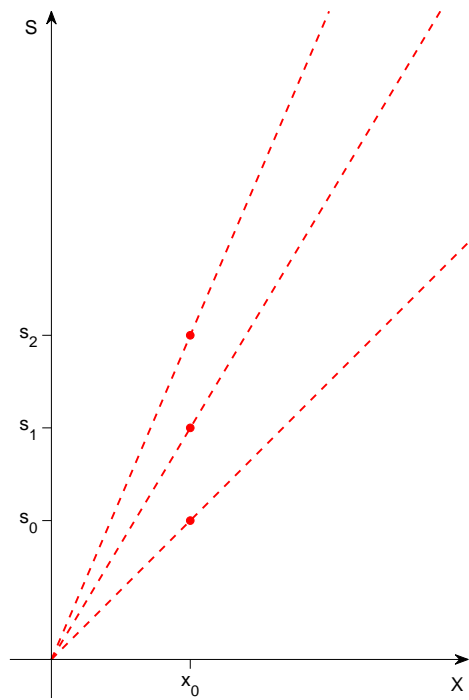


Figure 2: For any fixed position  $x_0$ , the size of a pattern or template at that position determines the slope of its  $s, x$  trajectory under scaling centered at  $x = 0$ . Note that the  $s, x$  trajectories under scaling are straight lines through the origin.

## 1.2 Magic window instead of M scaling functions

Let us start with a template at  $s = s_0, x = 0$ .  $s_0$  could be for instance the minimum possible receptive field size (given constraints such as the optics). Transforming it by shifts in  $x$  within  $-x_0, +x_0$  (the translation group, as the scale group, is only locally compact) generates a set  $S_0$ . Suppose that we want to ensure that what is recognizable at the highest resolution ( $s_0$ ) remains recognizable at all scales – that is distances from the object – up to  $s_{max}$ . The associated scale transformations of the set  $S_0$  yield the inverted truncated pyramid shown in Figure 3. Pooling over that set of transformed templates according to the magic algorithm will give uniform invariance to all scale transformations of a pattern over the range  $(s_0, s_{max})$ ; invariance to shifts will be at least within  $(-x_0, x_0)$ , depending on scale. Note that the above process of observing and storing a set of transformations of templates in order to be able to compute invariant representations may take place at the level of evolution or at the level of development of an individual system or as a combination of both.

The following definition holds: *the inverted truncated pyramid of Figure 3 is the locus  $S$  of the points such that their scaling between  $s_{min}$  and  $s_{max}$  gives points in  $S$ ; further all points  $P$  between  $-x_0$  and  $x_0$  are in  $S$ , eg  $P \in S$  if*

$$g_s P \in S, s \in [s_{min}, s_{max}], S_0 \in S$$

where  $S_0$  consists of all points at  $s_{min}$  between  $-x_0$  and  $x_0$ .

The region above follows naturally if scale invariance has a higher priority than shift invariance. Alternatively, the complementary constraint of uniform translation invariance for all points generates in  $s, x$  space a rectangular region between  $x_{min}$  and  $x_{max}$  and  $s_{min}$  and  $s_{max}$ . In this case, points in the region do not allow uniform scale invariance (in fact the scaling ranges from 0 to  $s_{max} - s_{min}$ ). The region of Figure 3 may be more “natural”: we conjecture that the region of Figure 3 may be optimal in terms of containing simultaneously maximum information about scale and space (the open problem is to make the notion of optimality precise). Notice that the best template in terms of simultaneously maximizing scale and shift invariance is a Gabor function (see [1]). Under scale and translation the Gabor template originates

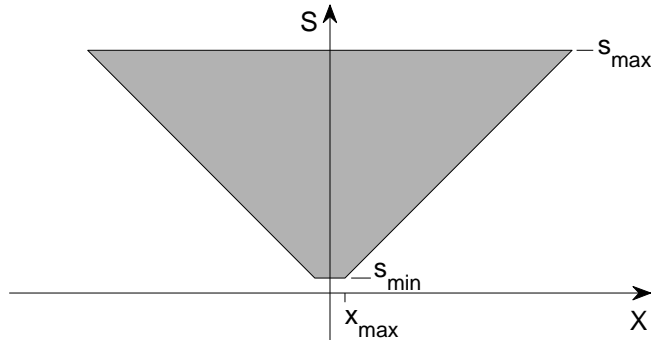


Figure 3: Let us start with a template at  $s_{min}$  and  $x = 0$ . Its transformations under (bounded) shift at  $s = s_{min}$  fill the line between  $-x_{max}$  and  $x_{max}$ . Its transformations under (bounded) scaling fill the space in the truncated, inverted pyramid shown in the figure. This is the space of bounded joint transformations in scale and space. The axes  $s, x$  are here in the same units.

a set of Gabor wavelets (a tight frame). In the region of Figure 3 scaling (between  $s_{min}$  and  $s_{max}$ ) each wavelet generates other elements of the frame.

Notice that the shape of the lower boundary of the inverted pyramid, though similar to standard empirical functions that have been used to describe M scaling, is different for small eccentricities. An example of an empirical function (Covey and Rolls, 1974 [4], described in [16]) is  $M^{-1} = M^0(1+ax)$ , where  $M$  is the cortical magnification factor,  $M^0$  and  $a$  are constants and  $x$  is eccentricity.

### 1.3 Inverted pyramid scale-space image fragments

A normal, non-filtered image may activate all or part of the Gabor filters within the inverted truncated pyramid of Figure 3. The pattern of activities is related to a multi-resolution wavelet decomposition of an image. We call this transform of the image an “inverted pyramid scale-space fragment” (“IP fragment” in short) and consider it as supported on a domain in the 2D space  $x, s$  that is contained in the inverted truncated pyramid of the figure. The fragment corresponding to a bandpass filtered image should be a more or less narrow horizontal slice in the  $s, x$  plane. The term fragment is borrowed from Ullman.



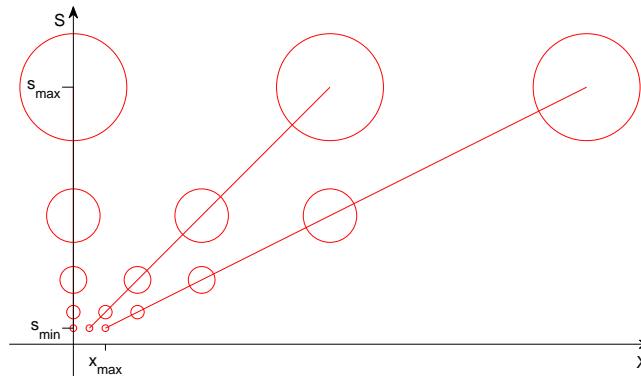


Figure 4: Under the assumption of Gabor filters and associated sampling for each scale at spatial intervals  $\Delta x = s$ , this graph depicts a subset of the resulting array of template units. Note that in this graph the sampling over  $s$  is arbitrary but the sampling over  $x$  follows the sampling theorem. Note that there will be no samples between the  $s$  axis and the line with slope 1 (when  $x$  and  $s$  are plotted in the same units). The center of the circles in the figure gives the  $s, x$  coordinates; the circles are icons symbolizing here the receptive fields.

#### 1.4 Sampling the templates: specific assumptions

In the following we *assume for simplicity* that the template is a Gabor filter (of one orientation; other templates may have different orientations). We assume that the Gabor filter and its transforms under translation and scaling are roughly bandpass and the sampling interval at one scale over  $x$  is  $s$ , implying half overlap of the filters in  $x$ . This is illustrated in Figure 4.

The above assumptions imply that for each array of filters of size  $s$  the first unit on the right of the central one (at  $x = 0$ ) is at  $x = s$ , where  $x$  and  $s$  are measured in the same units.

For the scale axis we follow the sampling pattern estimated by Marr et al.[11] with 5 “frequency channels” having  $2s = 1'20''$ ,  $3.1'$ ,  $6.2'$ ,  $11.7'$ ,  $21'$ . Filter channels as described above are supported by sampling by photoreceptors that starts in the center of the fovea at the Shannon rate, dictated by the diffraction-limited optics with a cutoff around 60 cycles/degree, and then decreases as a function of eccentricity.

## 1.5 Prediction: size of the foveola and slope of the magnification factor

In this “truncated pyramid” model of V1 (we assume that the Gabor templates correspond to simple cells in V1, called S1 in HMAX) the slope of the magnification factor  $M$  as a function of eccentricity depends on the size of the foveola – that is the region at the minimum scale  $s_{min}$ . The larger the foveola, the smaller the slope.

We submit that this “engineering model” nontrivially fits data about the size of the fovea, the slope of  $M$  and other data about the size of receptive fields in V1. In particular the size of the foveola, the size of the largest RFs in AIT and the slope of acuity as a function of eccentricity depend on each other: fixing one determines the other (after setting the range of spatial frequency channels, i.e., the range of RF sizes at  $x = 0$  in V1).

### 1.5.1 Back of the envelope estimates: size and sampling of foveola and fovea

As a back of the envelope calculation we assume here that  $s_{min} \approx 40''$  (from an estimate of  $1'20''$  for the diameter of the smallest simple cells [11], see also [12]). Data of Hubel and Wiesel [9, 10] (shown in Figure 5) and Gattass [7, 8] (shown in Figure 6) yield an estimate of the slope  $a$  for  $M$  in V1 (the slope of the line  $s_{min}(x) = ax$ ). Hubel and Wiesel cite the slope of the average RF diameter in V1 as  $a = 0.05$ ; Gattass quotes a slope of  $a \approx 0.16$  (in both cases the combination of simple and complex cells may yield a biased estimate relative to the “true” slope of simple cells alone). Our model of an inverted truncated pyramid predicts (using an estimate of  $a = 0.1$ ) from these estimates that the radius of the foveola (the bottom of the truncated pyramid) is  $R = 1'20/0.1 \approx 13'$  with a full extent of  $2R \approx 26'$  corresponding to about 40 cells separated by  $40''$  each. The size of the fovea (the top of the truncated pyramid) would then have  $2R \approx 6^\circ$  with 40 cells spaced  $\approx 21'$  apart; see Figure 7.

*These estimates depend on the actual range of receptive field sizes and could easily be wrong by factors of 2.* We have used liberally data from the macaque together with data from human psychophysics. Our main goal is to provide a logical interpretation of future data

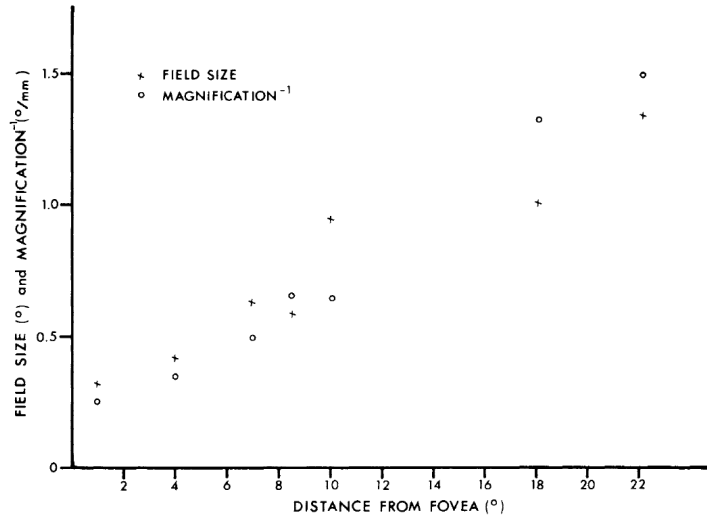


Figure 5: Data of Hubel and Wiesel for monkey V1 gives a slope for average RF diameter, relative to eccentricity, of  $a = 0.05$ . From [10].

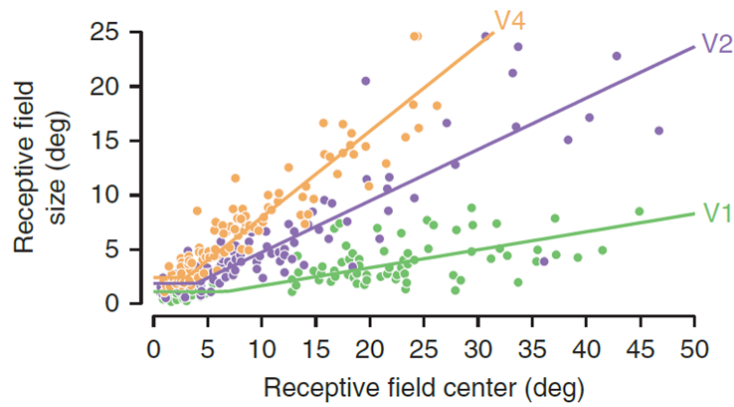


Figure 6: V1 data supports a model in which the maximum scale also depends on eccentricity, as in figure 9 (lower). Adapted from [6] (original monkey data from [7, 8]). Note that the y-axis here represents scale using RF diameter, not spacing.

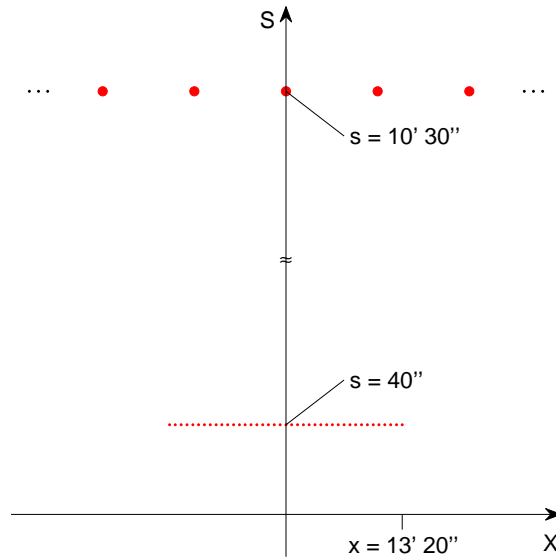


Figure 7: The figure shows a foveola of  $\approx 26'$  and a total of  $\approx 40$  units (20 on each side of the center of the fovea). It also shows sampling intervals at the coarsest scale in the fovea (assumed to be around  $2s = 21'$ [11] which would span  $\approx \pm 6^\circ$ ). Note that the size of a letter on the ophthalmologist's screen for 20/20 vision is around  $5'$ .

and a ballpark estimate of relevant quantities.

### 1.6 Sampling array in V1

It is possible to map the nonuniform sampling array into a square lattice; see Figure 8. (It is actually a cube when  $y$  is included; note that a further mapping of this cube onto a 2D cortical sheet would be similar to – but different from! — the roughly log-polar mapping commonly assumed to take place between the retina and V1.)

Once the  $s, x$  array has been transformed into a square lattice, the lattice behaves like  $x, y$  from the point of view of group transformations, that is shifts and “scaling” commute. An array spanning  $26'$  at the bottom and  $6^\circ$  at the top is a kind of image, but in pixel and scale space.

Each of the scale layers has the same number of units which is determined by the number of units in the fovea – that is, the number of units at the finest resolution. This remapping shows that S1 corresponds to a lattice of dimensions  $x, y, s, \theta$ , where the

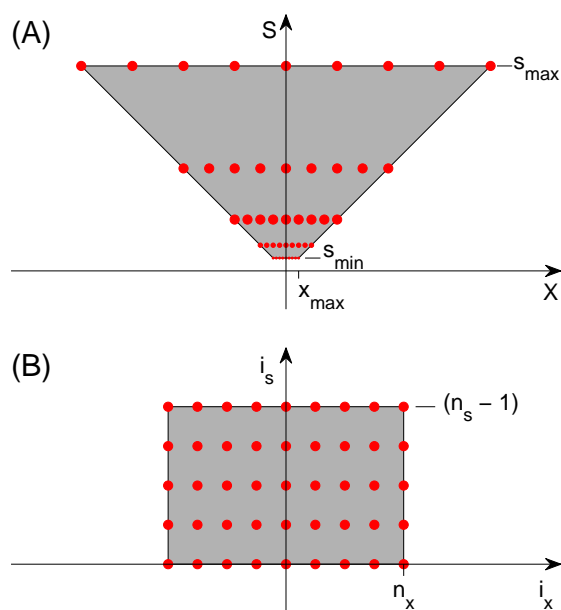


Figure 8: Since the inverted truncated pyramid (A) has the same number of sample points at every scale, it maps perfectly onto a square array (B) when  $x$  is replaced by  $i_x = x/s$ , i.e., the number of samples from the center.  $i_s$  is the scale band number. We call this the “magic map”.

dimension sizes are different (but roughly the same for  $x, y$ ); the topology is that of a cylinder with  $\theta$  being periodic.

Coordinates  $(i_s, i_x)$  within the square lattice are obtained from  $(s, x)$  as follows:

$$i_s = \log_f\left(\frac{s}{s_{min}}\right)$$

$$i_x = \frac{x}{s}$$

where  $f$  is the multiplicative factor between adjacent scale bands.  $i_s$  is the number of the scale band, starting with 0 for the finest band.  $i_x$  is the number of samples (in position) from the center, and may be positive or negative. The lattice will have size  $n_s$  in scale and  $2n_x + 1$  in position.

As we have described, the region of Figure 9 (upper) is obtained by setting a minimal scale  $s_{min}$ , shifting the template between  $-x_0$  and  $x_0$  and then scaling all the templates obtained so far up to  $s_{max}$ . One of many alternatives is to set a constant difference between the minimum and the maximum scale at each eccentricity ( $s_{max} - s_{min} = const.$ ); see Figure 9 (lower). Experimental data suggests this is a more likely possibility (large eccentricities are also represented in the visual system); see Figure 6, though several variations are of course equally possible.

## 1.7 Pooling in C1

Invariance is not provided directly by the array but by the pooling over it. Note that we are limiting ourselves in this section to the invariant recognition of *isolated objects*.

The range of invariance in  $x$  is limited for each  $s$  by the slope of the lower bound of the inverted pyramid (see Figure 8 plot A). The prediction is that the range of invariance  $\Delta x$  depends on scale  $s$  as

$$\Delta x \approx n_x s,$$

where  $n_x$  is the radius of the inverted pyramid in samples, and is constant for all scales. Thus small details (high frequencies) have a limited invariance range in  $x$  whereas larger details have larger invariance.  $n_x$  is obtained from the slope ( $a$ ) of the cortical magnification factor (section 1.5) as:

$$n_x = 1/a$$

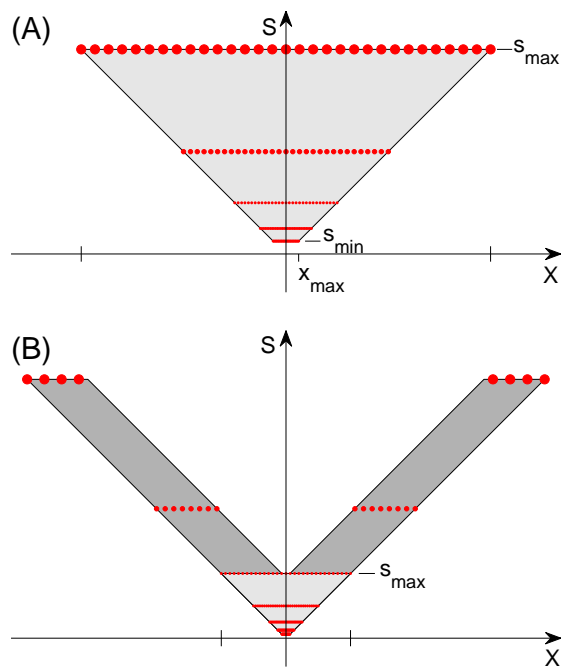


Figure 9: Regions of scaled and shifted templates. In the upper plot we have a constant maximum scale ( $s_{max}$ ). In the lower plot we consider a wider range of eccentricities, and there is a constant difference between the largest and smallest scale at each eccentricity. Note that the central region in the lower plot (lighter shading) is the same one shown in the upper plot.

## 2 Part II: Hierarchy and decimating the array

We follow the magic theory's extension to hierarchies of pooling, dot products, sampling, pooling. At each dot product (S unit) stage or pooling (C unit) stage, there can be a downsampling of the array of units (in  $x, y$  and possibly in  $s$ ) that follows from the low-pass-like effect of the dot-product or pooling operation.

Here we assume a specific strategy of downsampling in space by 2 at each stage. This choice is for simplicity but is roughly consistent with biological data. The results below can be easily changed using different criteria for downsampling. The following statement follows from considerations of invariance for image patches that are smaller than half the pooling region:

**Sampling theorem for pooling** *Pooling over  $2 \times 2$  units gives local invariance to scale and translation and permits downsampling of the array in  $x$  by a factor of 2.*

This simple result is discussed elsewhere. We assume here that each combined S-C stage brings about a downsampling by 2 of the  $x, y$  array. We call this process "decimation"; see Figure 10. Starting with V1, 4 stages of decimation (possibly at V2, V4, TEO and AIT) reduce the number of units at each scale from  $\approx 40$  to  $\approx 2$  spanning  $\approx 26'$  at the finest scale and  $\approx 6^\circ$  at the coarsest.

Pooling over scale in a similar way may also decimate the array down to  $\approx 1$  scale from V1 to IT. Neglecting orientations, in  $x, y, s$  the  $30 \times 30 \times 6$  array of units may be reduced to just a few units in  $x, y$  and 1 in scale. This picture is consistent with the invariance found in IT cells ([15]). According to the magic theory, different types of such units are needed, each for one of several templates at the top level.

### 2.1 Clutter and crowding

So far our discussion has assumed isolated objects in the image – which is the usual assumption in most theorems of M-theory. Clutter of course poses a problem. For instance, pooling over the whole  $s, x$  domain of Figure 3 could suffer from the presence of clutter-induced fragments in any location in the inverted pyramid. A more natural form of pooling is the (specific) layer-by-layer decimation described in the previous section (see also Fig-



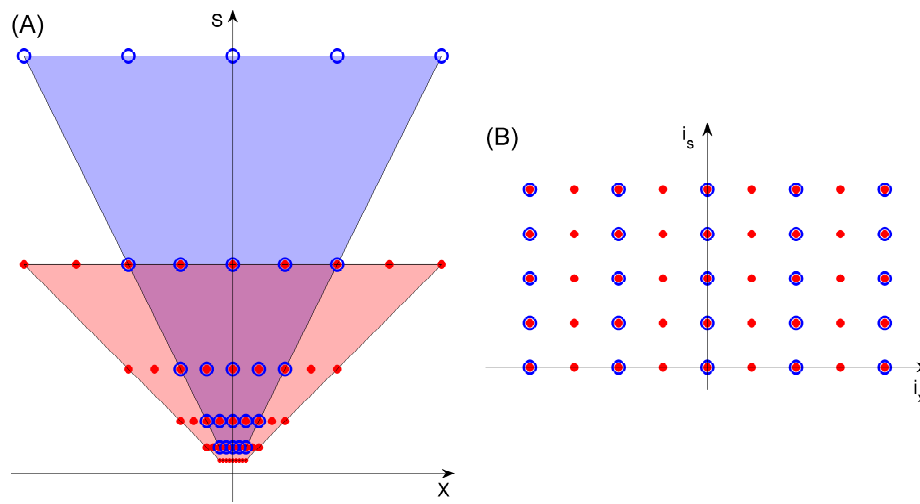


Figure 10: Pooling over  $2 \times 2$  lattices of the  $s, x$  array in V1 and subsampling reduces the square lattice; if the lattice in  $x$  is 40 units, 4 layers (V2, V4, TEO, AIT) of  $x$  pooling are sufficient to create cells that are 16 times larger than the largest in the fovea in V1 (probably around  $21'$  at the coarsest scale), yielding at the top cells with a RF diameter of up to  $\approx 5^\circ$ . Each area in the fovea would see a doubling of size with corresponding doubling of the slopes at the border (before remapping to a cube lattice). The index of the units at position  $x$  and scale  $s$  is given by  $i_x^s = 2^s i_x^1$ .

ure 10 B). Notice that the corresponding pooling range is uniform across eccentricities in the coordinates of Figure 10 B. The spatial pooling range on  $x$  depends on the area: for V1 it is the sampling interval between the red dots, for V2 it is the sampling interval between the blue dots: they should be roughly equivalent to the radius of the receptive field of the complex cells in V1 and V2 respectively. We assume now the following reasonable criterion for pooling to remain unaffected by clutter, that is interference-free from a flanking object: *the target and the flanking distractor must be separated at least by the pooling distance and thus by a complex cell receptive field*. Under this assumption, our theory predicts that the critical separation for avoiding crowding should be

$$\Delta x \approx bx \tag{1}$$

since the RF size increases linearly with eccentricity, with  $b$  depending on the cortical area responsible for the recognition signal (see Figure 10 A). Thus the theory “predicts” Bouma’s law of crowding! [3] The experimental value found by Bouma for crowding of  $b \approx 0.4$  suggests that the area is V2 since this is the slope found by Gattass for the dependence of V2 RF size on eccentricity. Studies of “metameric” stimuli by Freeman and Simoncelli also implicated V2 in crowding and peripheral vision deficiencies [6].

Simulations by Isik et al. (2011) with an HMAX-type model confirm the intuition that for the model to be interference-free the spatial separation between the target object and the flanking object must be  $\approx bx$ .

### 3 Discussion

#### 3.1 A magic role for V1

The most striking result of this paper is that the linear increase of RFs size with eccentricity, found in all primates, is required by M-theory to compute a scale and position invariant representation. Also predicted by the computational theory is the existence of the foveola and its link with the slope of  $M^{-1}$ — though not its (small) size that certainly depends on resource constraints (such as the size of the optical nerve).

An indirect result with potentially interesting implications is that the representation at the level of S1 (simple cells in V1) after the “magic map” of Figure 8 corresponds to a discrete *abelian* group for shift and “scale”. In other words, in this representation, transformations in  $i_s$  and  $i_x$  commute. This property is inherited by higher levels in the network and implies possibly interesting properties for the covariance function associated with the image representation at each level and for its eigenfunctions.

### 3.2 A framework for understanding “everyday” human recognition

This state of affairs means that there is a quite limited “field of vision” in a single glimpse. Most of an image of 5 by 5 degrees is seen at the coarsest resolution only, if fixation is in its center; at the highest resolution only a fraction of the image (up to 30’) can be recognized, and an even smaller part of it can be recognized in a position invariant way (the number above are rough estimates).

We introduced earlier the idea of “IP fragment” corresponding to the information captured from a single fixation of an image. Such a fragment is supported on a domain in the 2D space  $x, s$ , contained in the inverted truncated pyramid of Figure 3. There are two interesting remarks:

1. For normally-sized images and scenes with fixations well inside, the resulting IP-fragment will occupy the full spatial and scale extent of the inverted pyramid. Consider now the fragment corresponding to an object to be stored in memory (during learning) and recognized (at run time). Suppose the best situation for the learning stage (the other cases can be discussed in terms of the figure): the object is close so that both coarse and fine scales are present in its fragment. At run time then the object can be recognized whenever it is closer or farther away. The important point is that a look at this and the other possible situations (at the learning stage) suggest that the matching should always weight more the finest available frequencies (bottom of the pyramid). This is the finding of Schynz [14]. As implied by his work, top-down effects may modulate somewhat these weights (this could be done during pooling) depending on the task.

2. Assume that such a fragment is stored in memory for each fixation of a novel image. Then there is the following trade-off between learning and run-time recognition:
  - if a new object is learned from a single fixation, recognition may require multiple fixations at run time to match the memory item (given its limited position invariance).
  - if a new object is learned from multiple fixations, with different fragments stored in memory each time, run time recognition will need a lower expected number of fixations.

Notice that the above argument critically rely on “large” and position-independent scale invariance: scale does not need to be sampled – unlike position.

The fragments of an image stored in memory via multiple fixations could be organized into an egocentric map. Though the map is *not directly used* for recognition (this is a conjecture that should be experimentally verified), it is *probably needed* to plan fixations during the learning and especially during the recognition and verification stage (and thus indirectly used for recognition in the spirit of minimal images and related recent work by Ullman and coworkers). Such a map may be somewhat related to Marr’s  $2\frac{1}{2}$  sketch, for which no neural correlate has been found as yet.

### 3.2.1 IP-based recognition: a new paradigm for computer vision?

Learning and recognition relying on a set of IP fragments (see the tradeoff discussed above and its context) may provide a new powerful framework for object recognition. Computer vision seems so far to mostly rely on matching only the coarsest scale available in the pyramid resulting in performance which is different from human performance.

### 3.3 Some predictions

We collect here the key predictions of the theory:

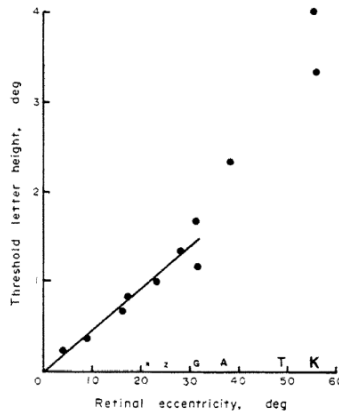


Figure 11: *Linear size dependency of letter recognition (Anstis, 1974 [?]).*

- The computational requirement of scale invariance implies linear dependency on eccentricity of receptive field sizes; it also predicts that the size of the foveola depends on the slope.
- Image patches that are recognized at some scale and eccentricity remain equally recognizable under scaling between  $s_{min}$  and  $s_{max}$ . This is a computational prediction of Anstis' findings [2]. Anstis found (see Figures 11, 12) that a letter just recognizable at some eccentricity remains equally recognizable under scaling: the size of receptive fields in V1 depends on eccentricity exactly to match the effect of scaling.
- Anstis did not measure close to the minimum letter size – which is around 5' for 20/20 vision. Our prediction is that if there is a range of receptive fields in V1 between  $s_{min}$  and  $s_{max}$  in the fovea then there is a finite range of scaling between  $s_{min}$  and  $s_{max}$  under which recognition is maintained (see 9). It is obvious that looking at the image from an increasing distance will at some point make it unrecognizable; it is somewhat less obvious that getting too close will also make it unrecognizable (this phenomenon was found in Ullman's minimal images; Ullman, personal comm.)
- Consider the experimental use of images such as letters of appropriate sizes that are bandpass filtered (with the filters mentioned in the text for human vision). Our predictions –

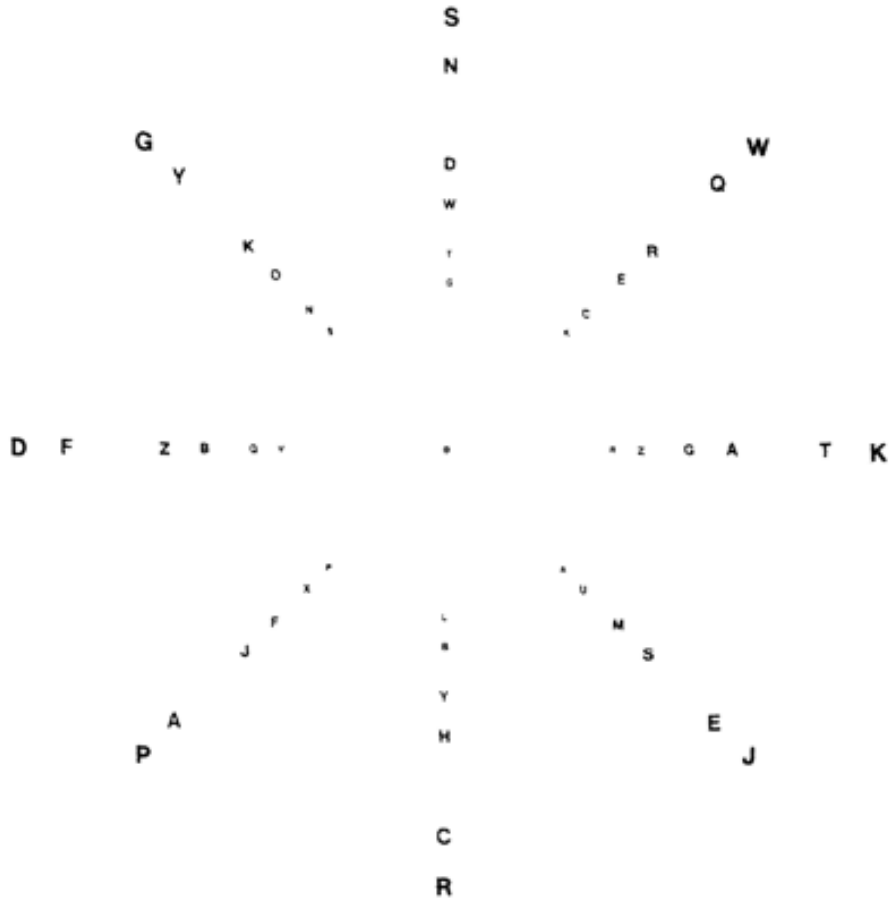


Figure 12: When fixating the central point, recognizability of a letter in the chart does not change under scaling (Anstis, 1974).

extending Anstis' findings – are that *there will be scale invariance for all frequencies between  $s_{min}$  and  $s_{max}$ ; there will be shift invariance that increases linearly with spatial wavelength and is at any spatial frequency at least between  $x_{min}$  and  $x_{max}$  (the bottom edge of the truncated pyramid).*

- The theory predicts a flat region of constant maximum resolution – which we suggest may correspond to the so-called *foveola*. This is different from the empirical equations that have been proposed in the past but it is unclear whether the difference is experimentally detectable, especially because the biological implementation may well be an approximation of the theory.
- The theory predicts that the slope of the line of minimum  $s$  for each  $x$  depends on the size of the foveola and vice-versa. Since the slope can be estimated relatively easily from a number of existing data, our prediction for the linear size of the foveola is around 40 minutes of arc, corresponding to about 30 simple cells of the smallest size (assumed to be  $\approx 1'20''$  of arc). Notice our definition of the fovea is in terms of the set of all scaled versions of the foveola between  $s_{min}$  and  $s_{max}$  spanning about 6 degrees of visual angle.
- The theory predicts crowding effects due to pooling. It predicts Bouma's law and its linear dependency on eccentricity [3].
- Since Bouma's constant has a value of about 0.4 (see also [6]), our theory predicts that an interference-free signature from V2 is important for recognition. This is consistent with M-theory independent requirement [1] that signals associated with image patches of increasing size are accessing memory and classification stages separately. The requirement follows from the need of recognizing "parts and whole". The V2 signal could directly or indirectly (via IT or V4 and IT) reach memory and classification.
- The angular size of the fovea remains the same at all stages of a hierarchical architecture (V1, V2, V4...), but the number of units per unit of visual angle decreases and the slope increases because the associated  $s$  increases (see Figure 10).

- Top-down control signals may be able to control the extent of pooling, or the pooling stage, used in computing a signature from a region of the visual field. This process may correspond to attentional suppression of the non-attended region.
- If we assume that some AIT neurons effectively pool over “all” positions and scales their invariant receptive field (over which consistent ranking of stimuli is maintained) should be smaller with higher frequency patterns than with low frequency ones. This prediction is consistent with the fact that some physiology data suggest that AIT neurons are somewhat less tolerant to position changes of small stimuli (Op de Beeck and Vogels 2000 [13], also described in [5]). A comparison across studies suggests that position tolerance is roughly proportional to stimulus size[5].

*Acknowledgments* Cheston Tan suggested the connection with the DiCarlo and Maunsell paper.

## References

- [1] F. Anselmi, J.Z. Leibo, L. Rosasco, J. Mutch, A. Tacchetti, T. Poggio. Magic Materials: a theory of deep hierarchical architectures for learning sensory representations *CBCL paper*, Massachusetts Institute of Technology, Cambridge, MA, April 1, 2013.
- [2] S.M. Anstis A chart demonstrating variations in acuity with retinal position *Vision Research*, 14, 589, 592, 1974.
- [3] H. Bouma Interaction Effects in Parafoveal Letter Recognition *Nature* 226 (5241) 177-178, 1970.
- [4] A. Cowey and E. T. Rolls Human cortical magnification factor and its relation to visual acuity *Experimental Brain Research* 21:447-454, 1974.
- [5] J. J. DiCarlo and J.H. R. Maunsell Anterior Inferotemporal Neurons of Monkeys Engaged in Object Recognition Can be Highly Sensitive to Object Retinal Position *J Neurophysiol* 89:3264-3278, 2003.



- [6] J. Freeman and E. Simoncelli Metamers of the ventral stream  
Nature Neuroscience 14:1195-1201, 2011.
- [7] Gattass, R and Gross, C G and Sandell, J H, Visual topography of V2 in the macaque. The Journal of comparative neurology, 201, 519–39, 1981
- [8] Gattass, R and Sousa, AP and Gross, CG Visuotopic organization and extent of V3 and V4 of the macaque J. Neurosci., 6, 1831–1845, 1988
- [9] D.H. Hubel and T.N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex  
*The Journal of Physiology* 160, 1962.
- [10] Hubel, D H and Wiesel, T N Uniformity of monkey striate cortex: a parallel relationship between field size, scatter, and magnification factor. The Journal of comparative neurology, 158, 295–305, 1974
- [11] D. Marr, T. Poggio, E. Hildreth Smallest channel in early human vision *JOSA*, Vol. 70, Issue 7, 868-870 1980
- [12] Mazer, James A and Vinje, William E and McDermott, Josh and Schiller, Peter H and Gallant, Jack L. Spatial frequency and orientation tuning dynamics in area V1. Proceedings of the National Academy of Sciences of the United States of America, 3, 99, 1645–50, 2002
- [13] H. Op de Beek and R. Vogels Spatial sensitivity of macaque inferior temporal neurons Journal of Comparative Neurology 426, 500-518, 2001.
- [14] P. G. Schyns and F. Gosselin Diagnostic Use of Scale Information for Componential and Holistic Recognition In: Peterson, M.A., Rhodes, G. (Eds.), Perception of Faces, Objects, And Scenes., Oxford University Press, pp. 120145, 2003.
- [15] Serre, T. and Kouh, M. and Cadieu, C. and Knoblich, U. and Kreiman, G. and Poggio, T. A theory of object recognition: computations and circuits in the feedforward path of the ventral stream in primate visual cortex CBCL Paper #259/AI Memo #2005-036, 2005

- [16] H. Strasburger, I. Rentschler, M. Jüttner Peripheral vision and pattern recognition: A review *Journal of Vision*, 11(5):13, 182, 2011
- [17] D. Whitney and D. M. Levi Visual crowding: a fundamental limit on conscious perception and object recognition *Trends in Cognitive Sciences*, April 2011, Vol. 15, No. 4

## Appendix: Notes and Remarks

We list a few experiments which are not critical predictions but may be worth considering:

- Create Anstis charts for different tasks (such as letter recognition, face recognition, symmetry detection...could they show different slopes and foveal sizes?)
- Do crowding parameters depend on task (say letter recognition vs face recognition)? They may provide information about the range of pooling at cortical area at level  $n$  that is responsible for specific task.
- Which experiments could best answer the empirical question of the shape of the upper boundary of the magic window?
- What is the scale and translation invariance of “minimal images”? It may be interesting to appropriately filter images to match one of the visual frequency channels.
- What are the eigenfunctions of the covariance induced by natural images at the level of S1, that is after the “magic map”?
- Can we measure the size of the foveola by using circular pattern with letters of different sizes?
- Psychophysics on position invariance did not yet provide a clear picture. Our theory suggests that this may be due to experimental designs that have ignored scale and spatial frequency variables. Could then sharper results be obtained by repeating experiments such as Nazir and O’Regan – testing position *and* scale invariance – with patterns which are band-pass filtered in spatial frequency?

We list a few remarks:

- Instead of bandpass filtering patterns in spatial frequency it may be possible to use spatial noise localized in space and frequency to block specific frequency channels (this is the dual of Schyns approach).

- *In a sense* evolution has given primacy to invariance to scale with respect to invariance to translation. A reason for it seems obvious for the human eye. Could this be tested considering other species with different tradeoffs between moving gaze and translating the body?
- Interestingly, the largest images recognizable at a glance map onto about  $30 \times 30$  simple cells at the finest resolution, spanning  $21'$ . (As a remark, the eye of *Drosophila* contains about 900 ommatidia.) They can be scaled up to span  $\approx 4^\circ$  at the coarsest resolution without changing the associated number of bits. There may be a connection with the framework of “minimal images” of Ullman. In terms of the theory there is a maximum image size at each resolution that can be recognized at a glimpse and be invariant to scaling, corresponding to the  $S$  cell size at that level of the hierarchy; the  $C$  cell pooling range provides invariance over that range. Notice that this state of affairs is independent of a multilevel hierarchical architecture; it follows from the *image hypercube* representation (see Figure 8) at each level (e.g., S1).
- Discrete transformations in  $i_s$  and  $i_x$  commute in the S1 layer after the “magic map” (see Figure 8). The property is maintained at higher layers.
- The size of the templates at the first layer (S1) determines the size of the minimal image patch. Of course patterns that are not very simple (eg bars) require patches that are at least two or three times the minimal  $1'20''$  (like letters on the ophtalmologist table).
- The pooling range (at the C1 level) affects the range of tolerance to clutter. The theory ensures invariance and uniqueness if all the patches that are pooled are without clutter (they may contain the object only).
- The pooling range at the first layer (C1) is likely to be the main determinant of the sampling interval at that level and at the S2 level.
- The set of templates for each C1 must be overcomplete; number of templates depends on the number of images to be discriminated  $n$  which depends (in an unknown way) on the

size of the image patch. One may conjecture that the total number of templates for an image patch of size  $n$  is less in  $k$  layers than in 1 layer.

- $M$  is in  $mmdeg^{-1}$  and  $M^{-1}$  is in  $mm^{-1}deg$ .  $M^{-1}$  is defined as the receptive field size (in degrees pr mm of cortex).