

Chapter VIII. Towards a biologically plausible computational model of ventral visual cortex

We have now come a long way since our initial steps towards defining the problem of visual recognition. We started with characterizing the spatial and temporal statistics of natural images (Lecture 2). We explored how neurons along ventral visual cortex respond to a variety of different stimulus conditions (Lectures 3, 5, 7, 8). We described the recognition impairments that arise through cortical lesions (Lecture 4) and the effect of applying currents to the neural circuitry (Lecture 9). We would like to put all of these separate bits and pieces of data into a coherent framework to rigorously understand how neuronal circuits help us recognize objects. Here we summarize some of the initial steps towards a theoretical understanding of the computational principles behind transformation-invariant visual recognition in the primate cortex.

8.1. Defining the problem

We start by defining what needs to be explained and the necessary constraints to solve the problem. A theory of visual object recognition, implemented by a computational model, should be able to explain the following phenomena and have the following characteristics:

1. *Selectivity*. The primate visual system shows a remarkable degree of selectivity and can differentiate among shapes that appear to be very similar at the pixel level (e.g. arbitrary 3D shapes created from paperclips, different faces, etc.). Critical to object recognition, a model should be able to discriminate among physically similar but distinct shapes.
2. *Transformation tolerance*. A trivial solution to achieve high selectivity would be to memorize all the pixels in the object. The problem with this type of algorithm is that it would not tolerate any changes in the image. An object can cast an infinite number of projections onto the retina. These changes arise due to changes in object position with respect to fixation, object scale, plane or depth rotation, changes in contrast or illumination, color, occlusion and others. The importance of combining selectivity and tolerance has been emphasized by many investigators (e.g. (Rolls, 1991; Olshausen et al., 1993; Logothetis and Sheinberg, 1996; Riesenhuber and Poggio, 1999; Deco and Rolls, 2004b; Serre et al., 2007b) among others).
3. *Speed*. Visual recognition is very fast, as emphasized by many psychophysical investigations (Potter and Levy, 1969; Kirchner and Thorpe, 2006; Serre et al., 2007a), scalp EEG measurements (Thorpe et al., 1996) and neurophysiological recordings in humans (Liu et al., 2009) and monkeys (e.g. (Richmond et al., 1983; Keysers et al., 2001; Hung et al., 2005) among others). This speed imposes an important constraint to the number of computational steps that the visual system can use for pattern recognition (Rolls, 1991; Serre et al., 2007b).

47 4. *Generic.* We can recognize a large variety of objects and shapes.
48 Estimates about the exact number of objects or object categories that primates
49 can discriminate vary widely depending on several assumptions and
50 extrapolations (e.g. (Standing, 1973; Biederman, 1987; Abbott et al., 1996; Brady
51 et al., 2008)). Certain types of shapes may be particularly interesting, they may
52 have more cortical real estate associated with them, they could be processed
53 faster and could be independently impaired. For example, there has been
54 extensive discussion in the literature about faces, their representation and how
55 they can be different from other visual stimuli. Yet, independently of precise
56 figures about the number of shapes that primates can discriminate and
57 independently also of whether natural objects and faces are special or not, it is
58 clear that there exists a generic system capable of discriminating among multiple
59 arbitrary shapes. For simplicity and generality, we focus first on this generic
60 shape recognition problem. Face recognition, or specialization for natural objects
61 versus other shapes constitute interesting and important specific instantiations
62 and sub problems of the general one that we try to address here.

63 5. *Implementable in a computational algorithm.* A successful theory of visual
64 object recognition needs to be described in sufficient detail to be implemented
65 through computational algorithms. This requirement is important because the
66 computational implementation allows us to run simulations and hence to
67 quantitatively compare the performance of the model against behavioral metrics.
68 The simulations also lend themselves to a direct comparison of the model's
69 computational steps and neurophysiological responses at different stages of the
70 visual processing circuitry. The algorithmic implementation forces us to rigorously
71 state the assumptions and formalize the computational steps; in this way,
72 computational models can be more readily compared than "armchair" theories
73 and models. The implementation can also help us debug the theory by
74 discovering hidden assumptions, bottlenecks and challenges that the algorithms
75 cannot solve or where performance is poor. There are multiple fascinating ideas
76 and theories about visual object recognition that have not been implemented
77 through computational algorithms. These ideas can be extremely useful and
78 helpful for the field and can inspire the development of computational models.
79 Yet, we emphasize that we cannot easily compare theories that can be and have
80 been implemented against other ones that have not.

81 6. *Restricted to primates.* Here we restrict the discussion to object
82 recognition in primates. There are strong similarities in visual object recognition
83 at the behavioral and neurophysiological levels between macaque monkeys (one
84 of the prime species for neurophysiological studies) and humans (e.g. (Myerson
85 et al., 1981; Logothetis and Sheinberg, 1996; Orban, 2004; Nielsen et al., 2006;
86 Kriegeskorte et al., 2008; Liu et al., 2009)).

87 7. *Biophysically plausible.* There are multiple computational approaches to
88 visual object recognition. Here we restrict the discussion to models that are
89 biophysically plausible. In doing so, we ignore a vast literature in Computer
90 Vision where investigators are trying to solve similar problems without direct
91 reference to the cortical circuitry. These engineering approaches are extremely
92 interesting and useful from a practical viewpoint. Ultimately, in the same way that

93 computers can become quite successful at playing chess without any direct
94 connection to the way humans play chess, computer vision approaches can
95 achieve high performance without mimicking neuronal circuits. Here we restrict
96 the discussion to biophysically plausible algorithms.

97 8. *Restricted to the visual system.* The visual system is not isolated from the
98 rest of the brain and there are plenty of connections between visual cortex and
99 other sensory cortices, between visual cortex and memory systems in the medial
100 temporal lobe and between the visual cortex and frontal cortex. It is likely that
101 these connections also play an important role in the process of visual recognition,
102 particularly through feedback signals that incorporate expectations (e.g. the
103 probability that there is a lion in an office setting is very small), prior knowledge
104 and experience (e.g. the object appears similar to another object that we are
105 familiar with), cross-modal information (e.g. the object is likely to be a musical
106 instrument because of the sound). To begin with and to simplify the problem, we
107 restrict the discussion to the visual system.

108 109 **8.2. Visual recognition goes beyond identifying objects in single images**

110
111 We emphasize that visual recognition is far more complex than the
112 identification of specific objects. Under natural viewing conditions, objects are
113 embedded in complex scenes and need to be separated from their background.
114 How this segmentation occurs constitutes an important challenge in itself.
115 Segmentation depends on a variety of cues including sharp edges, texture
116 changes and object motion among others. Some object recognition models
117 assume that segmentation must occur prior to recognition. There is no clear
118 biological evidence for segmentation prior to recognition and therefore this
119 constitutes a weakness in such approaches. We do not discuss segmentation
120 here (see (Borenstein et al., 2004; Sharon et al., 2006) for recent examples of
121 segmentation algorithms).

122 Most object recognition models are based on studying static images.
123 Under natural viewing conditions, there are important cues that depend on the
124 temporal integration of information. These dynamic cues can significantly
125 enhance recognition. Yet, it is clear that we can recognize objects in static
126 images and therefore many models focus on the reduced version the pattern
127 recognition problem using static objects. Here we also focus on static images.

128 We can perform a variety of complex tasks that rely on visual information
129 that are different from identification. For example, we can put together images of
130 snakes, lions and dolphins and categorize them as animals. Categorization is a
131 very important problem in vision research and it also constitutes a formidable
132 challenge for computer-based approaches. Here we focus on the question of
133 object identification.

134 135 **8.3. Modeling the ventral visual stream – Common themes**

136
137 Several investigators have proposed computational models that aim to
138 capture some of the essential principles behind the transformations along the

139 primate ventral visual stream. Before discussing some of those models in more
140 detail, we start by providing some common themes that are shared by many of
141 these models.

142 The input to the models is typically an image, defined by a matrix that
143 contains the grayscale value of each pixel. Object shapes can be discriminated
144 quite well in grayscale images and, therefore, most models ignore the added
145 complexities of color processing (but eventually it will also be informative and
146 important to add color to these models). Because the focus is often on the
147 computational properties of ventral visual cortex, several investigators often
148 ignore the complexities of modeling the computations in the retina and LGN; the
149 pixels are meant to coarsely represent the output of retinal ganglion cells or LGN
150 cells. This is of course one of the many oversimplifications in several
151 computation models given that we know that images go through a number of
152 transformations before retinal ganglion cells convey information to the LGN and
153 on to cortex (Meister, 1996).

154 Most models have a hierarchical and deep structure that aims to mimic the
155 approximately hierarchical architecture of ventral visual cortex (Felleman and
156 Van Essen, 1991; Maunsell, 1995). The properties of deep networks has
157 received considerable attention in the computational world, even if the
158 mathematics of learning in deep networks that include non-linear responses is far
159 less understood than shallow counterparts (Poggio and Smale, 2003). It seems
160 that neocortex and computer modelers have adopted a *Divide and Conquer*
161 strategy whereby a complex problem is divided into many simpler tasks.

162 Most computational models assume, explicitly or implicitly, that cortex is
163 cortex, and hence that there exist canonical microcircuits and computations that
164 are repeated over and over throughout the hierarchy (Riesenhuber and Poggio,
165 1999; Douglas and Martin, 2004; Serre et al., 2007b).

166 As we ascend through the hierarchical structure of the model, units in
167 higher levels typically have larger receptive fields, respond to more complex
168 visual features and show an increased degree of tolerance to transformations of
169 their preferred features.

170

171 **8.4. A panoply of models**

172

173 We summarize here a few important ideas that have been developed to
174 describe visual object recognition. The presentation here is neither an exhaustive
175 list nor a thorough discussion of each of these approaches. For a more detailed
176 discussion of several of these approaches, see (Ullman, 1996; LeCun et al.,
177 1998; Riesenhuber and Poggio, 2002; Deco and Rolls, 2004a; Serre et al.,
178 2005b).

179 Straightforward template matching does not work for pattern recognition.
180 Even shifting a pattern by one pixel would pose significant challenges for an
181 algorithm that merely compares the input with a stored pattern on a pixel-by-pixel
182 fashion. As noted at the beginning of this chapter, a key challenge to recognition
183 is that an object can lead to infinite number of retinal images depending on its
184 size, position, illumination, etc. If all objects were always presented in a

185 standardized position, scale, rotation and illumination, recognition would be
186 considerably easier (DiCarlo and Cox, 2007; Serre et al., 2007b). Based on this
187 notion, several approaches are based on trying to transform an incoming object
188 into a canonical prototypical format by shifting, scaling and rotating objects (e.g.
189 (Ullman, 1996)). The type of transformations required is usually rather complex,
190 particularly for non-affine transformations. While some of these problems can be
191 overcome by ingenious computational strategies, it is not entirely clear (yet) how
192 the brain would implement such complex calculations nor is there currently any
193 clear link to the type of neurophysiological responses observed in ventral visual
194 cortex.

195 A number of approaches are based on decomposing an object into its
196 component parts and their interactions. The idea behind this notion is that there
197 could be a small dictionary of object parts and a small set of possible interactions
198 that act as building blocks of all objects. Several of these ideas can be traced
199 back to the prominent work of David Marr (Marr and Nishihara, 1978; Marr, 1982)
200 where those constituent parts were based on generalized cone shapes. The
201 artificial intelligence community also embraced the notion of structural
202 descriptions (Winston, 1975). In the same way that a mathematical function can
203 be decomposed into a sum over a certain basis set (e.g. polynomials or sine and
204 cosine functions), the idea of thinking about objects as a sum over parts is
205 attractive because it may be possible and easier to detect these parts in a
206 transformation-invariant manner (Biederman, 1987; Mel, 1997). In the simplest
207 instantiations, these models are based on merely detecting a conjunction of
208 object parts, an approach that suffers from the fact that part rearrangements
209 would not impair recognition but they should (e.g. a house with a garage on the
210 roof and the chimney on the floor). More elaborate versions include part
211 interactions and relative positions. Yet, this approach seems to convert the
212 problem of object recognition to the problem of object part recognition plus the
213 problem of object parts interaction recognition. It is not entirely obvious that
214 object part recognition would be a trivial problem in itself nor is it obvious that *any*
215 object can be uniquely and succinctly described by a universal and small
216 dictionary of simpler parts. It is not entirely trivial how recognition of complex
217 shapes (e.g. consider discriminating between two faces) can be easily described
218 in terms of a structural description of parts and their interactions. Computational
219 implementations of these structural descriptions have been sparse (see however
220 (Hummel and Biederman, 1992)). More importantly, it is not entirely apparent
221 how these structural descriptions relate to the neurophysiology of the ventral
222 visual cortex (see however (Vogels et al., 2001)).

223 A series of computational algorithms, typically rooted in the neural network
224 literature (Hinton, 1992), attempt to build deep structures where inputs can be
225 reconstructed (for a recent version of this, see e.g. (Hinton and Salakhutdinov,
226 2006). In an extreme version of this approach, there is no information loss along
227 the deep hierarchy and backward signals carry information capable of re-creating
228 arbitrary inputs in lower visual areas. There are a number of interesting
229 applications for such “auto-encoder” deep networks such as the possibility of
230 performing dimensionality reduction. From a neurophysiological viewpoint,

231 however, it seems that the purpose of cortex is precisely the opposite, namely, to
232 lose information in biologically interesting ways. It is not clear why one build an
233 entire network to copy the input (possibly with fewer units). In other words, as
234 emphasized at the beginning of this chapter, it seems that a key goal of ventral
235 visual cortex is to be able to extract biologically relevant information (e.g. object
236 identity) in spite of changes in the input at the pixel level.

237 Particularly within the neurophysiology community, there exist several
238 “metric” approaches where investigators attempt to parametrically define a space
239 of shapes and then record the activity of neurons along the ventral visual stream
240 in response to these shapes (Tanaka, 1996; Brincat and Connor, 2004; Connor
241 et al., 2007). This dictionary of shapes can be more or less quantitatively defined.
242 For example, in some cases, investigators start by presenting different shapes in
243 search of a stimulus that elicits strong responses. Subsequently, they manipulate
244 the “preferred” stimulus by removing different parts and evaluating how the
245 neuronal responses are modified by these transformations. While interesting,
246 these approaches suffer from the difficulties inherent in considering arbitrary
247 shapes that may or may not constitute truly “preferred” stimuli. Additionally, in
248 some cases, the transformations examined only reveal anthropomorphic biases
249 about what features could be relevant. Another approach is to define shapes
250 parametrically. For example, Brincat and colleagues considered a family of
251 curvatures and modeled responses in a six-dimensional space defined by a sum
252 of Gaussians with parameters given by the curvature, orientation, relative
253 position and absolute position of the contour elements in the display. This
254 approach is intriguing because it has the attractive property of allowing
255 investigators to plot “tuning curves” similar to the ones used to represent the
256 activity of units in earlier visual areas. Yet, it also makes strong assumptions
257 about the type of shapes preferred by the units. Expanding on these ideas,
258 investigators have tried to start from generic shapes and use genetic algorithms
259 whose trajectories are guided by the neuronal preferences (Yamane et al., 2008).
260 What is particularly interesting about this approach is that it seems to be less
261 biased than the former two. The key limitation here is the recording time and this
262 type of algorithm, particularly with small data sets, may converge onto local
263 minima or even not converge at all. Genetic algorithms can be more thoroughly
264 examined in the computational domain. For example, investigators can examine
265 a huge variety of computational models and let them “compete” with each other
266 through evolutionary mechanisms (Pinto et al., 2009). To guide the evolutionary
267 paths, it is necessary to define a cost function; for example, evolution can be
268 constrained by rewarding models that achieve better performance in certain
269 recognition tasks. This type of approach can lead to models with high
270 performance (although it also suffers from difficulties related to local minima).
271 Unfortunately, it is not obvious that better performance necessarily implies any
272 better approximation to the way in which cortex solves the visual recognition
273 problem.

274

275 **8.5. Bottom-up hierarchical models of the ventral visual stream**

276

277 A hierarchical network model can be described by a series of layers
278 $i = 0, 1, \dots, N$. Each layer contains $n(i) \times n(i)$ units arranged in a matrix. The
279 activity of each unit in each layer can be represented by the matrix \mathbf{x}_i
280 ($\mathbf{x}_i \in \mathbb{R}^{n(i) \times n(i)}$). In several models, $x_i(j, k)$ (i.e., the activity of unit at position j, k in
281 layer i) is a scalar value interpreted as the firing rate of the unit. The initial layer is
282 defined as the input image; \mathbf{x}_0 represents the (grayscale) values of the pixels a
283 given image.

284 Equations 1 and 2 above constitute the initial steps for many object
285 recognition models and capitalize on the more studied parts of the visual system,
286 the pathway from the retina to primary visual cortex. The output of Equation 2,
287 after convolving the output of center-surround receptive fields with a Gabor
288 function, can be thought of as a first order approximation to the edges in the
289 image. As noted above, our understanding of ventral visual cortex beyond V1 is
290 far more primitive and it is therefore not surprising that this is where most models
291 diverge. In a first order simplification, we can generically describe the
292 transformations along the ventral visual stream as:

293
$$\mathbf{x}_{i+1} = f_i(\mathbf{x}_i) \qquad \text{Equation 11.1}$$

294 This assumes that the activity in a given layer only depends on the activity
295 pattern in the previous layer. This assumption implies that at least three types of
296 connections are ignored: (i) connections that “skip” a layer in the hierarchy (e.g.
297 synapses from the LGN to V4 skipping V1); (ii) top-down connections (e.g.
298 synapses from V2 to V1 (Virga, 1989)) and (iii) connections within a layer (e.g.
299 horizontal connections between neurons with similar preferences in V1 (Callaway,
300 1998)) are not included in **Equation 11.1**.

301 The subindex i in the function f indicates that the transformation from one layer to
302 another is not necessarily the same. A simple form that f could take is a linear
303 function:

304
$$\mathbf{x}_{i+1} = \mathbf{W}_i \mathbf{x}_i \qquad \text{Equation 11.2}$$

305 where the matrix \mathbf{W}_i represents the linear weights that transform activity in layer i
306 into activity in layer $i+1$. Not all neurons in layer i need to be connected to all
307 neurons in layer $i+1$; in other words, many entries in \mathbf{W} can be 0. This simple
308 formulation finds some empirical evidence; for example, Hubel and Wiesel
309 proposed that the oriented filters in primary visual cortex could arise from a linear
310 summation of the activity of neurons in the lateral geniculate nucleus with
311 appropriately aligned center-surround receptive fields (Hubel and Wiesel, 1962).
312 Unfortunately, neurons are far more intricate devices and non-linearities abound
313 in their response properties. For example, Hubel and Wiesel also described the
314 activity of so-called complex cells that are also orientation tuned but show a non-
315 linear response as a function of spatial frequency or bar length.

316 It is tacitly assumed by most modelers that there exist general rules, often
317 summarized in the epithet “cortex is cortex”, such that only a few such
318 transformations are allowed. One of the early models that aimed to describe
319 object recognition, inspired by the neurophysiological findings of Hubel and
320 Wiesel, was the neocognitron (Fukushima, 1980) (see also earlier theoretical
321 ideas in (Sutherland, 1968)). This model had two possible operations, a linear
322 tuning function (performed by “simple” cells) and a non-linear OR operation

323 (performed by “complex” cells). These two operations were alternated and
324 repeated through the multiple layers in the deep hierarchy. This model
325 demonstrated the feasibility of such linear/non-linear cascades in achieving scale
326 and position tolerance in a letter recognition task. Several subsequent efforts
327 (Olshausen et al., 1993; Wallis and Rolls, 1997; LeCun et al., 1998; Riesenhuber
328 and Poggio, 1999; Amit and Mascaró, 2003; Deco and Rolls, 2004b) were
329 inspired by the Neocognitron architecture.

330 One such effort to expand on the computational abilities of the
331 Neocognitron in the computational model developed in the Poggio group
332 (Riesenhuber and Poggio, 1999; Serre et al., 2005b; Serre et al., 2007b). This
333 model is characterized by a purely feed-forward and hierarchical architecture. An
334 image, represented by grayscale values, is convolved with Gabor filters at
335 multiple scales and positions to mimic the responses of simple cells in primary
336 visual cortex. Like other efforts, the model consists of a cascade of linear and
337 non-linear operations. These operations come in only two flavors in the model: a
338 tuning operation and soft-max operation. Both operations can be expressed in
339 the following form:

$$340 \quad x_{i+1}[k] = g \left(\frac{\sum_{j=1}^n w[j,k] x_i^p[j]}{\alpha + \left(\sum_{j=1}^n x_i^q[j] \right)^r} \right) \quad \text{Equation 11.3}$$

341 where $x_{i+1}[k]$ represents the activity of unit k in layer $i+1$, $w[j,k]$ represents the
342 connection weight between unit j in layer i and unit k in layer $b+1$, p , q , r are
343 integer constants, a is a normalization constant and g is a nonlinear function (e.g.
344 sigmoid). Depending on the values of p , q and r different interesting behaviors
345 can be obtained. In particular, taking $r=1/2$, $p=1$, $q=2$, leads to a *tuning operation*:

$$346 \quad x_{i+1}[k] = g \left(\frac{\sum_{j=1}^n w[j,k] x_i[j]}{\alpha + \sqrt{\sum_{j=1}^n x_i^2[j]}} \right) \quad \text{Equation 11.3'}$$

347 The responses of the unit are controlled by the weights \mathbf{w} . As emphasized above,
348 tuning is a central aspect of any computational model of visual recognition,
349 allowing units along the hierarchy to respond to increasingly more elaborate
350 features. Taking $\mathbf{w}=1$, $p=q+1$, $r=1$, leads to a softmax operation, particularly for
351 large values of q :

$$352 \quad x_{i+1}[k] = g \left(\frac{\sum_{j=1}^n x_i^{q+1}[j]}{\alpha + \sum_{j=1}^n x_i^q[j]} \right) \quad \text{Equation 11.3''}$$

353 Of note, the unit with response $x_{i+1}[k]$ shows similar response tuning to the units
354 with response $x_i[j]$ for $j=1, \dots, n$. Yet, the higher-level unit shows a stronger

355 degree of tolerance to those aspects of the response that differentiate the
356 responses of different units with similar tuning in layer i . For example, different
357 units in layer i may show identical feature preferences but have slightly different
358 receptive fields. A winner-take-all unit in layer $i+1$ that takes as input the
359 responses of those units would inherit the same feature preferences but would
360 reveal a larger receptive and tolerate changes in the position of the feature within
361 the larger receptive field. Both operations can be implemented through relatively
362 simple biophysical circuits (Kouh and Poggio, 2004).

363 This and similar architectures have been subjected to several tests
364 including comparison with psychophysical measurements (e.g. (Serre et al.,
365 2007a)), comparison with neurophysiological responses (e.g. (Deco and Rolls,
366 2004b; Lampl et al., 2004; Hung et al., 2005; Serre et al., 2005b; Cadieu et al.,
367 2007) and quantitative evaluation of performance in computer vision recognition
368 tasks (e.g. (LeCun et al., 1998; Serre et al., 2005a; Mutch and Lowe, 2006)).

369

370 **8.6. Top-down signals in visual recognition**

371

372 In spite of the multiple simplifications, the success of bottom-up
373 architecture in describing a large number of visual recognition phenomena
374 suggest that they may not be a bad first cut. As emphasized above, bottom-up
375 architectures constitute only an approximation to the complexities and wonders
376 of neocortical computation. One of the several simplifications in bottom-up
377 models is the lack of top-down signals. We know that there are abundant back-
378 projections in neocortex (e.g. (Felleman and Van Essen, 1991; Callaway, 2004;
379 Douglas and Martin, 2004)). The functions of top-down connections have been
380 less studied at the neurophysiological level but there is no shortage of
381 computational models illustrating the rich array of computations that emerge with
382 such connectivity. Several models have used top-down connections to guide
383 attention to specific locations or specific features within the image (e.g.
384 (Olshausen et al., 1993; Itti and Koch, 2001))(Tsotsos, 1990; Deco and Rolls,
385 2005; Rao, 2005; Compte and Wang, 2006; Chikkerur et al., 2009). The
386 allocation of attention to specific parts of an image can significantly enhance
387 recognition performance by alleviating the problems associated with image
388 segmentation and with clutter.

389 Top-down signals can also play an important role in recognition of
390 occluded objects. When only partial object information is available, the system
391 must be able to perform object completion and interpret the image based on prior
392 knowledge. Attractor networks have been shown to be able to retrieve the
393 identity of stored memories from partial information (e.g. (Hopfield, 1982)). Some
394 computational models have combined bottom-up architectures with attractor
395 networks at the top of the hierarchy (e.g. (Deco and Rolls, 2004b)).

396 During object completion, top-down signals could play an important role by
397 providing prior stored information that influences the bottom-up sensory
398 responses. Several proposals have argued that visual recognition can be
399 formulated as a Bayesian inference problem (Mumford, 1992; Rao et al., 2002;
400 Lee and Mumford, 2003; Rao, 2004; Yuille and Kersten, 2006; Chikkerur et al.,

401 2009). Considering three layers of the visual cascade (e.g. LGN, V1 and higher
402 areas), and denoting activity in those layers as \mathbf{x}_0 , \mathbf{x}_1 and \mathbf{x}_h respectively, then
403 the probability of obtaining a given response pattern in V1 depends both on the
404 sensory input as well as feedback from higher areas:

$$405 \quad P(\mathbf{x}_1|\mathbf{x}_0) = \frac{P(\mathbf{x}_0|\mathbf{x}_1)P(\mathbf{x}_1|\mathbf{x}_h)}{P(\mathbf{x}_0|\mathbf{x}_h)} \quad \text{Equation 11.8}$$

406 where $P(\mathbf{x}_1|\mathbf{x}_h)$ represents the feedback biases conveying prior information. An
407 intriguing idea proposed by Rao and Ballard argues that top-down connections
408 provide predictive signals whereas bottom-up signals convey the difference
409 between the sensory input and the top-down predictions (Rao and Ballard, 1999).

410

411 **8.7. Applications**

412

413 There is no shortage of applications where automatic or semi-automatic
414 algorithms are being explored in computer vision. Here are a few examples:

415 (A) Intelligent content-based search. Searching for images in the web by
416 content will open the doors to a large number of applications. Facebook
417 users can already experiment with prototypes that let them search for
418 people. One may be able to look for images that are similar to a search
419 query in terms of content. Blind people may be able to point their phones
420 and find out where they are and how to navigate.

421 (B) Prototype cars that can navigate automatically rely heavily on algorithms
422 to detect pedestrians, other cars, other objects and road conditions.

423 (C) ATM machines may be able to recognize their customers. Cars and
424 houses may recognize their owners.

425 (D) Security screening in places like airports may benefit from automatic
426 recognition systems.

427 (E) Several clinical problems are based on pattern recognition and
428 computers may soon help doctors to make informed decisions based on
429 their understanding of patterns.

430

431 **8.8. Computer vision tasks**

432

433 Algorithms have been developed to address several interrelated problems
434 in machine vision. While some of the boundaries are blurred in several
435 applications, it is useful to think about the following tasks:

436 (A) Object detection. For example, a digital camera may require detecting the
437 presence and location of a face in an image for focusing. Face detection
438 may thrive without solving the problem of recognition.

439 (B) Object segmentation. In natural images, it may be of interest to separate
440 an object from the background. For example, it may be important to
441 detect the location of a tumor in an image. Or to detect the presence of a
442 tank in a camouflaged image.

443 (C) Object recognition. Recognizing objects can often be thought of as
444 associating the image with labels. These labels may refer to the identity of

445 the object (e.g. given a face, who is it?) or the object's category (e.g. is
446 there an animal in this image?).

447 (D) Object verification. In some cases, it may be of interest to evaluate
448 whether two images are the same or not.

449

450 **8.9. Object segmentation**

451

452 Given a natural scene, humans (and other species) are quite good at
453 being able to characterize and localize different objects embedded in complex
454 backgrounds. The fact that this is not a trivial problem is highlighted by the
455 ubiquitous use of camouflage in the animal world. Particularly for objects that are
456 not moving, matching colors, contrast and textures can help animals avoid
457 predators or at least buy sufficient time for escape. Basic aspects of
458 segmentation may depend on adequately detecting edges. However, more
459 complex problems often involve a deeper understanding of the interrelationships
460 among different object parts. A typical case involves recognizing a zebra as a
461 whole animal as opposed to thinking of each stripe as a separate object.

462 Some algorithms require recognition prior to segmentation while other
463 algorithms use segmentation to guide recognition in complex scenes. To avoid
464 this chicken-and-egg dilemma, it is tempting to speculate that certain aspects of
465 bottom-up recognition and segmentation could occur (or at least) start
466 independently of each other, using overlapping neuronal circuits. Top-down
467 signals may then combine segmentation and recognition in synergistic fashion.
468 For examples of object segmentation algorithms see (Borenstein et al., 2004).

469

470 **8.10. A general scheme for object recognition**

471

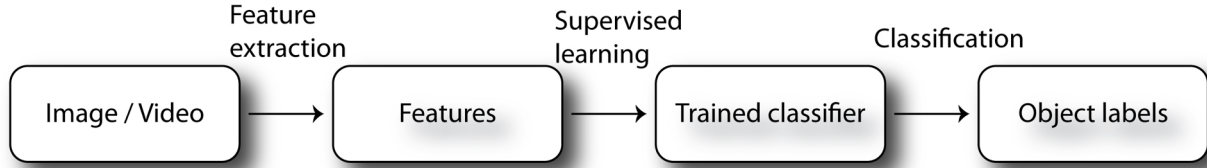
472 Figure 12.1 illustrates a typical approach in computer vision efforts.
473 Consider a series of N labeled images (\mathbf{x}_i, y_i) where $i=1, \dots, N$, \mathbf{x} is a matrix
474 representing the image and y is a label (e.g. face present or not). A set of
475 features \mathbf{f} is extracted from the images: $\mathbf{f}_i = g(\mathbf{x}_i)$. Those features may include
476 properties such as edges, principal components, etc. How those features are
477 chosen is one of the key aspects that differentiates computer vision algorithms. A
478 supervised learning scheme is then used to learn the map between those
479 features and labels (Poggio and Smale, 2003; Meyers and Kreiman, 2011;
480 Singer and Kreiman, 2011). For example, a support vector machine (SVM)
481 classifier with a linear kernel may be used to learn the structure of the data and
482 labels. A cross-validation procedure is followed by separating the data into a
483 training set and a test set to avoid overfitting. After training, the algorithm is
484 evaluated with the images in the test set. By using different algorithms applied to
485 the same data, the merits of alternative approaches can be quantitatively
486 compared.

487

488 **8.11. A successful example: digit recognition**

489

Figure 12.1. A general scheme for object recognition. Features are extracted from an image (or video). Those features are used to train a classifier via supervised learning. The resulting classification boundary is used with novel images (different from the ones used during training) to assign object labels to images.



490
491
492
493
494
495
496
497

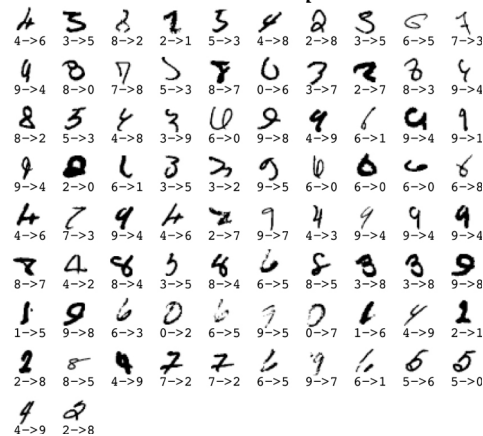
Recognizing hand-written digits constitutes an example where computers have reached high accuracy, almost comparable to human levels (e.g. (LeCun et al., 1998)). Figure 12.2 shows an example of the errors made by an early attempt at recognizing hand-written digits. The overall error rate of this algorithm was <2%. Several of those errors are not trivial to recognize and humans could make mistakes as well.

498 **8.12. Image recognition competitions**

499
500
501
502
503
504
505
506
507
508

There are several computer vision competitions with large data sets consisting of labeled images. One such competition is called the Imagenet Large Scale Visual Recognition Challenge (ILSVRC) (Russakovsky et al., 2014). The 2014 instantiation of the object classification part of this challenge included 1000 object classes, 1,281,413 images for training (732-1300 images per class) and 100,000 images for testing (100 images per class). This competition also includes other tasks beyond classification including object detection and localization. To give an idea of performance, the winning team in the object classification part of the challenge achieved an error rate slightly above 6%. This

Figure 12.2. Example of digit recognition mistakes by the algorithm in LeCun et al 1988. Below each digit, the image shows the true label and the computer label.



509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524

is quite impressive considering that there were 1000 classes. It should be noted, though, that the results of these competitions are often reported in a somewhat strange way by allowing the models 5 changes to get it right and reporting the results as correct if any of those 5 predictions are correct. This makes it a bit more difficult to directly compare against human performance (Borji and Itti, 2014). Another aspect of machine vision that has also been highlighted is the difficulty in interpreting how the machine classifies objects and

525 investigators have reported puzzling examples where minimal changes to an
526 image drastically change the predicted class (Szegedy et al., 2014). With that
527 said, the results are still quite remarkable and they show rapid progress in
528 teaching machines to recognize objects.

529

530 **8.7. The road ahead**

531 If you can do it, a computer can do it too. Significant progress has been made
532 over the last decade in teaching computers to perform multiple tasks that were
533 traditionally thought to be the domain of humans. Any desktop computer can play
534 chess competitively and the best computers can beat the world's chess
535 champion. IBM's Watson has thrived in the trivia-like game of Jeopardy. And
536 while imperfect, Siri and related systems are making enormous strides in
537 becoming the world's best assistants.

538

539 In the domain of vision, computational algorithms are already able to perform
540 certain tasks such as recognizing digits in a fully automatic fashion at human
541 performance level and demonstrate reasonable performance in other tasks such
542 as detecting faces for focusing on in digital cameras. In several other tasks,
543 humans still outperform the most sophisticated current algorithms but the gap
544 between machines and humans in vision tasks is closing rapidly. Here we
545 provide an overview of several computer vision systems, particularly in the
546 context of pattern recognition problems and describe what machine vision
547 systems can and cannot do,

548

549 Significant progress has been made towards describing visual object
550 recognition in a principled and theoretically sound fashion. Yet, the lacunas in our
551 understanding of the functional and computational architecture of ventral visual
552 cortex are not small. The preliminary steps have distilled important principles of
553 neocortical computation including deep networks that can divide and conquer
554 complex tasks, bottom-up circuits that perform rapid computations, gradual
555 increases in selectivity and tolerance to object transformation.

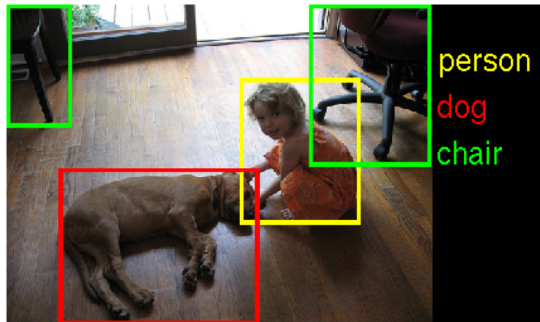
556

557 In stark contrast with the pathway from the retina to primary visual cortex,
558 we do not have a quantitative description of the feature preferences of neurons
559 along the ventral visual pathway. And several computational models do not make
560 clear, concrete and testable predictions towards systematically characterizing
561 ventral visual cortex at the physiological levels. Computational models can
562 perform several complex recognition tasks and compete against non-biological
563 computer vision approaches. Yet, for the vast majority of recognition tasks, they
564 still fall significantly below human performance.

565

566 The next several years are likely to bring many new surprises in the field.
567 We will be able to characterize the system at unprecedented resolution at the
568 experimental level and we will be able to evaluate sophisticated and
569 computationally intensive theories in realistic times. In the same way that the
570 younger generations are not surprised by machines that can play chess quite

Figure 12.3. Example labeled image from
the validation set in the 2013 ILSVRC
competition. [Image source:
[http://www.image-
net.org/challenges/LSVRC/2013/](http://www.image-net.org/challenges/LSVRC/2013/)]



competitively, the next generation
may not be surprised by intelligent
devices that can “see” like we do.

8.13. References

- 586
587
588 image understanding. *Psychological Review* 24:115-147.
589 Borenstein E, Sharon E, Ullman S (2004) Combining Top-Down and Bottom-Up
590 Segmentation. In: *IEEE Conference on Computer Vision and Pattern*
591 *Recognition (CVPR)*. Washington, DC.
592 Borji A, Itti L (2014) Human vs. computer in scene and object recognition. In: *CVPR*.
593 Brady TF, Konkle T, Alvarez GA, Oliva A (2008) Visual long-term memory has a
594 massive storage capacity for object details. *Proceedings of the National*
595 *Academy of Science U S A* 105:14325-14329.
596 Brincat SL, Connor CE (2004) Underlying principles of visual shape selectivity in
597 posterior inferotemporal cortex. *Nature Neuroscience* 7:880-886.
598 Cadieu C, Kouh M, Pasupathy A, Connor C, Riesenhuber M, Poggio T (2007) A model
599 of V4 shape selectivity and invariance. *Journal of Neurophysiology* 98:1733-
600 1750.
601 Callaway EM (1998) Local circuits in primary visual cortex of the macaque monkey.
602 *Annu Rev Neurosci* 21:47-74.
603 Callaway EM (2004) Feedforward, feedback and inhibitory connections in primate
604 visual cortex. *Neural Netw* 17:625-632.
605 Chikkerur S, Serre T, Poggio T (2009) A Bayesian inference theory of attention:
606 neuroscience and algorithms. In: (MIT-CSAIL-TR, ed). Cambridge: MIT.
607 Compte A, Wang XJ (2006) Tuning curve shift by attention modulation in cortical
608 neurons: a computational study of its mechanisms. *Cerebral cortex* 16:761-
609 778.
610 Connor CE, Brincat SL, Pasupathy A (2007) Transformation of shape information in
611 the ventral pathway. *Current opinion in neurobiology* 17:140-147.
612 Deco G, Rolls ET (2004a) *Computational Neuroscience of Vision*. Oxford Oxford
613 University Press.
614 Deco G, Rolls ET (2004b) A neurodynamical cortical model of visual attention and
615 invariant object recognition. *Vision research* 44:621-642.

- 616 Deco G, Rolls ET (2005) Attention, short-term memory, and action selection: a
617 unifying theory. *Prog Neurobiol* 76:236-256.
- 618 DiCarlo JJ, Cox DD (2007) Untangling invariant object recognition. *Trends Cogn Sci*
619 11:333-341.
- 620 Douglas RJ, Martin KA (2004) Neuronal circuits of the neocortex. *Annu Rev Neurosci*
621 27:419-451.
- 622 Felleman DJ, Van Essen DC (1991) Distributed hierarchical processing in the
623 primate cerebral cortex. *Cerebral cortex* 1:1-47.
- 624 Fukushima K (1980) Neocognitron: a self organizing neural network model for a
625 mechanism of pattern recognition unaffected by shift in position. *Biological*
626 *Cybernetics* 36:193-202.
- 627 Hinton G (1992) How neural networks learn from experience. *Scientific American*
628 267:145-151.
- 629 Hinton GE, Salakhutdinov RR (2006) Reducing the dimensionality of data with
630 neural networks. *Science* 313:504-507.
- 631 Hopfield JJ (1982) Neural networks and physical systems with emergent collective
632 computational abilities. *PNAS* 79:2554-2558.
- 633 Hubel DH, Wiesel TN (1962) Receptive fields, binocular interaction and functional
634 architecture in the cat's visual cortex. *The Journal of physiology* 160:106-154.
- 635 Hummel JE, Biederman I (1992) Dynamic binding in a neural network for shape
636 recognition. *Psychol Rev* 99:480-517.
- 637 Hung CP, Kreiman G, Poggio T, DiCarlo JJ (2005) Fast Read-out of Object Identity
638 from Macaque Inferior Temporal Cortex. *Science* 310:863-866.
- 639 Itti L, Koch C (2001) Computational modelling of visual attention. *Nat Rev Neurosci*
640 2:194-203.
- 641 Keysers C, Xiao DK, Foldiak P, Perret DI (2001) The speed of sight. *Journal of*
642 *Cognitive Neuroscience* 13:90-101.
- 643 Kirchner H, Thorpe SJ (2006) Ultra-rapid object detection with saccadic eye
644 movements: visual processing speed revisited. *Vision research* 46:1762-1776.
- 645 Kouh M, Poggio T (2004) A general mechanism for tuning: gain control circuits and
646 synapses underlie tuning of cortical neurons. In: (Memo MA, ed). Cambridge:
647 MIT.
- 648 Kriegeskorte N, Mur M, Ruff DA, Kiani R, Bodurka J, Esteky H, Tanaka K, Bandettini
649 PA (2008) Matching categorical object representations in inferior temporal
650 cortex of man and monkey. *Neuron* 60:1126-1141.
- 651 Lampl I, Ferster D, Poggio T, Riesenhuber M (2004) Intracellular measurements of
652 spatial integration and the MAX operation in complex cells of the cat primary
653 visual cortex. *J Neurophysiol* 92:2704-2713.
- 654 LeCun Y, Bottou L, Bengio Y, Haffner P (1998) Gradient-based learning applied to
655 document recognition. *Proc of the IEEE* 86:2278-2324.
- 656 Lee TS, Mumford D (2003) Hierarchical Bayesian inference in the visual cortex. *J Opt*
657 *Soc Am A Opt Image Sci Vis* 20:1434-1448.
- 658 Liu H, Agam Y, Madsen JR, Kreiman G (2009) Timing, timing, timing: Fast decoding
659 of object information from intracranial field potentials in human visual
660 cortex. *Neuron* 62:281-290.

- 661 Logothetis NK, Sheinberg DL (1996) Visual object recognition. Annual Review of
662 Neuroscience 19:577-621.
- 663 Marr D (1982) Vision. San Francisco: Freeman publishers.
- 664 Marr D, Nishihara HK (1978) Representation and recognition of the spatial
665 organization of three-dimensional shapes. Proc R Soc Lond B Biol Sci
666 200:269-294.
- 667 Maunsell JHR (1995) The brain's visual world: representation of visual targets in
668 cerebral cortex. Science 270:764-769.
- 669 Meister M (1996) Multineuronal Codes in Retinal Signaling. PNAS 93:609-614.
- 670 Mel B (1997) SEEMORE: Combining color, shape and texture histogramming in a
671 neurally inspired approach to visual object recognition. Neural Computation
672 9:777.
- 673 Meyers EM, Kreiman G (2011) Tutorial on Pattern Classification in Cell Recordings.
674 In: Understanding visual population codes (Kriegeskorte N, Kreiman G, eds).
675 Boston: MIT Press.
- 676 Mumford D (1992) On the computational architecture of the neocortex. II. The role
677 of cortico-cortical loops. Biol Cybern 66:241-251.
- 678 Mutch J, Lowe D (2006) Multiclass Object Recognition with Sparse, Localized
679 Features. In: CVPR, pp 11-18. New York.
- 680 Myerson J, Miezin F, Allman J (1981) Binocular rivalry in macaque monkeys and
681 humans: a comparative study in perception. Behavioral Analysis Letters
682 1:149-159.
- 683 Nielsen KJ, Logothetis NK, Rainer G (2006) Discrimination strategies of humans and
684 rhesus monkeys for complex visual displays. Current Biology 16:814-820.
- 685 Olshausen BA, Anderson CH, Van Essen DC (1993) A neurobiological model of visual
686 attention and invariant pattern recognition based on dynamic routing of
687 information. Journal of Neuroscience 13:4700-4719.
- 688 Orban GA, Van Essen, D., Vanduffel, W. (2004) Comparative mapping of higher visual
689 areas in monkeys and humans. Trends in Cognitive Sciences 8:315-324.
- 690 Pinto N, Doukhan D, DiCarlo JJ, Cox DD (2009) A high-throughput screening
691 approach to discovering good forms of biologically inspired visual
692 representation. PLoS Comput Biol 5:e1000579.
- 693 Poggio T, Smale S (2003) The mathematics of learning: dealing with data. Notices of
694 the AMS 50:537-544.
- 695 Potter M, Levy E (1969) Recognition memory for a rapid sequence of pictures.
696 Journal of experimental psychology 81:10-15.
- 697 Rao RP (2004) Bayesian computation in recurrent neural circuits. Neural Comput
698 16:1-38.
- 699 Rao RP (2005) Bayesian inference and attentional modulation in the visual cortex.
700 Neuroreport 16:1843-1848.
- 701 Rao RP, Ballard DH (1999) Predictive coding in the visual cortex: a functional
702 interpretation of some extra-classical receptive-field effects. Nature
703 neuroscience 2:79-87.
- 704 Rao RPN, Olshausen BA, Lewicki MS, eds (2002) Probabilistic Models of the Brain:
705 Perception and Neural Function. Cambridge: MIT Press.

- 706 Richmond B, Wurtz R, Sato T (1983) Visual responses in inferior temporal neurons
707 in awake Rhesus monkey. *Journal of Neurophysiology* 50:1415-1432.
- 708 Riesenhuber M, Poggio T (1999) Hierarchical models of object recognition in cortex.
709 *Nature Neuroscience* 2:1019-1025.
- 710 Riesenhuber M, Poggio T (2002) Neural mechanisms of object recognition. *Current*
711 *opinion in neurobiology* 12:162-168.
- 712 Rolls E (1991) Neural organization of higher visual functions. *Current opinion in*
713 *neurobiology* 1:274-278.
- 714 Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang S, Karpathy A, Khosla
715 A, Bernstein M, Berg A, Fei-Fei L (2014) ImageNet Large Scale Visual
716 Recognition Challenge. In: CVPR: arXiv:1409.0575, 2014.
- 717 Serre T, Wolf L, Poggio T (2005a) Object Recognition with Features Inspired by
718 Visual Cortex. In: IEEE Computer Society Conference on Computer Vision and
719 Pattern Recognition (CVPR). San Diego: IEEE Computer Society Press.
- 720 Serre T, Oliva A, Poggio T (2007a) Feedforward theories of visual cortex account for
721 human performance in rapid categorization. *PNAS* 104:6424-6429.
- 722 Serre T, Kouh M, Cadieu C, Knoblich U, Kreiman G, Poggio T (2005b) A theory of
723 object recognition: computations and circuits in the feedforward path of the
724 ventral stream in primate visual cortex. In, pp CBCL Paper #259/AI Memo
725 #2005-2036. Boston: MIT.
- 726 Serre T, Kreiman G, Kouh M, Cadieu C, Knoblich U, Poggio T (2007b) A quantitative
727 theory of immediate visual recognition. *Progress In Brain Research* 165C:33-
728 56.
- 729 Sharon E, Galun M, Sharon D, Basri R, Brandt A (2006) Hierarchy and adaptivity in
730 segmenting visual scenes. *Nature* 442:810-813.
- 731 Singer J, Kreiman G (2011) Introduction to Statistical Learning and Pattern
732 Classification. In: *Visual Population Codes* (Kriegeskorte N, Kreiman G, eds).
733 Boston: MIT Press.
- 734 Standing L (1973) Learning 10,000 pictures. *Quarterly Journal of Experimental*
735 *Psychology* 25:207-222.
- 736 Sutherland NS (1968) Outlines of a theory of visual pattern recognition in animals
737 and man. *Proc R Soc Lond B Biol Sci* 171:297-317.
- 738 Szegedy C, Zaremba W, Sutskever I, Bruna J, Erhan D, Goodfellow I, Fergus R (2014)
739 Intriguing properties of neural networks. In: *International Conference on*
740 *Learning Representations*.
- 741 Tanaka K (1996) Inferotemporal cortex and object vision. *Annual Review of*
742 *Neuroscience* 19:109-139.
- 743 Thorpe S, Fize D, Marlot C (1996) Speed of processing in the human visual system.
744 *Nature* 381:520-522.
- 745 Tsotsos J (1990) Analyzing Vision at the Complexity Level. *Behavioral and Brain*
746 *Sciences* 13-3:423-445.
- 747 Ullman S (1996) *High-Level Vision*. Cambridge, MA: The MIT Press.
- 748 Virga A, Rockland, KS (1989) Terminal Arbors of Individual "Feedback" Axons
749 Projecting from Area V2 to V1 in the Macaque Monkey: A Study Using
750 Immunohistochemistry of Anterogradely Transported Phaseolus vulgaris-
751 leucoagglutinin. *The Journal of Comparative Neurology* 285:54-72.

- 752 Vogels R, Biederman I, Bar M, Lorincz A (2001) Inferior temporal neurons show
753 greater sensitivity to nonaccidental than to metric shape differences. *J Cogn*
754 *Neurosci* 13:444-453.
- 755 Wallis G, Rolls ET (1997) Invariant face and object recognition in the visual system.
756 *PROGRESS IN NEUROBIOLOGY* 51:167-194.
- 757 Winston P (1975) Learning structural descriptions from examples. In: *The*
758 *psychology of computer vision* (Winston P, ed), pp 157-209. London:
759 McGraw-Hill.
- 760 Yamane Y, Carlson ET, Bowman KC, Wang Z, Connor CE (2008) A neural code for
761 three-dimensional object shape in macaque inferotemporal cortex. *Nature*
762 *neuroscience* 11:1352-1360.
- 763 Yuille A, Kersten D (2006) Vision as Bayesian inference: analysis by synthesis?
764 *Trends Cogn Sci* 10:301-308.
765