

available at www.sciencedirect.comwww.elsevier.com/locate/brainres**BRAIN
RESEARCH**

Research Report

Why is the fusiform face area recruited for novel categories of expertise? A neurocomputational investigationMatthew H. Tong^a, Carrie A. Joyce^b, Garrison W. Cottrell^{a,*}^aComputer Science and Engineering, University of California, San Diego, 9500 Gilman Dr. La Jolla, CA 92093-0114, USA^bID Analytics, Inc., San Diego, CA, USA

ARTICLE INFO

Article history:

Accepted 18 June 2007

Keywords:

Fusiform face area

Face processing

Visual expertise

Computational modeling

ABSTRACT

What is the role of the Fusiform Face Area (FFA)? Is it specific to face processing, or is it a visual expertise area? The expertise hypothesis is appealing due to a number of studies showing that the FFA is activated by pictures of objects within the subject's domain of expertise (e.g., cars for car experts, birds for birders, etc.), and that activation of the FFA increases as new expertise is acquired in the lab. However, it is incumbent upon the proponents of the expertise hypothesis to explain how it is that an area that is initially specialized for faces becomes recruited for new classes of stimuli. We dub this the "visual expertise mystery." One suggested answer to this mystery is that the FFA is used simply because it is a fine discrimination area, but this account has historically lacked a mechanism describing exactly how the FFA would be recruited for novel domains of expertise. In this study, we show that a neurocomputational model trained to perform subordinate-level discrimination within a visually homogeneous class develops transformations that magnify differences between similar objects, in marked contrast to networks trained to simply categorize the objects. This magnification generalizes to novel classes, leading to faster learning of new discriminations. We suggest this is why the FFA is recruited for new expertise. The model predicts that individual FFA neurons will have highly variable responses to stimuli within expertise domains.

© 2007 Elsevier B.V. All rights reserved.

1. Introduction

There has been a great deal of progress in understanding how complex objects, in particular, human faces, are processed by the cortex. At the same time, there is a great deal of controversy about the role of various cortical areas, especially the Fusiform Face Area (FFA) (Kanwisher et al., 1997; Kanwisher, 2000; Tarr and Gauthier, 2000). Is the FFA a "module," specific to the domain of faces, or is it instead specific to the process of fine level discrimination? Several fMRI studies showed high acti-

vation in the FFA only to face stimuli and not other objects (Kanwisher et al., 1997; Kanwisher, 2000). Furthermore, studies involving patients with associative prosopagnosia, the inability to identify individual faces (Farah et al., 1995), and visual object agnosia, the inability to recognize non-face objects (Moscovitch et al., 1997), seem to indicate a clear double dissociation between face and object processing. Prosopagnosic patients had lesions encompassing either right hemisphere or bilateral FFA, while object agnostic patients' lesions did not (De Renzi et al., 1994).

* Corresponding author. CSE Department 0404, UCSD, La Jolla, CA 92093-0404, USA. Fax: +1 858 534 7029.
E-mail address: gary@ucsd.edu (G.W. Cottrell).

Gauthier et al. (1997, 1999a) have challenged the notion of the face specificity of the FFA by pointing out that the earlier studies failed to equate the level of experience subjects had with non-face objects with the level of experience they had with faces. Gauthier et al. (2000) showed that the FFA was activated when car and bird experts were shown pictures of the animals in their area of expertise. Furthermore, they illustrated that, if properly trained, individuals can develop expertise on novel, non-face objects (e.g., “Greebles”), and subsequently show increased FFA activation to them (Gauthier and Tarr, 1997; Gauthier et al., 1999b). Crucially, the same 2 or 3 voxels that are most active for faces also show the largest increase in activity over the course of expertise training on non-face stimuli, suggesting that the FFA is recruited as subjects learn to visually discriminate novel homogeneous stimuli, and is automatically engaged when the subject is an expert (Tarr and Gauthier, 2000). Hence the theory is that the FFA is a *fine level discrimination area* (this is still controversial - see Grill-Spector et al. (2004) and Rhodes et al. (2004) for competing evidence). However, the idea that the FFA is a fine level discrimination area still does not answer the question of what mechanism would explain how an area that presumably starts life as a face processing region is recruited for these other types of stimuli. This is a job for modeling.

Before addressing this question, it is important to define the notion of an “expert.” We use Gauthier’s operational definition

of the term: experts are as fast to verify that a picture of an object is a particular individual (subordinate level) as they are to verify their category membership (basic level). For example, a bird expert would be as fast and as accurate at verifying that a picture of a bird is an “Indigo Bunting” as at identifying it as a “bird.” On the other hand, a novice will show the fastest reaction time at the basic level, and is slower at both subordinate and superordinate level (Tanaka and Taylor, 1991). The basic level was first identified by Rosch as the level at which objects tend to share the same shape and function, and tend to correspond to the first word we use to describe an object (a picture of a chair is labeled “chair” rather than “furniture” or “office chair”). When training a subject in a novel category, the downward shift in reaction times in these two tasks is taken as evidence of expertise.

Previously, we have demonstrated that developmentally appropriate conditions (low spatial frequency input and learning subordinate/individual level classification) are sufficient for our neurocomputational model to specialize for faces (Dailey and Cottrell, 1999). Here, we investigate what properties the FFA might possess that would result in its recruitment for non-face, subordinate level discrimination tasks.

We compare the properties of two kinds of cortical models: “expert networks” trained to make subordinate level categorizations (“Is this Bob, Carol, Ted or Alice?”, top path of Fig. 1), and “basic networks” trained to make category level classifications (“Is this a face, cup, can, or book?”, bottom path of Fig. 1)

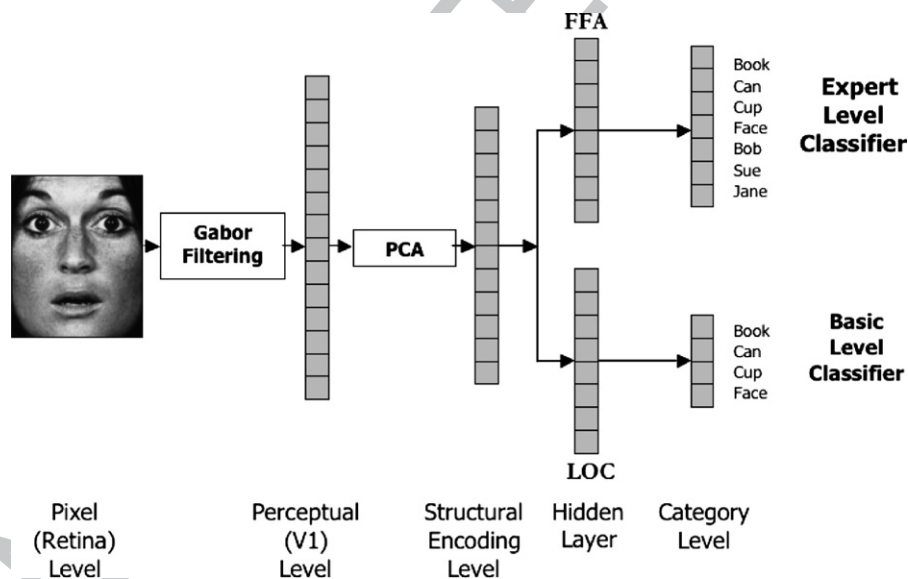


Fig. 1 – Network architecture. Input images are 64×64 grayscale images. The first layer of processing consists of Gabor filters (wavelets) at 8 different orientations ($0, \pi/8, \pi/4, 3\pi/8, \pi/2, 5\pi/8, 3\pi/4, \text{ and } 7\pi/8$) and 5 different scales (see Farah et al. (1995) for details). We keep the magnitudes of these filters (i.e., 40 numbers) from an 8×8 grid of 64 points, resulting in a 2560-dimensional representation of the image, which we term the perceptual level. The filter magnitudes are z-scored (shifted and scaled so they have 0 mean and unit standard deviation) on an individual basis across the data set before applying PCA. The top 40 components, again z-scored, were then used as input to a one hidden layer network. The hidden layer models the representations used for basic level categorization or fine-level discrimination, depending on the task. For basic networks, classification at the output nodes was at the basic level (i.e., four outputs, one per category) for all stimuli during pre-training and at the subordinate level (10 additional outputs) for Greebles following pre-training. For expert networks, one category (cars, cups, books, faces) was learned at the subordinate level and all other at the basic level during pre-training. Following pre-training, Greebles were learned at the subordinate level.



Fig. 2 – Example network stimuli⁹ 64×64, 8-bit, grayscale photos of books, cans, cups, faces, and Greebles used in the network simulations. Greebles are a created class of objects often used in studying expertise due to their novelty to subjects⁵. Three different images of each individual were used in training. Faces of the same person varied in expression, while images of other individual objects varied slightly in placement of the object in the image.

113 on the stimuli shown in Fig. 2. We then show that expert
 114 networks learn individuation of novel categories faster than
 115 basic networks. Thus, if cortical networks compete to solve
 116 tasks, this learning advantage suggests that the FFA, as a fine
 117 level discrimination network, would be recruited to perform
 118 novel fine-level discrimination tasks over a network that has
 119 no previous experience with such processing. An advantage of
 120 computational modeling is that the “first expertise” domain of
 121 the networks needs not be faces: our results do not depend on
 122 the order in which domains are learned, suggesting there is
 123 nothing special about faces.

124 Similar to previous work (Dailey and Cottrell, 1999; Dailey
 125 et al., 2002; Palmeri and Gauthier, 2004; Reisenhuber and
 126 Poggio, 1999), the model uses layers of processing from low
 127 level features to high level categories: (1) a Gabor filter layer
 128 models cortical responses of early visual cortex (Daugman,
 129 1985); (2) a principal components layer (learnable via Hebbian
 130 methods; Sanger, 1989) models object representations as
 131 correlations between Gabor filter responses; (3) a hidden
 132 layer models a task-specific feature representation (represent-
 133 ing subordinate or basic level processing, depending on the
 134 task), trained by back-propagation (Rumelhart et al., 1986); and
 135 (4) a categorization layer that controls the level of discrimina-
 136 tion between the stimuli, either subordinate or basic level.
 137 Minor variations of this model have accounted for a variety of
 138 behavioral face processing data (Cottrell et al., 2002; Dailey and
 139 Cottrell, 1999; Dailey et al., 2002). By analyzing the hidden
 140 layers of the two types of networks, we found that expert
 141 networks spread out the representations of similar objects in
 142 order to distinguish them. Conversely, basic networks repre-
 143 sent invariances among category members, and hence com-
 144 press them into a small region of representational space. The
 145 transformation performed by expert networks (i.e., magnifying
 146 differences) generalizes to new categories, leading to faster
 147 learning. The simulations predict that FFA neurons will have
 148 highly variable responses across members of an expert
 149 category.

2. Results and discussion 150

2.1. Network training 152

153 Training of the networks occurred in two phases. During the 153
 154 pre-training phase, two kinds of networks were trained. Basic- 154
 155 level networks were trained to differentiate a set of stimuli 155
 156 (cups, cans, books, and faces) (see Fig. 2) at the category level. 156
 157 Expert-level networks also had to perform this basic-level 157
 158 categorization, but were also required to differentiate one of 158
 159 these classes at the subordinate level. Hence there were four 159
 160 kinds of expert networks—“cup experts,” “can experts,” “book 160
 161 experts,” and “face experts.” During the second phase of 161
 162 training, a novel stimulus type, “Greebles¹,” was introduced 162
 163 and both basic and expert networks were trained to identify 163
 164 Greebles and to recognize individual Greebles. Training was 164
 165 also continued on the prior tasks. This reflects the fact that 165
 166 exposure to the new area of expertise is added to the daily 166
 167 routine of interacting with the world. This is also true in 167
 168 human experiments in creating experts in the lab, where 168
 169 training typically occurs for an hour a day over one to two 169
 170 weeks (Gauthier and Tarr, 1997). Not performing this inter- 170
 171 leaving would be equivalent to taking a human subject “out of 171
 172 the world,” and allowing them only visual exposure to the 172
 173 objects of expertise, a situation that seems unrealistic at best. If 173
 174 our model was not trained in such an interleaved fashion, face 174
 175 expertise would decay over the course of training. This may 175
 176 seem like an unrealistic prediction of the model. However, it is 176
 177 worth noting that it has recently been reported that, for one 177
 178 class of experts, this prediction would seem to hold up. Kung 178

¹ Greebles are a fictitious category of objects created by Isabel Gauthier for her Ph.D. thesis. They were constructed to have some properties similar to human categories—they have family resemblances, they have a “gender,” and are symmetric. They have gender labels, family labels, and individual names. Two examples are shown in the last column of Fig. 2.

179 et al. (2007) examined bird experts' FFA activity with respect to
 180 their degree of expertise. Expertise was measured by d' on a
 181 same/different species test with bird images. As might be
 182 expected from previous studies of visual expertise, they found
 183 that with increasing levels of bird expertise, the FFA was more
 184 activated by bird images. However, they also found that with
 185 increasing levels of bird expertise, the FFA was less activated by
 186 faces. This finding suggests that the FFA is plastic in its
 187 responsiveness depending on the kind of expertise that is
 188 most prominent in a particular subject. Our model would
 189 exhibit similar differences if it was trained more frequently on
 190 Greebles than on its original domain of expertise.

191 Basic networks learned their pre-training task the fastest
 192 and maintained the lowest error (RMSE, see Methods) until
 193 between 1280 and 5120 training epochs (one pass through the
 194 training set), when the various expert networks caught up (can,
 195 cup, book, and face experts in that order) (see Fig. 3).
 196 Conversely, the basic-level networks took by far the longest
 197 to learn the novel task (Fig. 4), obtaining no significant benefit
 198 from additional pre-training cycles. A linear trend analysis
 199 shows that all of the expert networks (but not the basic
 200 networks) learned the novel task faster if they were given more
 201 pre-training on their initial expert task, with faces benefiting
 202 the most from additional pre-training (an F -test for non-zero
 203 slope with $n=100$ for each test (10 networks at 10 time steps)
 204 yields $p=0.2962$ for basic networks and $p<0.0001$ for expert
 205 networks). Thus, for the networks learning a harder pre-
 206 training task (expert-level classification), more pre-training
 207 lead to faster learning on the secondary, expert-level task. In
 208 this study, we alternately used faces, cups, cans, and books as
 209 the primary expertise task, and Greebles as the novel
 210 (secondary) expertise task. However, we have replicated
 211 these results consistently with a variety of primary and
 212 secondary expertise tasks. For example, a network with prior

213 expertise with books learns expertise with faces faster than a
 214 network with only basic level experience with books.

215 The networks learn both the primary and secondary tasks,
 216 but are they experts? We model human subjects' reaction time
 217 as the uncertainty of the maximally activated output (see
 218 Methods). Fig. 5b shows the entry-level shift for Greebles in a
 219 network that was trained to be a face expert during pre-
 220 training (note that subordinate face model reaction times are
 221 already as low as basic level face reaction time). This curve is
 222 quite similar to the entry-level shift shown by a human subject
 223 trained in our lab to individuate Greebles (Fig. 5a). Therefore,
 224 according to the criterion used for human subjects, the
 225 networks have attained expert status.

2.2. Internal representations

226 We hypothesized that the learning advantage for expert
 227 networks was due to the larger amount of information that
 228 must be carried by the internal representations formed during
 229 training. We can visualize the representations by performing a
 230 Principal Components Analysis (PCA) of the hidden unit
 231 activations over the data and then project the data onto a
 232 two-dimensional subspace. We perform this over the training
 233 time of the network in order to see how the representations
 234 develop. This is shown in Fig. 6, in which the second and third
 235 principal components of the hidden unit activation to each
 236 input pattern are plotted against one another (the first PC just
 237 captures the magnitudes of the weights growing over time).
 238 Note the larger separation for the expert network on both
 239 subordinate and basic level categories as pre-training progresses.
 240 On the other hand, while the basic network separates the
 241 classes, it also compresses each class into a small blob in
 242 the space. Furthermore, we can project the (so far untrained)
 243 Greeble patterns into the same space, and the plot shows that
 244

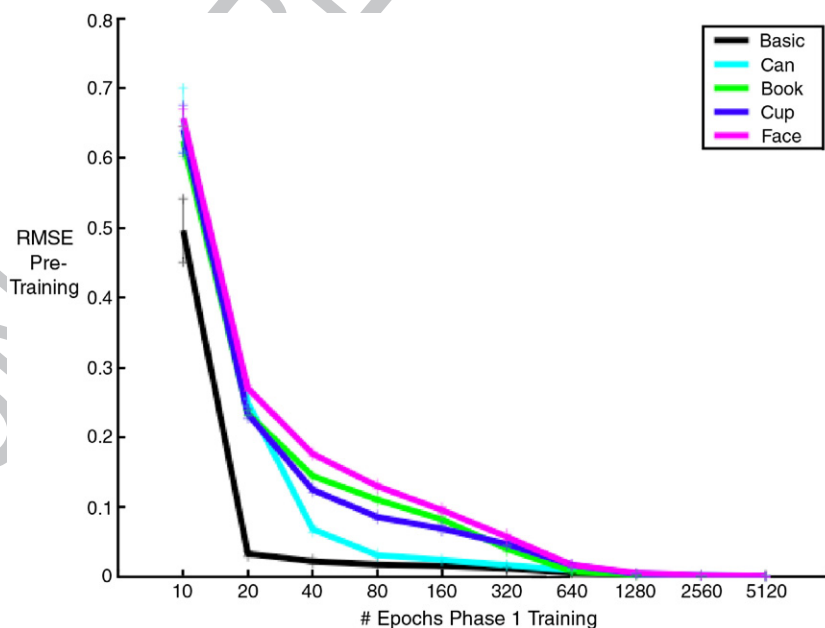


Fig. 3 – Root Mean Squared Error (RMSE) on the training set over training time for the primary task. The basic level categorization task is the easiest.

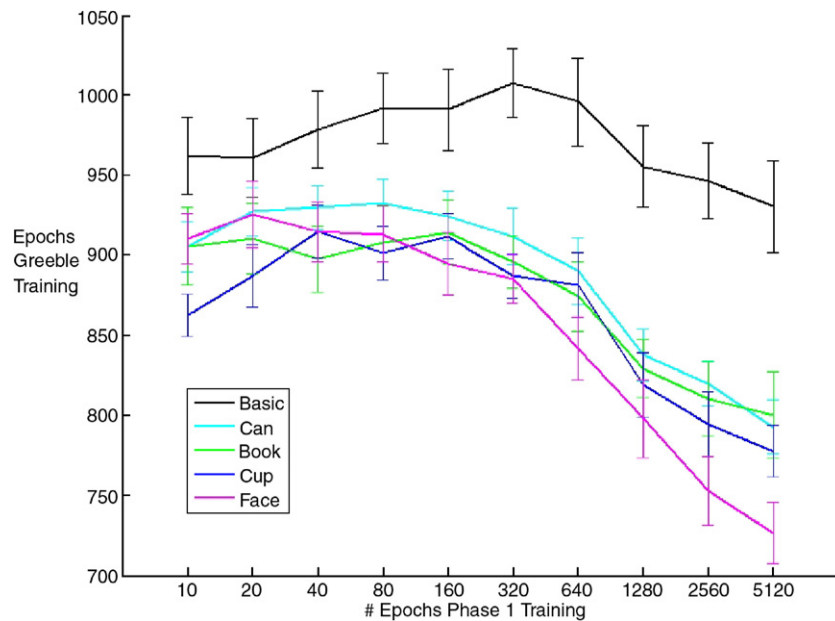


Fig. 4 – Amount of time to learn Greebles as a function of number of epochs of pre-training on the first task. Training concluded when the RMSE of the Greebles fell below 0.05. Networks at the basic level always took longer to learn Greebles than all other networks and did not benefit significantly from increased experience with the basic level task. All expert level networks benefited from more pre-training, especially faces. Error bars denote ± 1 standard error.

245 these are also more separated by the expert network—the
 246 spread of representations of homogeneous classes generalizes
 247 to a novel category. This is the fundamental reason for speeded
 248 learning of Greebles: it is easier to “pick off” each Greeble if they
 249 are in different locations in feature space to begin with.

250 A neurophysiological correlate to the above results is that
 251 the spread of representations will correspond to increased
 252 variability of single-unit responses across a homogeneous
 253 category in an expert network, and hence, in the FFA. Referring
 254 to the PCA visualization in Fig. 6, the two dimensions in that
 255 graph correspond to two “virtual unit” responses to the sti-
 256 muli. Since the points are more spread out in expert networks,
 257 this means that these units have higher variability of response
 258 across a class. We can visualize this in the single unit re-
 259 cordings shown in Fig. 7, which shows the actual activation
 260 levels of several hidden units in basic and expert networks to
 261 individual stimuli. As is clear from the figure, there is greater
 262 variability across a single class of stimuli in an expert network
 263 versus a basic network, and the greatest variability is for the
 264 class being discriminated. An analysis of variance with 5 levels
 265 of category (Expert networks shown stimuli from their domain
 266 of expertise (called Expert), Expert networks shown stimuli
 267 outside their domain of expertise (but trained at the basic-
 268 level, called Expert-basic), Expert networks shown the un-
 269 trained Greeble stimuli (Expert-Greeble), Basic networks
 270 shown stimuli from the trained basic set (Basic), and Basic
 271 networks shown the untrained Greeble stimuli (Basic-Greeble))
 272 and 11 levels of training epoch (0, 10, 20, 40, 80, 160, 320, 640,
 273 1280, 2560, 5120) was performed to determine effects. For this
 274 ANOVA, the mean variance over the relevant stimuli was the
 275 observation; thus for the expert networks 40 observations of
 276 each mean variance were available (4 types of networks, 10
 277 runs of each), while for basic networks 10 observations were

used; this yields a total of 1540 points in the ANOVA. There was
 a main effect of category [$F(4,1485)=992.91, p<0.0001$] such
 that the Expert category showed the most variance followed in
 order by the Expert-Basic, Expert-Greeble, Basic, and Basic-
 Greeble categories in order. There was also a main effect of
 epoch [$F(10,1485)=1216.73, p<0.0001$], with the least variance
 exhibited initially with variance significantly increasing across
 training epochs. There was also a significant interaction of
 category with epoch [$F(40,1485)=43.51, p<0.0001$].

To examine how this develops over time, we plot the
 average variance of response of the hidden units across a class
 over training in Fig. 8. As expected based on the PCA visuali-
 zation, the greatest variability is to the category learned at the
 subordinate level, and this variability of response extends to
 the non-expert categories as well. That is, in expert networks,
 there is more variability of response to every stimulus category
 than in networks that simply do basic-level categorization.
 Furthermore, this variability in response extended to the
 completely novel Greeble category. Note that Fig. 8 shows the
 response to untrained Greeble stimuli. When Greebles are then
 trained, the variance of response to them then increases above
 the levels shown in Fig. 8 (data not shown).

A post hoc right-tailed two-sample t-test was performed on
 the final epoch to determine the significance of the final
 ordering; all orderings were significant ($p<0.00001$, with $n=40$
 or $n=10$ measures of mean variances for expert and basic
 networks respectively) except for expert networks shown basic
 and Greeble stimuli ($p=0.8840$). All networks were initialized
 with weights drawn from the same distribution and show only
 the small differences in variance of response due to differences
 in stimuli classes, so this result is due to the effects of training
 with the pre-training stimulus prior to (and during) training
 the novel stimulus. Finally (data not shown), becoming a

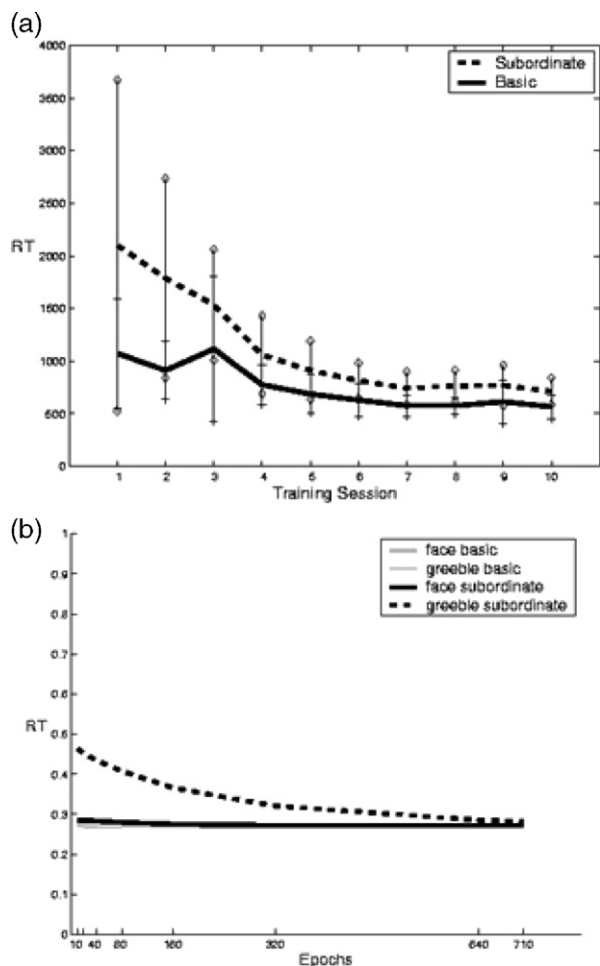


Fig. 5 – Entry level shift. (a) Typical reduction in reaction time for basic versus subordinate level judgments with training in a human subject learning to discriminate Greebles. (b) Reduction in reaction time in a neural network over training for subordinate versus basic level categorization. Reaction time is measured as the uncertainty of the maximum output ($1 - \max_{\text{output}}$).

311 Greeble expert increased variability in all networks. This
 312 caused the originally basic-level networks to resemble the
 313 other expert networks in that now their variability was higher
 314 to all categories. Based on these results, the model predicts
 315 that neurons involved in fine level discrimination, as is the
 316 hypothesis concerning the FFA, will show greater variability
 317 across stimuli that the subject possesses expertise in. This
 318 variability of response will be greater than in areas outside the
 319 FFA.

320 It is possible that these results are simply due to a scaling
 321 difference between the two types of networks, if the weights in
 322 a basic level network are simply smaller overall than in an
 323 expert network. To control for this possible artifact, we
 324 computed the variance of the object classes *relative* to the
 325 variance between classes of the internal representation. We
 326 find that the relative variance of the representation of
 327 discriminated classes in expert networks is significantly
 328 higher than in networks where these same stimuli are simply
 329 being categorized. As the PCA visualization suggests, we find

that if we average together the variability of all classes
 categorized at the basic level, and compute the ratio of this to
 the between class variance (the variance of the means), there is
 still a significant difference (using a right-tailed paired t-test
 with $n=10$, $p<0.0001$ for all pairings of expert to basic). This
 demonstrates that the expert networks are unnecessarily
 spreading out the classes they do not need to discriminate.
 Finally, we find that objects that are novel to the network
 (Greebles) also have a higher spread in expert networks (again
 using a right-tailed paired t-test with $n=10$, $p<0.0001$ for all
 pairings of expert to basic).

In the simulations discussed above, networks that learned
 a subordinate level task, and therefore exhibited a high degree
 of hidden unit variability, learned a secondary subordinate
 level task faster than basic level networks that exhibit little
 hidden unit variability. This suggests that the amount
 variance a network exhibits in response to a category prior
 to training on that category should be predictive of how fast
 that network will learn to discriminate that category.

To test this hypothesis, we performed a regression on the
 amount of variability of feature responses to Greebles prior to
 Greeble training, versus the number of epochs it takes the
 network to learn the Greeble task. There is a strong negative
 linear correlation between these two variables ($r=-0.6317$,
 $p<0.0001$), such that those networks exhibiting the lowest
 variance also take the longest to learn the Greeble task (Fig. 9).

At this point the careful reader will have noticed that the
 main effect of being an expert network is a higher variability of
 response to stimuli from the categories of expertise, and then
 wondered how this could possibly account for the increased
 BOLD signal seen in fMRI experiments in the FFA for expertise
 stimuli. One might assume we should be measuring increases
 in mean firing rates, rather than variance. However, we suggest
 that an increased variance in firing rates for neurons over a
 class of stimuli should correlate with higher mean firing rate,
 by the following argument: Biological neurons find encodings
 of the world that tend to maximize sparsity in order to
 minimize their firing rates while maintaining high levels of
 discriminability. In the interests of simplicity, our model
 contains no such bias for sparsity, and our artificial neurons
 utilize their full range of firing rates with equal probability.
 Furthermore, since both positive and negative weights are
 allowed, the actual activation of a neuron says nothing about
 its sensitivity to a particular type of stimuli; sensitivity is
 instead displayed by changes in activation, which is related to
 variance in sparse encodings. In the case of biological neurons
 with base rates near zero, an increased firing rate will result in
 a net increase in the variance. In particular, if the probability
 density function of a neuron's firing rate r follows a steep
 exponential distribution (one possible model of sparse coding),
 as in:

$$f(r, \lambda) = \lambda e^{-\lambda r} (r \geq 0)$$

then as the variance increases (given by λ^{-2}), so does the mean
 (λ^{-1}). While our model's activations do not follow this
 distribution, we argue that a more realistic model that did
 use sparse coding would also show the same increase in
 variance to stimuli of expertise. Indeed, it seems obvious now

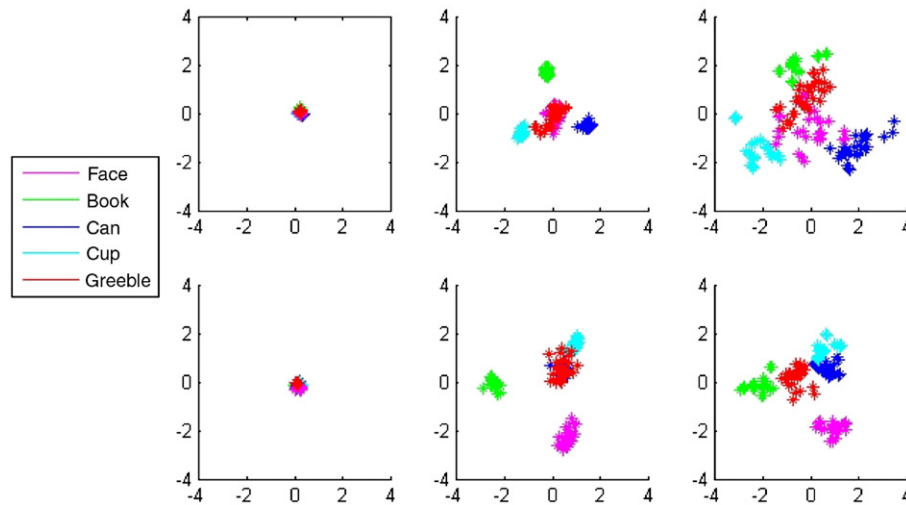


Fig. 6 – Visualization of the hidden unit representation. The figure shows the second and third principal components (the first PC simply describes a growth in activation magnitude) of the hidden unit activation to images from the training set of two types of networks, a face expert (top row) and a basic-level network (bottom row) over training time. Samples are taken at 0 epochs (column 1), 80 epochs (column 2), and 5120 epochs (column 3) of training on the first task. Colors correspond to different object categories. Both networks separate the categories over training, but the face expert (top) also spreads out the representations within each class, with the largest spread for the category learned at the subordinate level (faces). This difference in representation corresponds to a difference in variability of response of the hidden units between the expert networks and the basic networks: the farther apart each point is, the larger the difference in unit response. To demonstrate the spread of the unseen, novel stimuli (shown in red), Greebles were presented to the networks and their hidden unit activations were projected onto the principal components.

387 that within-class variability in such a model directly corre- 396
 388 sponds to having a different pattern of activation for different 397
 389 stimuli, an essential component of the ability to discriminate. 398
 390 Hence, while the goal of our model was to describe the 399
 391 recruitment of the FFA to other domains of expertise due to 400
 392 the FFA's relative fitness for such tasks compared with other 401
 393 areas, our model also does show the kind of sensitivity to 402
 394 domains of expertise that correlate with findings from the fMRI 403
 395 literature. A more literal correlation of mean firing rates would 404

require several additional assumptions in our model that go 396
 well beyond the scope of this paper. 397
 A second concern may arise due to the fact that, at least in 398
 the principal components plots of Fig. 6, it would appear that 399
 the same dimensions that are sensitive to faces are also 400
 sensitive to other stimuli. Is there really such an overlap in 401
 representation in the FFA? Recent work by Grill-Spector et al. 402
 (2006, 2007) suggests that there is. After localizing the FFA 403
 using standard fMRI, high-resolution fMRI was used to 404

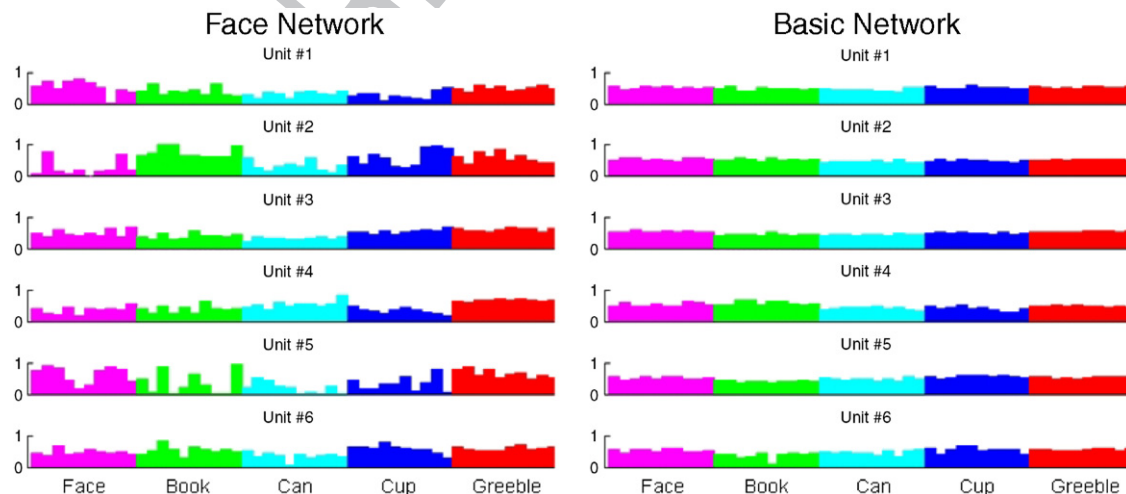


Fig. 7 – Single unit recordings of randomly chosen units from the hidden layer of an expert (face) network (left) and a basic network (right), showing the higher variability of the expert network feature responses. Each histogram shows the response of one unit to 10 stimuli from five different categories.

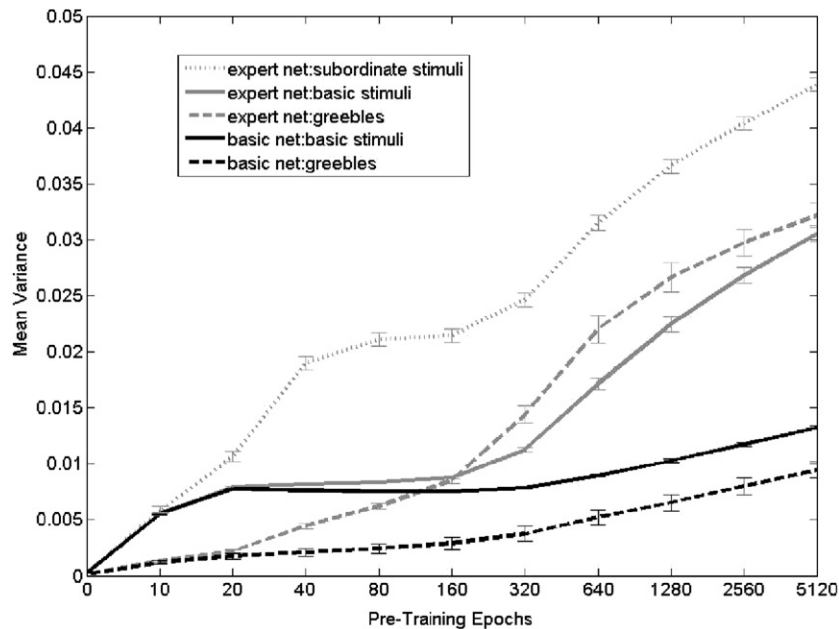


Fig. 8 – Mean variance of hidden unit activations over training. While variance in all networks increased with training, the increase was largest for expert networks, and for categories learned at the subordinate level. This variability transferred to the unlearned Greeble category. Error bars denote ± 1 standard error.

405 measure the BOLD response from 1 mm² voxels in the FFA.
 406 These voxels were assessed for their selectivity to faces, cars,
 407 animals, and abstract sculptures. In the original paper, it
 408 appeared that voxels were highly selective for each of these
 409 categories, but that face voxels were simply more numerous.
 410 However, in response to critiques of the analysis technique

(Baker et al., 2007; Simmons et al., 2007), a more accurate
 411 assessment of sensitivity was applied. The result was that, 412
 413 while most voxels were most selective for faces, they were also
 414 sensitive to other categories as well. While this does not prove
 415 anything about individual neuron tuning, it does suggest that
 416 the FFA is not just responsive to faces; it is a much more

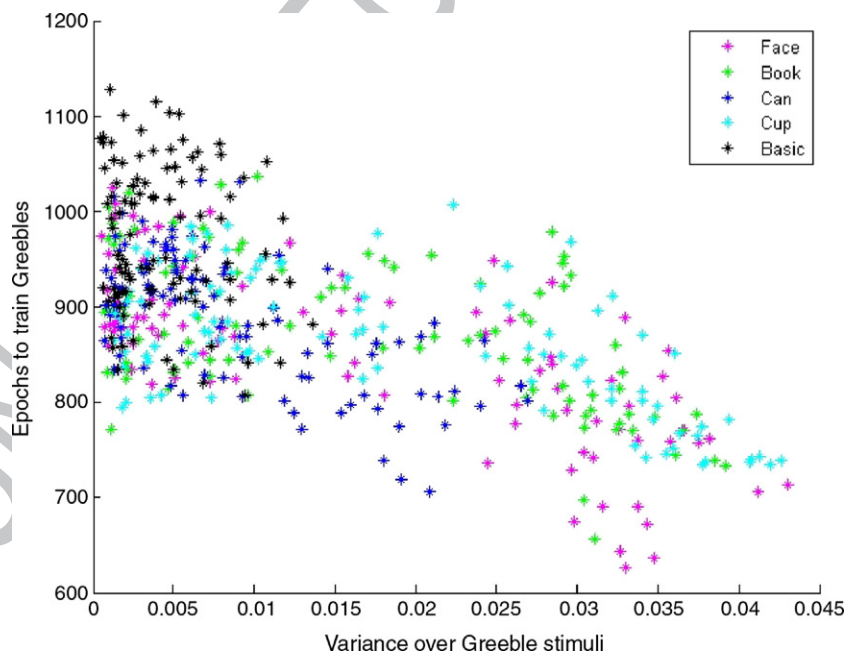


Fig. 9 – Time to learn Greebles over Greeble pre-training activation variance. As the variance of the hidden layer activations over the Greeble stimuli increases, the training required to learn Greebles decreases. This correlation is strong ($r = -0.6317$, $p < 0.0001$). This variance is taken before the networks are trained with Greebles and represents the initial spread of Greebles in representational space.

heterogeneous area than was originally thought. This analysis is also consistent with the idea that individual neurons may respond to faces and other categories of stimuli, and hence is consistent with our model's suggestion that minor re-tuning of the neural responses in the FFA is sufficient to account for the responses to new areas of expertise.

Finally, there is still a great deal of controversy whether there is a "face area" at all. Work by Haxby and colleagues (Haxby et al., 2001; Hanson et al., 2004) has shown that it is possible to accurately classify the stimulus class being observed by a subject using a standard machine learning pattern classifier applied to several different regions of cortex, that may or may not include the FFA. However, these experiments do not address the foremost role that we hypothesize for the FFA—fine level discrimination of homogeneous categories. It is not surprising that one can determine at a basic level what is being observed from multiple brain areas. Indeed, we would predict that from our model. What has not been shown that one can determine *who* is being observed from widely distributed brain activations. Thus, these data are not inconsistent with the putative role of the FFA as a fine level discrimination area.

3. Conclusions

Several effects were observed in these simulations: (1) networks can become experts, by the behavioral definition of the entry-level shift in reaction times; (2) expert networks learn the Greeble expertise task faster than basic-level categorizers; (3) this can be attributed to the spread of representations in expert networks: Greebles are more separated by these features than by the basic-network features; and (4) this feature variability to the Greeble category prior to training on it is predictive of the ease with which it will be learned. The results imply some specific hypotheses about phenomena that might be observable in human and/or primate subjects. First, though, let us be clear about what these results do *not* imply. We interpret these results to be relevant to competing cortical areas, not to different subjects learning different tasks. Thus, our results should not be interpreted to mean that subjects that have just learned a hard discrimination task should be more successful at learning a new discrimination task than subjects who have learned a simple discrimination task. Indeed, it is usually the case that it takes longer to learn novel categories of visual stimuli like these than it would if the network was starting from initial random weights. The point is rather that fine level discrimination areas are better at learning new fine level discriminations than simple object categorization areas.

What the results do suggest is that if the FFA is performing fine-level discrimination, then that task requires it to develop representations of the stimuli that separate them in representational space—the neural responses are highly differentiated. That is, similar objects have the differences between them magnified by the expert networks. On the other hand, networks that simply categorize objects map those objects into small, localized regions in representation space (this is in the space of neural firing patterns, and should not be confused with spatially localized representations). The magnifying transform of the expert networks generalizes to a novel

category, and this generalization leads to faster learning; hence, the recruitment of the FFA for Greeble expertise. We have suggested that the hidden layer of the expert networks of our model corresponds roughly to the FFA based on the equivalency of tasks and have shown that the nature of the task is sufficient to cause the recruitment of the FFA based on a shared need for fine-level discrimination; however, the actual brain is of far greater complexity than our model, and some of the changes observed in the hidden layer may turn out to be distributed among several brain areas.

An advantage of using simulations is that we were also able to show that this expertise effect is not limited to face experts. To put it in a somewhat fanciful way, the results suggest that if our parents were cans, then the Fusiform Can Area would be recruited for Greeble expertise. Furthermore, other simulations show that this learning advantage is not limited to novel Greeble expertise, nor is it dependent on the difference in the number of distinctions the two networks are making (Tong et al., 2005; Tran et al., 2004).

These simulations also make a prediction concerning the physiological responses of FFA neurons. They predict that, at the physiological level (perhaps using intracranial electrode arrays), cells in the FFA should show more variability across stimuli within a category than cells in other high-order visual object areas, and that this variability would be particularly high for categories for which the viewer possesses expertise (e.g., human and/or monkey faces). This is a falsifiable prediction of the model, and hence we look forward to our model being put to the test.

4. Experimental procedures

4.1. Training and testing

Neural networks were trained on a subordinate level classification task following various pre-training regimens. The image preprocessing steps, network configurations, and simulation procedures are described below.

The stimulus set consisted of 300 64×64 8-bit grayscale images of human faces, books, cans, cups, and Greebles (60 images per class, 5 images of 12 individuals, see Fig. 2). The five images of each Greeble were created by randomly moving the Greeble 1 pixel in the vertical/horizontal plane, and rotating up to ±3° in the image plane. Pictures of objects were taken under constant lighting and camera position, varying object position slightly over different images. Pictures of faces were frontal images of people making different facial expressions while camera angle and lighting remained constant (Cottrell and Metcalfe, 1991).

The images were preprocessed by applying Gabor wavelet filters of five scales and eight orientations as a simple model of complex cell responses in visual cortex, extracting the magnitudes, and reducing dimensionality to 40 via principal component analysis (PCA). We have found that the particular number of principal components used does not make any significant differences in our results for ranges from 30 to 50. Greeble images were not used to generate the principal components in order to model subjects' lack of experience with this category.

532 A standard feed-forward neural network architecture (40
533 input units, 60 hidden units) was used (see Fig. 1). The hidden
534 layer units used the standard logistic sigmoid function while
535 the outputs were linear. Networks were trained using back-
536 propagation of error with a learning rate of 0.01 and a
537 momentum of 0.5.

538 During pre-training, all networks were trained to perform
539 basic level categorization on all 4 non-Greeble categories. The
540 expert networks were additionally taught to perform subordi-
541 nate level categorization of one of the four categories. Non-
542 expert networks (basic level task only) had 4 output nodes
543 corresponding to book, can, cup, and face. Expert networks
544 (subordinate level task) had 14 outputs: 4 for the basic
545 categories and 1 for each of the 10 individuals (e.g., can1,
546 can2, ... can10, for a can expert). In phase 2, the pre-trained
547 networks learned subordinate level Greeble categorization
548 along with their original task. Eleven output nodes were added:
549 1 for the basic level Greeble categorization and 1 for each
550 Greeble individual. The network then learned a 15-way (basic
551 network) or 25-way (expert network) classification task. All
552 networks were trained on the same 30 images (3 images of 10
553 individuals) per class during pre-training. Thus, any differ-
554 ences in representation are due to the task, not experience
555 with exemplars. To test for generalization, 29 images were
556 used (1 new image of each of the expert category individuals
557 (10+10), plus 3 images of new basic level exemplars per
558 category). All networks generalized well.

559 Ten networks, each with different random initial weights,
560 were trained on each of the 5 pre-training tasks (basic, face
561 expert, can expert, cup expert, book expert) for 5120 epochs.
562 Image sets were randomized. Intermediate weights of each
563 network were stored every $5 \cdot 2^n$ epochs, for $n=1:10$. Phase 2
564 training was performed at each of these points ("copying" the
565 network at that point) to observe the time course of expertise
566 effects. Training concluded when the RMSE of the Greebles fell
567 below 0.05. Thus, there were a total of 50 phase 2 networks on
568 which to perform the analyses.

569 4.2. Analysis

570 The linear trend analysis on the time to learn the novel Greeble
571 identification task as a function of phase one training time was
572 performed using an *F*-test on a least-squares linear regression
573 to test for non-zero slopes. For each of the five networks, there
574 were 10 points at each of the 10 sampled epochs, yielding
575 $n=100$. The time scale used was logarithmic. Although the data
576 were non-linear, this nevertheless quantified the trend of the
577 networks as they were exposed to additional training.

578 Reaction times of the networks were modeled as the
579 uncertainty of the appropriate output. That is, for the Greeble
580 basic versus Greeble subordinate comparison in Fig. 4b, we
581 used $RT = 1 - \text{activation}$, where activation refers to the Greeble
582 output unit for the basic RT, and activation refers to the output
583 corresponding to the *i*th Greeble for the subordinate RT. Both
584 of these are averaged over all 10 Greebles for one network
585 chosen at random for the graph in Fig. 2.

586 The principal components analysis of the hidden layer was
587 performed on a network by recording the hidden unit
588 activations for every training pattern at every point during
589 which weights were saved (the initialization and the 10 stages

of phase one training). The 60×60 covariance matrix of these
data was formed, and the eigenvectors computed. A randomly
chosen set of examples from each class at each time point
were then projected onto the second and third eigenvector
and plotted. A representative set of Greeble stimuli were also
presented to the network (without training them), and their
hidden unit vectors were projected into the subspace.

The variability plots were formed by computing the
variance of each of the 60 hidden unit activations over the
appropriate class of stimuli at different training epochs. Five
levels of category were of interest: Expert networks shown
stimuli from their domain of expertise, Expert networks shown
stimuli outside their domain of expertise (but trained at the
basic-level), Expert networks shown the untrained Greeble
stimuli, Basic networks shown stimuli from the trained basic
set, and Basic networks shown the untrained Greeble stimuli.
The variance of these was tracked over 11 time samples (the
variance of the randomly initialized networks and the ten
stages of training). The variance over the 60 hidden units was
then averaged for each of the 10 networks in a given category
and epoch. As there were four categories of experts, there were
40 samples for each epoch for the expert networks, while there
were only 10 for the basic networks, yielding a total of 1540
samples of average variance. To compensate for uneven cell
sizes, an ANOVA using type 3 sum of squares was performed to
measure the effects of these 5 categories and 11 epochs. We
also computed the ratio of the average variability within a class
to the variability between classes, to measure the spread of
representations in the two types of networks, performing a two
sample *t*-test on the variance ratio after phase 1 training was
complete ($n=40$ for expert networks, $n=10$ for basic).

5. Uncited references

- | | |
|-----------------------------|-----|
| Buhmann et al., 1990 | 622 |
| Joyce and Cottrell, 2004 | 624 |
| Sugimoto and Cottrell, 2001 | 625 |

Acknowledgments

We would like to thank the reviewers for their helpful
comments on an earlier version of this paper. This work was
supported by a grant from the McDonnell Foundation to the
Perceptual Expertise Network (15573-S6), NSF grant DGE-
0333451 to GWC, NSF grant SBE-0542013 to GWC, and NIMH
Grant R01 MH57075 to GWC.

REFERENCES

- | | |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----|
| Baker, C.I., Hutchinson, T.L., Kanwisher, N., 2007. Nat. Neurosci. 10, 3-4. | 635 |
| Buhmann, J., Lades, M., von der Malsburg, C., 1990. Size and distortion invariant object recognition by hierarchical graph matching. Proceedings of the International Joint Conference on Neural Networks (IJCNN), San Diego, pp. II-411-II-416. | 636 |
| Cottrell, G.W., Metcalfe, J., 1991. Empath: face, gender and emotion recognition using holons. In: Lippman, Richard P., Moody, John, | 637 |
| | 638 |
| | 639 |
| | 640 |
| | 641 |
| | 642 |
| | 643 |
| | 644 |

- 644 Touretzky, David S. (Eds.), *Advances in Neural Information*
645 *Processing Systems 3*, pp. 564–571.
- 646 Cottrell, G.W., Branson, K.M., Calder, A.J., 2002. Do expression and
647 identity need separate representations? *Proceedings of the*
648 *24th Annual Conference of the Cognitive Science Society.*
649 Lawrence Erlbaum, Mahwah.
- 650 Dailey, M.N., Cottrell, G.W., 1999. The organization of face and
651 object recognition in modular neural network models. *Neural*
652 *Netw.* 12, 1053–1074.
- 653 Dailey, M.N., Cottrell, G.W., Padgett, C., Adolphs, R., 2002. EMPATH:
654 a neural network that categorizes facial expressions. *J. Cogn.*
655 *Neurosci.* 14 (8), 1158–1173.
- 656 Daugman, J.G., 1985. Uncertainty relation for resolution in space,
657 spatial frequency, and orientation optimized by two-dimensional
658 visual cortical filters. *J. Opt. Soc. Amer. A* 2, 1160–1169.
- 659 De Renzi, E., Perani, D., Carlesimo, G., Saveri, M., Fazio, F., 1994.
660 Prosopagnosia can be associated with damage confined to the
661 right hemisphere—an MRI and PET study and a review of the
662 literature. *Psychologia* 32 (8), 893–902.
- 663 Farah, M.J., Levinson, K.L., Klein, K.L., 1995. Face perception and
664 within-category discrimination in prosopagnosia.
665 *Neuropsychologia* 33 (6), 661–674.
- 666 Gauthier, I., Tarr, M.J., 1997. Becoming a “greeble” expert: exploring
667 mechanisms for face recognition. *Vis. Res.* 37, 1673.
- 668 Gauthier, I., Anderson, A.W., Tarr, M.J., Skudlarski, P., Gore, J.C.,
669 1997. Levels of categorization in visual recognition studied
670 with functional MRI. *Curr. Biol.* 7, 645–651.
- 671 Gauthier, I., Behrmann, M., Tarr, M.J., 1999a. Can face recognition
672 really be dissociated from recognition? *J. Cogn. Neurosci.* 11,
673 349–370.
- 674 Gauthier, I., Tarr, M.J., Anderson, A.W., Skudlarski, P., Gore, J.C.,
675 1999b. Activation of the middle fusiform “face area” increases
676 with expertise in recognizing novel objects. *Nat. Neurosci.* 2,
677 568–573.
- 678 Gauthier, I., Skudlarski, P., Gore, J.C., Anderson, A.W., 2000.
679 Expertise for cars and birds recruits brain areas involved in face
680 recognition. *Nat. Neurosci.* 3 (2), 191–197.
- 681 Grill-Spector, K., Knouf, N., Kanwisher, N., 2004. The fusiform face
682 area subserves face perception, not generic within category
683 identification. *Nat. Neurosci.* 7, 555–562.
- 684 Hanson, S.J., Matsuka, T., Haxby, J.V., 2004. Combinatorial codes in
685 ventral temporal lobe for object recognition: Haxby (2001)
686 revisited: is there a face area? *Neuroimage* 23 (1), 156–166.
- 687 Haxby, James V., Haxby, M. Ida, Furey, Maura L., Alumit, Ishai,
688 Schouten, Jennifer L., Pietro, Pietrini, 2001. Distributed and
689 overlapping representations of faces and objects in ventral
690 temporal cortex. *Science* 293 (5539), 2425–2430
691 (28 September).
- 692 Joyce, C.A., Cottrell, G.W., 2004. Solving the visual expertise mystery.
693 In: Bowman, Howard, Labiouse, Christophe (Eds.), *Connectionist*
694 *Models of Cognition and Perception II: Proceedings of the Eighth*
695 *Neural Computation and Psychology Workshop.* World
696 Scientific. 697
- 697 Kanwisher, N., 2000. Domain specificity in face perception. *Nat.*
698 *Neurosci.* 3, 759–762.
- 699 Kanwisher, N., McDermott, J., Chun, M.M., 1997. The fusiform face
700 area: a module in human extrastriate cortex specialized for
701 face perception. *J. Neurosci.* 17, 4302–4311.
- 702 Kung, Chun-Chia, Ellis, Colin, Tarr, Michael J., 2007. Dynamic
703 reorganization of fusiform gyrus: long-term bird expertise
704 reduces face selectivity. Poster Presented at the 2007 Meeting
705 of Cognitive Neuroscience Society. May 2007, New York,
706 New York.
- 707 Moscovitch, M., Winocur, G., Behrmann, M., 1997. What is special
708 about face recognition? Nineteen experiments on a person
709 with visual object agnosia and dyslexia but normal face
710 recognition. *J. Cogn. Neurosci.* 9 (5), 555–604.
- 711 Palmeri, T., Gauthier, I., 2004. Visual object understanding. *Nat.*
712 *Rev., Neurosci.* 3, 291–303.
- 713 Reisenhuber, M., Poggio, T., 1999. Hierarchical models of object
714 processing in cortex. *Nat. Neurosci.* 2, 1019–1026.
- 715 Rhodes, G., Byatt, G., Michie, P.T., Puce, A., 2004. Is the fusiform
716 face area specialized for faces, individuation, or expert
717 individuation? *J. Cogn. Neurosci.* 16, 189–203.
- 718 Rumelhart, D.E., Hinton, G.E., Williams, R.J., 1986. Learning
719 representations by back-propagating errors. *Nature* 323,
720 533–536.
- 721 Sanger, T.D., 1989. Optimal unsupervised learning in a
722 single-layer linear feedforward neural network. *Neural Netw.*
723 2, 459–453.
- 724 Simmons, W.K., Bellgowan, P.S.F., Martin, A., 2007. Measuring
725 selectivity in fMRI data. *Nat. Neurosci.* 10, 4–5.
- 726 Sugimoto, M., Cottrell, G.W., 2001. Visual expertise is a general
727 skill. *Proceedings of the 23rd Annual Cognitive Science*
728 *Conference*, pp. 994–999.
- 729 Tanaka, J.W., Taylor, M., 1991. Object categories and expertise: is
730 the basic level in the eye of the beholder? *Cogn. Psychol.* 23,
731 457–482.
- 732 Tarr, M.J., Gauthier, I., 2000. FFA: a flexible fusiform area for
733 subordinate-level visual processing automatized by expertise.
734 *Nat. Neurosci.* 3, 764–769.
- 735 Tong, M.H., Joyce, C.A., Cottrell, G.W., 2005. Are Greebles special?
736 Or, why the fusiform fish area (if we had one) would be
737 recruited for sword expertise. *Proceedings of the 27th Annual*
738 *Conference of the Cognitive Science Society.* Lawrence
739 Erlbaum, Mahwah.
- 740 Tran, B., Joyce, C.A., Cottrell, G.W., 2004. Visual expertise depends
741 on how you slice the space. *Proceedings of the 26th Annual*
742 *Conference of the Cognitive Science Conference.* Lawrence
743 Erlbaum, Mahwah, NJ.