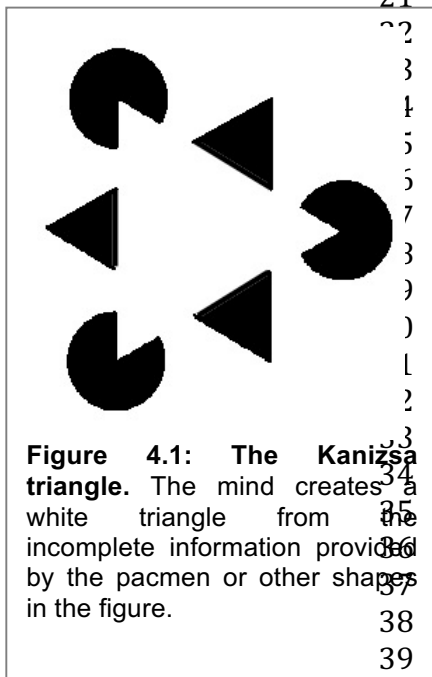


Chapter IV. Psychophysical studies of visual object recognition

We want to understand the neural mechanisms responsible for visual object recognition and we want to instantiate these mechanisms into computational algorithms that resemble and perhaps eventually surpass human performance. In order to untangle the mechanisms orchestrating visual recognition and build adequate computational models, we need to define visual recognition capabilities at the behavioral level. What shapes can humans recognize and when and how? Under what conditions do humans make mistakes? How fast can humans recognize complex objects? How much experience and what type of experience with the world is required to learn to recognize objects?

We can learn about visual object recognition by carefully quantifying human performance under a variety of well-controlled visual tasks. A discipline with the peculiar and attractive name of “Psychophysics” aims to rigorously characterize, quantify and understand behavior during cognitive tasks.

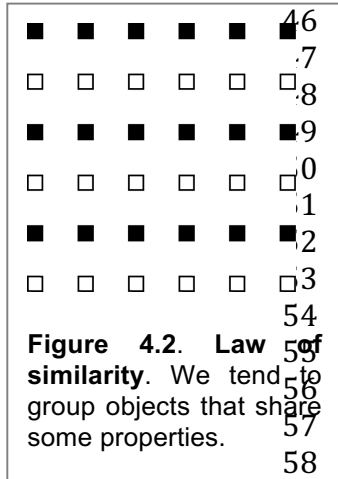
4.1. What you get ain't what you see



It is clear that what we end up perceiving is a significantly transformed version of the pattern of photons impinging on the retina. Our brains filter and process visual inputs to understand the physical world by constructing an interpretation that is consistent with our experiences. This observation may seem counterintuitive at first: our perception is a sufficiently reasonable representation of the outside world to allow us to navigate, to grasp objects, to interpret where things are going and whether a friend is happy or not. It is extremely tempting to assume that our visual system is actually capturing a perfect rendering of the outside world.

Visual illusions constitute strong examples of the dissociation between what is in the real world and what we end up perceiving. A

simple example of the dissociation between inputs and percepts is given by the blind spot. If you close one eye, there is a part of the visual field that is not mapped onto retinal ganglion cells, the spot where these cells leave the retina to form the optic nerve. It is possible to distinguish this blind spot by closing one eye, fixating on a given spot and slowly moving a finger from the center to the



periphery until part of it disappears from view (but not in its entirety which would imply that you moved your finger completely outside of your visual field).

Visual illusions are not the exception, rather they illustrate the fundamental principle that our perception is a construct, a confabulation, inspired by the visual inputs. There is a lot of information in the world that we just do not see. As a simple example, we do not perceive with our eyes information in the ultraviolet portion of the light spectrum (but other animals do). Another simple example is when we are watching a movie. A movie is nothing more than a sequence of frames, typically presented at a rate of 30

60 frames per second or more. Our brains do not perceive this rate and instead we
61 interpret objects as moving on the screen.

62
63 In addition to not being able to perceive a lot of what's happening in the
64 real world, our brains invent a lot of information that does not exist. Consider for
65 example, the Kanizsa triangle illustrated in Figure 4.1. We perceive a large white
66 triangle in the center of the image and we can trace each of the sides of said
67 triangle. Yet, those edges are composed of illusory contours: in between the
68 edge of a pacman and the adjacent small black triangle, there is no white edge.

69 4.2. Gestalt laws of grouping

70
71 One of the early and founding attempts at establishing basic principles of
72 visual perception originated from the German philosophers and experimental
73 psychologists in the late nineteenth century. The so-called *Gestalt* laws (in
74 German "gestalt" means shape) provide basic constraints about how patterns of
75 light are integrated into perceptual sensations (Reagan, 2000). These rules arose
76 from attempts to understand the basic perceptual principles that lead to
77 interpreting objects as wholes rather than the constituent isolated lines or
78 elements that give rise to them.

- 80 ■ *Law of closure.* We complete lines and extrapolate to complete known
81 patterns or regular figures. An example of this is given by the famous
82 Kanizsa triangle. Our mind creates a triangle in the middle of the image
83 from incomplete information (**Figure 4.1**).
- 84 ■ *Law of similarity.* We tend to group similar objects together. Similarity
85 could be defined by shape, color, size or brightness (**Figure 4.2**)
- 86 ■ *Law of proximity.* We tend to group objects based on their distance
87 (**Figure 4.3**).
- 88 ■ *Law of symmetry.* We tend to group symmetrical images.

- 89 ■ *Law of continuity*. We tend to continue regular patterns (**Figure 4.4**).
- 90 ■ *Law of common fate*. Elements with the same moving direction tend to be
- 91 grouped.

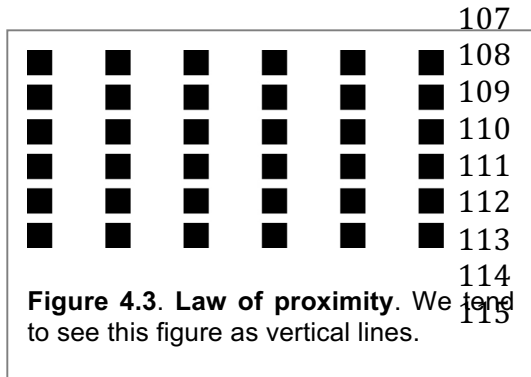
92
93
94
95

These laws are usually summarized by pointing out that the forms (Gestalten) are more than the mere sum of the component parts.

96 4.3. Holistic processing of faces

97
98

99 An interesting example of the processing and interpretation of a whole
100 image beyond what can be discerned from the individual components is the f
101 holistic processing of faces. Three main observations have been put forward to
102 document the holism of face processing. First is the inversion effect (Yin, 1969;
103 Valentine, 1988), which describes how difficult it can be to distinguish local
104 changes in a face when it is turned upside down (this is also called the “Thatcher
105 effect” alluding to the images of Britain’s prime minister originally used to
106 demonstrate the perceptual illusion). The second observation is the composite



face 2, one can create a novel face that appears to be perceptually distinct from the two original ones (Young et al., 1987). A third argument for holistic processing is the parts and wholes effect: changing a local aspect of a face distorts the overall perception of the entire face (Tanaka and Farah, 1993).

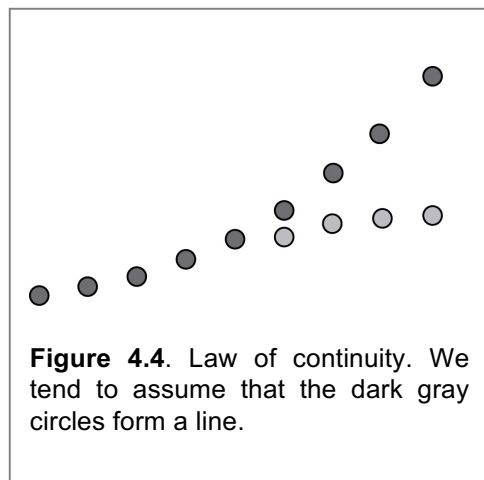
116 4.4. Tolerance to object transformations

117
118

119 A hallmark of visual recognition is our ability to identify and categorize
120 objects in spite of large transformations in the image. An object can cast an *infinite* number
121 of projections onto the retina due to changes
122 in position, scale, rotation, illumination, color,
123 etc. This invariance to image transformations
124 is critical to recognition. Our visual
125 recognition capabilities would be quite
126 useless without the ability to abstract away
127 those changes.

128
129
130
131

To further illustrate the critical role of tolerance to object transformations in visual recognition, consider a very simple algorithm



132 that we will refer to as “the rote memorization machine”. This algorithm receives
133 inputs from a digital camera and remembers every single pixel. It can remember
134 the Van Gogh sunflowers, it can remember a picture with your face taken two
135 weeks ago on Monday at 2:30pm, it can remember exactly what your car looked
136 like three years ago on a Saturday at 5:01pm. While such extraordinary memory
137 might seem quite remarkable at first, it turns out that this would constitute a
138 rather brittle approach to recognition. This algorithm would not be able to
139 recognize your car in the parking lot today, because you may see it under a
140 different illumination, a different angle, and with different amounts of dust than in
141 any of the memorized photographs. This problem is beautifully illustrated in a
142 short story by Argentinian fiction writer Jorge Luis Borges in “[Funes the](#)
143 [memorious](#)”, relating the story of a character who has infinite memory due to a
144 brain accident. Borges concludes: “To think is to forget differences, generalize,
145 make abstractions”.

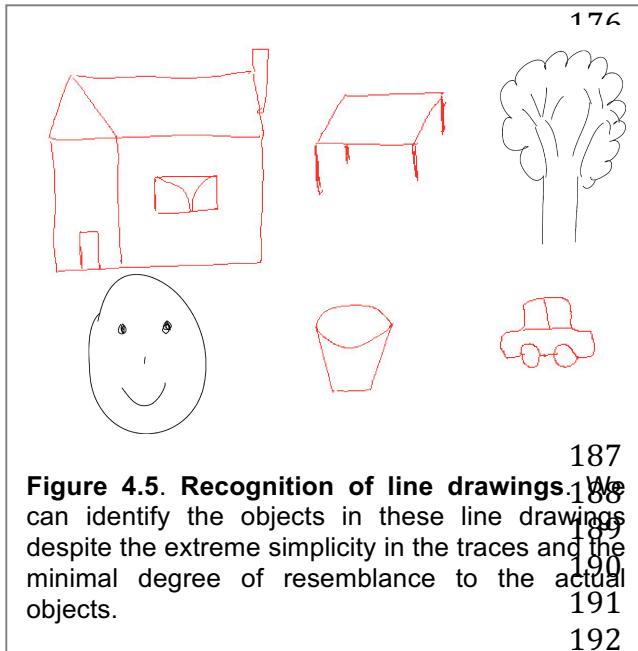
146

147 Our visual system is able to abstract away many of those image
148 transformations to recognize objects. The visual system shows a significant
149 degree of robustness to changes in many image properties, including the
150 following:

151

- 152 •Scale (e.g. you can recognize an object at different sizes). You can easily
153 demonstrate the strong degree of tolerance for object transformations.
154 For example, take a piece of text with 12pt font size, hold it at arm’s
155 length and focus on any given letter, say “A”. The A will subtend a
156 fraction of one degree of visual angle (approximately the size of your
157 thumb at arm’s length).
- 158 •Position with respect to fixation (e.g. we can recognize an object placed
159 at different distances to the fixation point)
- 160 •2D rotation (e.g. we can recognize an object after turning our head
161 sideways or rotating the object within the plane)
- 162 •3D rotation (e.g. we can recognize an object from different viewpoints)
- 163 •Color (e.g. we can recognize the objects in a photograph whether it’s in
164 color, sepia, grayscale)
- 165 •Illumination (e.g. consider illuminating an object from the left, right, top or
166 bottom)
- 167 •Cues (e.g. an object’s shape can be determined by edges, by motion
168 cues, by completion without sharp edges)
- 169 • Clutter (e.g. we can recognize objects despite the presence of other
170 objects in the image)
- 171 •Occlusion (e.g. we can recognize objects from partial information)
- 172 •Other non-rigid transformations (e.g. we can recognize faces even with
173 changes in expression, aging, even from the line drawing sketches in
174 **Figure 4.5!**)

175



A particularly intriguing example of tolerance is given by the capability to recognize caricatures and line drawings. At the pixel level, these images seem to bear little resemblance to the actual objects and yet, we can recognize them quite efficiently, sometimes even better than the real images!

4.5. Speed of visual recognition

Visual recognition *seems* almost instantaneous. Several investigators have shown that we

can recognize complex objects in a small fraction of a second.

One of the original studies by Mary Potter consisted of showing a sequence of images in a rapid sequence (RSVP, rapid serial visual presentation) and showing that subjects could detect the individual images even when presented at rates of 8 per second (Potter and Levy, 1969). Complex objects can be recognized when presented tachistoscopically for < 50 ms without a mask, even in the absence of any prior expectation or other knowledge (Vernon, 1954).

Part of the delays in reaction time measurements are associated with the behavioral response. In an attempt to constrain the amount of time required for visual recognition, Thorpe and colleagues recorded evoked response potentials from scalp electroencephalographic (EEG) signals while subjects performed a go/no-go animal categorization task (Thorpe et al., 1996). They found that frontal cortex electrodes showed a differential signal at about 150 ms; they argued that visual discrimination of animals versus non-animals in complex scenes should happen before that time. Kirchner *et al* used eye movements to elicit rapid responses and showed that subjects could make a saccade to discriminate the presence of a face or non-face stimulus in slightly more than 100 ms (Kirchner and Thorpe, 2006). These observations place a strong constraint into the mechanisms that underlie visual recognition.

Such speed in object recognition also suggests that the mechanisms that integrate information in time must occur rather rapidly. Under normal viewing conditions, all parts of an object reach the eye more or less simultaneously (in the absence of occlusions and object movement). By disrupting such synchronous access, one can probe the speed of temporal integration in vision. In a behavioral experiment to quantify the speed of integration, Jed Singer presented different parts of an object asynchronously (Figure 4.6), like breaking

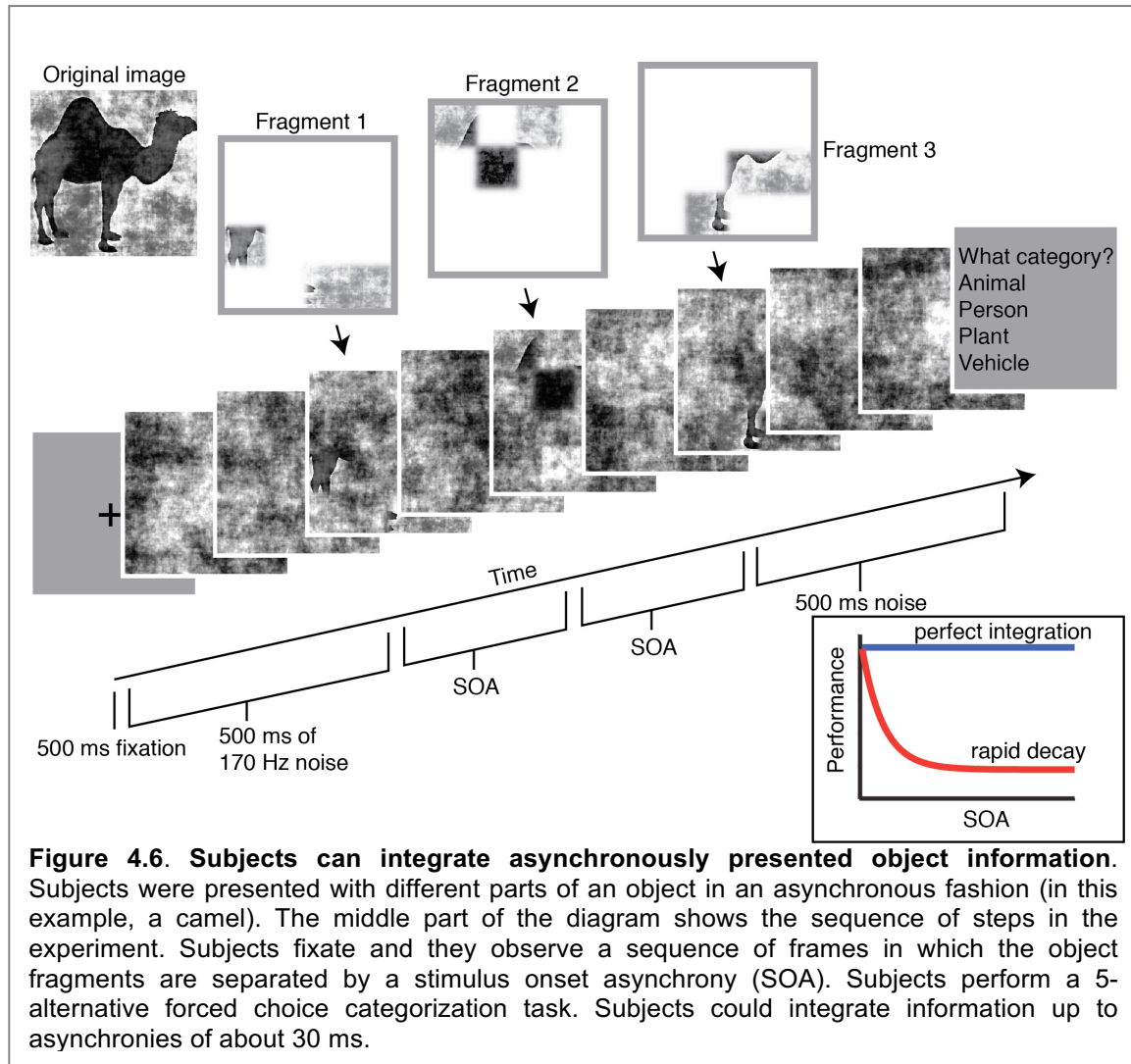
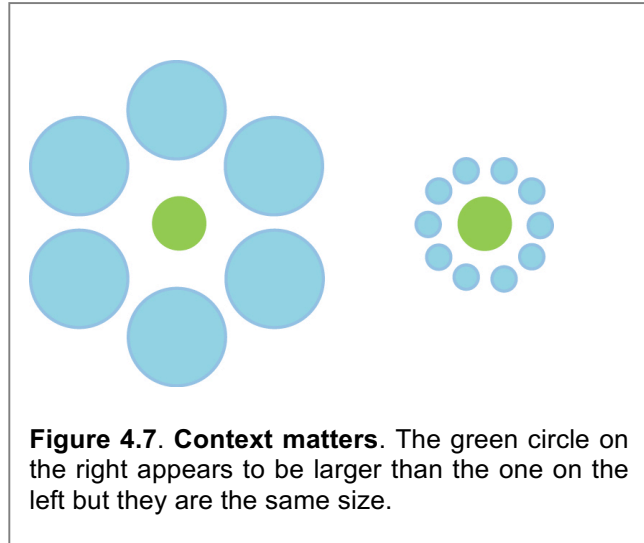


Figure 4.6. Subjects can integrate asynchronously presented object information. Subjects were presented with different parts of an object in an asynchronous fashion (in this example, a camel). The middle part of the diagram shows the sequence of steps in the experiment. Subjects fixate and they observe a sequence of frames in which the object fragments are separated by a stimulus onset asynchrony (SOA). Subjects perform a 5-alternative forced choice categorization task. Subjects could integrate information up to asynchronies of about 30 ms.

222 Humpty Dumpty and trying to put the pieces back together again. He reasoned
223 that if there was a long interval between the presentation of different parts,
224 subjects would be unable to interpret what the object was, but if the parts were
225 presented very close in time, the brain would easily be able to integrate them
226 back to a unified perception of the object.

227
228 Another striking example of temporal integration is the phenomenon
229 known as [anorthoscopic perception](#), defined as perception of a whole object in
230 cases where only a part of which is seen at a given time, perhaps one of the very
231 earliest attempts at cinema. In classical experiments, an image is seen through a
232 slit and the image moves rapidly allowing the viewer to catch only a small part of
233 the whole at any given time. The brain integrates all the snapshots and puts them
234 together to create a perception of a whole object moving. The power of temporal
235 integration is emphasized in cases where an actor is placed in a completely dark
236 room wearing black with only a few sources of information placed along his body.

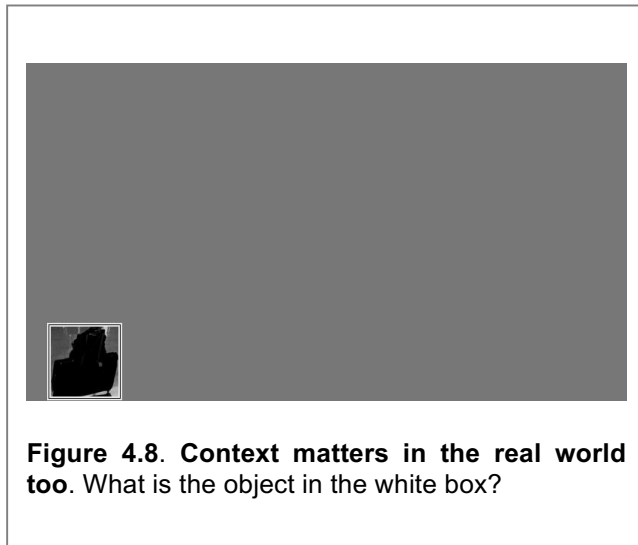
237 With just a handful of points [it is possible to infer the actor's motion patterns](#)(Johansson, 1973).
238 [it is possible to infer the actor's motion patterns](#)(Johansson, 1973).
239 [it is possible to infer the actor's motion patterns](#)(Johansson, 1973).
240 Related studies have shown that it
241 is possible to dynamically group
242 and segment information purely
243 based on temporal
244 integratoin(Anstis, 1970; Kellman
245 and Cohen, 1984).
246



247 **4.6. Beyond pixels –** 248 **contextual effects**

249 In addition to temporal
250 integration, visual recognition also
251 exploits the possibility of integrating spatial information. Several visual
252 illusions demonstrate the existence of strong contextual effects in visual object
253 recognition. For example, it is significantly more difficult to recognize faces when
254 they are upside down (see “Holistic processing” above). In a simple yet elegant
255 demonstration, the perceived size of a circle can be strongly influenced by the
256 size of its neighbors (**Figure 4.7**). Another extremely simple example is the
257 Muller-Lyer illusion: the perceived length of a line with arrows at the two ends
258 depends on the directions of the two arrows. Several entertaining examples of
259 contextual effects have been
260 reported (e.g. (Sinha and Poggio,
261 1996; Eagleman, 2001)). These
262 strong contextual dependences
263 illustrate that the visual system
264 spatially integrates information
265 and the perception of local
266 features may depend on the
267 global surrounding properties.
268

269 Such contextual effects
270 are not restricted to visual
271 illusions and psychophysics
272 demos like the one in Figure 4.7.
273 Consider **Figure 4.8**. What is the
274 object in the white box? It is
275 typically very hard to answer this question with any degree of certainty. Now, turn
276 your attention to **Figure 4.9**. What is the object in the white box? This is a much
277 easier question! Even though the pixels inside the white box are identical in both
278 figures, the surrounding contextual information dramatically changes the
279 probability of correctly detecting the object. These contextual effects are very fast
280 and can be triggered by presenting even simpler and blurred version of the
281 background information(Wu et al., Submitted). These contextual effects also
282



283 emphasize that perception constitutes an interpretation of the input in the light of
284 context and experience.

285

286 **4.5 The value of experience**

287

288 Our percepts are also
289 influenced by previous visual
290 experience. This observation holds
291 at multiple different temporal
292 scales. At short time scales,
293 several visual illusions show the
294 powerful effects of visual
295 adaptation. One such illusion is the
296 waterfall effect: after staring at a
297 waterfall for a minute or so, and
298 then shifting the gaze to other static
299 objects, those objects appear to be
300 moving upward. At longer time
301 scales, the interpretation of an
302 image could depend on whether

303 one has seen the image before. A typical example is the Dalmatian dog: for the
304 first-time observer the image consists of a smudge of black and white spots.
305 However, after recognizing the dog, people can immediately spot him the next
306 time. Other similar examples are Mooney images (Mooney, 1957).

307

308 **References**

309

- 310 Anstis SM (1970) Phi movement as a subtraction process. *Vision research*.
311 Eagleman DM (2001) Visual illusions and neurobiology. *Nat Rev Neurosci* 2:920-
312 926.
313 Johansson G (1973) Visual perception of biological motion and a model for its
314 analysis. *Perception and Psychophysics* 14:201-211.
315 Kellman PJ, Cohen MH (1984) Kinetic subjective contours. *Perception and*
316 *Psychophysics*.
317 Kirchner H, Thorpe SJ (2006) Ultra-rapid object detection with saccadic eye
318 movements: visual processing speed revisited. *Vision research* 46:1762-
319 1776.
320 Mooney CM (1957) Age in the development of closure ability in children. *Can J*
321 *Psychol* 11:219-226.
322 Potter M, Levy E (1969) Recognition memory for a rapid sequence of pictures.
323 *Journal of experimental psychology* 81:10-15.
324 Reagan D (2000) *Human perception of objects*: Sinauer.
325 Sinha P, Poggio T (1996) I think I know that face. *Nature* 384:404.
326 Tanaka JW, Farah MJ (1993) Parts and wholes in face recognition. *The*
327 *Quarterly journal of experimental psychology A, Human experimental*
328 *psychology* 46:225-245.

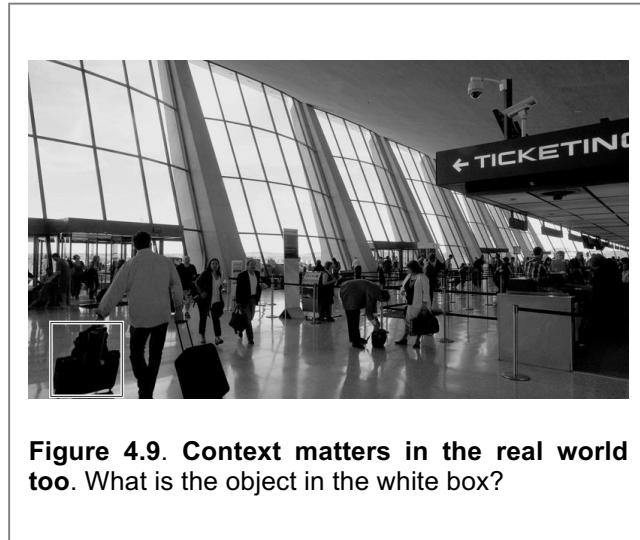


Figure 4.9. Context matters in the real world too. What is the object in the white box?

- 329 Thorpe S, Fize D, Marlot C (1996) Speed of processing in the human visual
330 system. *Nature* 381:520-522.
- 331 Valentine T (1988) Upside-down faces: a review of the effect of inversion upon
332 face recognition. *Br J Psychol* 79 (Pt 4):471-491.
- 333 Vernon M (1954) *Visual perception*. Cambridge: Harvard University Press.
- 334 Wu E, Wu K, Kreiman G (Submitted) Learning Scene Gist with Convolutional
335 Neural Networks to Improve Object Recognition.
- 336 Yin R (1969) Looking at upside-down faces. *Journal of experimental psychology*
337 81:141-145.
- 338 Young AW, Hellawell D, Hay DC (1987) Configurational information in face
339 perception. *Perception* 16:747-759.
- 340