

# Fitness Beats Truth in the Evolution of Perception

Chetan Prakash<sup>a\*</sup>, Kyle D. Stephens<sup>b</sup>, Donald D. Hoffman<sup>b</sup>,

Manish Singh<sup>c</sup>, Chris Fields<sup>d</sup>

<sup>a</sup>*Department of Mathematics, California State University, San Bernardino, CA 92407;*

<sup>b</sup>*Department of Cognitive Sciences, University of California, Irvine*

<sup>c</sup>*Department of Psychology and Center for Cognitive Science, Rutgers University, New Brunswick,*

<sup>d</sup>*243 West Spain Street, Sonoma, CA 95476*

## Abstract

Does natural selection favor veridical perceptions—those which accurately, though perhaps not exhaustively, depict objective reality? Prominent vision scientists and evolutionary theorists claim that it does. Here we formalize this claim using the tools of evolutionary game theory and Bayesian decision theory. We then present and prove a "Fitness-Beats-Truth (FBT) Theorem" which shows that the claim is false. We find that increasing the complexity of objective reality, or perceptual systems, or the temporal dynamics of fitness functions, increases the selection pressures against veridical perceptions. We illustrate the FBT Theorem with a specific example in which veridical perception minimizes expected fitness payoffs. We conclude that the FBT Theorem supports the "interface theory of perception," which proposes that our senses have evolved to hide objective reality and guide adaptive behavior. It also supports the assertion of some proponents of embodied

---

\* Corresponding author. [cprakash@csusb.edu](mailto:cprakash@csusb.edu); +1 (909) 537-5390

23 cognition that “representing the animal-independent world is not what action-oriented  
24 representations are supposed to do; they are supposed to guide action” (Chemero, 2009).

25

26 **Keywords:** Replicator Dynamics; Interface Theory of Perception; Evolutionary Game  
27 Theory; Universal Darwinism; Sensation

28

29

## 30 **1. Introduction**

31 It is standard in the perceptual and cognitive sciences to assume that more accurate  
32 perceptions are fitter perceptions and, therefore, that natural selection drives perception to  
33 be increasingly veridical, i.e. to reflect the objective world in an increasingly accurate  
34 manner. This assumption forms the justification for the prevalent view that human  
35 perception is, for the most part, veridical. For example, in his classic book *Vision*, Marr  
36 (1982) argued that:

37 “We ... very definitely do compute explicit properties of the real visible surfaces out  
38 there, and one interesting aspect of the evolution of visual systems is the gradual  
39 movement toward the difficult task of representing progressively more objective  
40 aspects of the visual world”. (p. 340)

41 Similarly, in his book *Vision Science*, Palmer (1999) states that:

42 “Evolutionarily speaking, visual perception is useful only if it is reasonably  
43 accurate ... Indeed, vision is useful precisely because it is so accurate. By and large,  
44 *what you see is what you get*. When this is true, we have what is called veridical

45 perception ... perception that is consistent with the actual state of affairs in the  
46 environment. This is almost always the case with vision.”

47 In discussing perception within an evolutionary context, Geisler and Diehl (2003) similarly  
48 assume that:

49 “In general, (perceptual) estimates that are nearer the truth have greater utility  
50 than those that are wide off the mark.”

51 In their more recent book on human and machine vision, Pizlo et al. (2014) go so far as to  
52 say that:

53 “...veridicality is an essential characteristic of perception and cognition. It is  
54 absolutely essential. *Perception and cognition without veridicality would be like*  
55 *physics without the conservation laws.*” (p. 227, emphasis theirs.)

56 If human perception is in fact veridical, it follows that the objective world shares the  
57 attributes of our perceptual experience. Our *perceived* world is three-dimensional, and is  
58 inhabited by objects of various shapes, colors, and motions. Perceptual and cognitive  
59 scientists thus typically assume that the *objective* world is so inhabited. In other words,  
60 they assume that the vocabulary of our perceptual representations is the correct vocabulary  
61 for describing the objective world and, moreover, that the specific attributes we perceive  
62 typically reflect the actual attributes of the objective world. These assumptions are  
63 embodied within the standard Bayesian framework for visual perception, which we  
64 consider in the next section.

65 Some proponents of embodied cognition reject the claim that perception is normally  
66 veridical. For instance, Chemero (2009) argues that “... perceptual systems evolved to guide  
67 behavior. Neither humans nor beetles have action-oriented representations that represent

68 the animal-independent world exactly correctly. Indeed, representing the animal-  
69 independent world is not what action-oriented representations are supposed to do; they  
70 are supposed to guide action. So the set of human affordances, that is, action-oriented  
71 representeds, is just as tightly geared to human needs and sensorimotor capacities as those  
72 of other types of animal. This leaves us with a multiplicity of conflicting sensorimotor  
73 systems, each of which is appropriate for guiding the adaptive behavior of animals whose  
74 systems they are.” The FBT Theorem, which we present below, supports Chemero’s claim. It  
75 is supported, in turn, by specific examples of non-veridical perceptions, such as those  
76 discussed by Loomis (2004) and Koenderink et. Al. (2010).

77

## 78 **2. The standard Bayesian framework for visual perception**

79 The standard approach to visual perception treats it as a problem of inverse optics: The  
80 “objective world”—taken to be 3D scenes consisting of objects, surfaces, and light sources—  
81 projects 2D images onto the retinas. Given a retinal image, the visual system’s goal is to infer  
82 the 3D scene that is most likely to have projected it (e.g. Adelson & Pentland, 1996; Feldman,  
83 2013; Knill & Richards, 1996; Mamassian, Landy, & Maloney, 2002; Shepard, 1994; Yuille &  
84 Bülthoff, 1996). Since a 2D image does not uniquely specify a 3D scene, the only way to infer  
85 a 3D scene is to bring additional assumptions or “biases” to bear on the problem—based on  
86 prior experience (whether phylogenetic or ontogenetic). For example, in inferring 3D shape  
87 from image shading, the visual system appears to make the assumption that the light source  
88 is more likely to be overhead (e.g. Kleffner & Ramachandran, 1992). Similarly, in inferring  
89 3D shape from 2D contours, it appears to use the assumption that 3D objects are maximally  
90 compact and symmetric (e.g. Li et al., 2013).

91 Formally, given an image  $x_0$ , the visual system aims to find the “best” (generally taken to  
92 mean “most probable”) scene interpretation in the world. In probabilistic terms, it must  
93 compare the posterior probability  $\mathbb{P}(w|x_0)$  of various scene interpretations  $w$ , given the  
94 image  $x_0$ . By Bayes’ Rule, the posterior probability is given by:

95 
$$\mathbb{P}(w|x_0) = \frac{\mathbb{P}(x_0|w) \cdot \mathbb{P}(w)}{\mathbb{P}(x_0)}$$

96 Since the denominator term  $\mathbb{P}(x_0)$  does not depend on  $w$ , it plays no essential role in  
97 comparing the relative posterior probabilities of different scenes interpretations  $w$ . The  
98 posterior probability is thus proportional to the product of two terms: The first is the  
99 likelihood  $\mathbb{P}(x_0|w)$  of any candidate scene interpretation  $w$ ; this is the probability that the  
100 candidate scene  $w$  could have projected (or generated) the given image  $x_0$ . Because any 2D  
101 image is typically consistent with many different 3D scenes, the likelihood will often be  
102 equally high for a number of candidate scenes. The second term is the prior probability  
103  $\mathbb{P}(w)$  of a scene interpretation; this is the probability that the system implicitly assigns to  
104 different candidate scenes, even prior to observing any image. For example, the visual  
105 system may implicitly assign higher prior probabilities to scenes where the light source is  
106 overhead, or to scenes that contain compact objects with certain symmetries. Thus, when  
107 multiple scenes have equally high likelihoods (i.e. are equally consistent with the image),  
108 the prior can serve as a disambiguating factor.

109 Application of Bayes’ Rule yields a probability distribution on the space of candidate  
110 scenes—the posterior distribution. A standard way to pick a single “best” interpretation  
111 from this distribution is to choose the world scene that has the maximal posterior  
112 probability—one that, statistically speaking, has the highest probability of being the  
113 “correct” one, given the image  $x_0$ . This is the maximum-a-posteriori or MAP estimate. More  
114 generally, the strategy one adopts for picking the “best” answer from the posterior

115 distribution depends on the choice of a loss (or gain) function, which describes the  
116 consequences of making “errors,” i.e. picking an interpretation that deviates from the “true”  
117 (but unknown) world state by varying extents. The MAP strategy follows under a Dirac-  
118 delta loss function—no loss for the “correct” answer (or “nearly correct” within some  
119 tolerance), and equal loss for everything else. Other loss functions (such as the squared-  
120 error loss) yield other choice strategies (such as the mean of the posterior distribution; see  
121 e.g. Mamassian et al., 2002). But we focus on the MAP estimate here because, in a well-  
122 defined sense, it yields the highest probability of picking the “true” scene interpretation  
123 within this framework.

124 This standard Bayesian approach embodies the “veridicality” or “truth” approach to visual  
125 perception. By this we do not mean, of course, that the Bayesian observer *always* gets the  
126 “correct” interpretation. Given the inductive nature of the problem, that would be a  
127 mathematical impossibility. It is nevertheless true that:

- 128 (i) The space of hypotheses or interpretations from which the Bayesian observer  
129 chooses is assumed to correspond to the objective world. That is, the vocabulary  
130 of perceptual experiences is assumed to be the right vocabulary for describing  
131 objective reality.
- 132 (ii) Given this setup, the MAP strategy maximizes (statistically speaking) the  
133 probability of picking the “true” world state.

134

### 135 **3. Evolution and Fitness**

136 The Bayesian framework, summarized above, focuses on estimating the world state that has  
137 the highest probability of being the “true” one, given some sensory inputs. This estimation

138 involves no notion of evolutionary fitness.<sup>2</sup> In order to bring evolution and fitness into the  
139 picture, we think of organisms as gathering fitness points as they interact with their  
140 environment. Thus each element  $w$  of the world  $W$  has associated with it a fitness value. In  
141 general, however, the fitness value depends not only on the world, but also on the organism  
142  $o$  in question (e.g., lion vs. rabbit), its state  $s$  (e.g., hungry vs. satiated), and the action class  $a$   
143 in question (e.g., feeding vs. mating). Given such a fitness landscape, natural selection favors  
144 perceptions and choices that yield more fitness points.

145 We may thus define a *global fitness function* as a (non-negative) real-valued function  $f(w, o,$   
146  $s, a)$  of these four variables. However, once we fix an organism, its state and a given action  
147 class, i.e., once we fix  $o, s$  and  $a$ , a *specific fitness function* is simply a (non-negative) real-  
148 valued function  $f: W \rightarrow [0, \infty)$  defined on the world  $W$ .

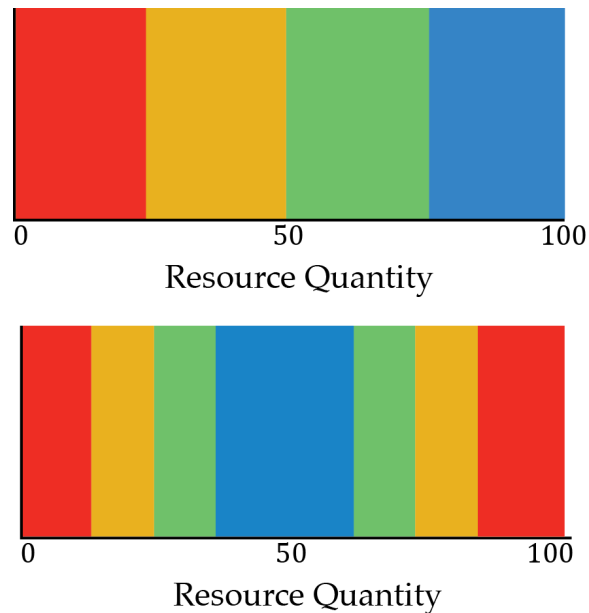
149 In order to compare the fitness of different perceptual and/or choice strategies, one pits  
150 them against one another in an evolutionary *resource game* (for simulations exemplifying  
151 the results of this paper, see, e.g., Mark, Marion, & Hoffman, 2010; Marion, 2013; and Mark,  
152 2013). In a typical game, two organisms employing different strategies compete for  
153 available territories, each with a certain number of resources. The first player observes the  
154 available territories, chooses what it estimates to be its optimal one, and receives the fitness  
155 payoff for that territory. The second player then chooses *its* optimal territory from the  
156 remaining available ones. The two organisms thus take turn in picking territories, seeking  
157 to maximize their fitness payoffs.

158 In this case, the quantity of resources in any given territory is the relevant world attribute.  
159 That is,  $W$  is here interpreted as depicting different quantities of some relevant resource.

---

<sup>2</sup> As noted above, Bayesian approaches often involve a loss (or gain) function. However, this is quite distinct from a fitness function, as defined below. Specifically, loss functions are functions of two variables  $l(x, x^*)$ , where  $x^*$  is the “true” world state, and  $x$  is a hypothetical estimate arrived at by the observer. A fitness function is, however, not a function of the observer’s *estimate*  $x$ .

160 We can then consider a perceptual map  $P : W \rightarrow X$ , where  $X$  is the set of possible sensory  
161 states, together with an ordering on it:  $P$  picks out the “best” element of  $X$  in a sense relevant  
162 to the perceptual strategy. One may, for instance, imagine a simple organism whose  
163 perceptual system has only a small number of distinct sensory states. Its perceptual map  
164 would then be some way of mapping various quantities of the resource to the small set of  
165 available sensory states. As an example, Figure 1 shows two possible perceptual mappings,  
166 i.e. two ways of mapping the quantity of resources (here, ranging from 0 through 100) to  
167 four available sensory categories (here depicted here by the four colors R, Y, G, B).



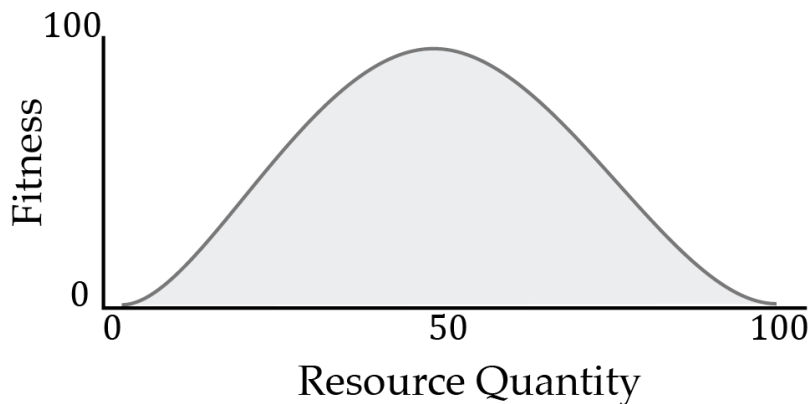
168

169 *Figure 1.* A simple example showing two different perceptual mappings  $P : W \rightarrow X$  from  
170 world states,  $W = [1, 100]$  to sensory states  $X = \{R, Y, G, B\}$ .

171 In addition, there is a fitness function on  $W$ ,  $f : W \rightarrow [0, \infty)$ , which assigns a non-negative  
172 fitness value to each resource quantity. One can imagine fitness functions that are  
173 monotonic (e.g. fitness may increase linearly or logarithmically with the number of  
174 resources), or highly non-monotonic (e.g. fitness may peak for a certain number of



175 resources, and decrease in either direction). Non-monotonic fitness functions (such as the  
176 one shown in Figure 2) are in fact quite common: too little water and one dies of thirst, too  
177 much water and one drowns. Similar arguments apply to the level of salt, or to the  
178 proportion of oxygen and indeed any number of other resources. Indeed, given the  
179 ubiquitous need for organisms to main homoeostasis, one expects non-monotonic fitness  
180 functions to be prevalent. (Moreover, from a purely mathematical point of view, the set of  
181 monotonic fitness functions is an extremely small subset of the set of all functions on a  
182 given domain. That is to say, there are “many more” non-monotonic functions than  
183 monotonic ones; hence a random sampling of fitness functions is much more likely to yield a  
184 non-monotonic one.)



185

186 *Figure 2.* An example of a non-monotonic fitness function  $f: W \rightarrow [0, \infty)$ . Fitness is maximal  
187 for an intermediate value of the resource quantity and decreases in either direction. Given  
188 the ubiquitous need for organisms to main homoeostasis, one expects that such fitness  
189 functions are quite common.

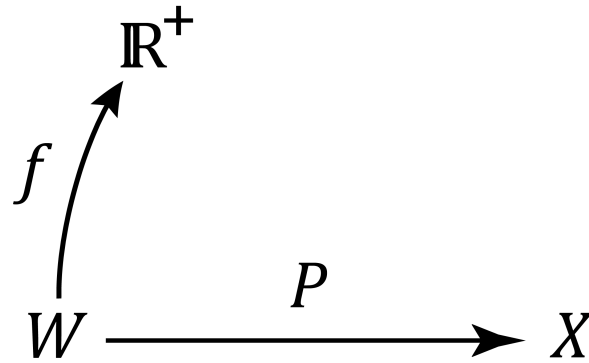
190

191 **4. Comparing perceptual strategies: “Truth” vs. “Fitness-only”**

192 In the context of these evolutionary games, in which perceptual strategies compete for  
193 resource acquisition, we take as fixed and known to the organism: the specific fitness  
194 function, its prior (in a particular state and for a particular action class) and its perceptual  
195 map (see Figure 3). On any given trial, the organism observes a number of available  
196 territories through its sensory states, say  $x_1, x_2, \dots, x_n$ . Its goal is to pick one of these  
197 territories, seeking to maximize its fitness payoff. One can now consider two possible  
198 resource strategies:

199 **The “Truth” strategy:** For each of the  $n$  sensory states, the organism estimates the world  
200 state or territory - the Bayesian MAP estimate - that has the highest probability of being the  
201 “true” one, given that sensory state. It then compares the fitness values for those estimated  
202 world states. Finally, it makes its choice of territory based on the sensory state  $x_i$  that yields  
203 the highest fitness. Its choice is thus mediated through MAP estimate of the world state.

204 **The “Fitness-only” strategy:** In this strategy, the organism makes no attempt to estimate  
205 the “true” world state corresponding to each sensory state. Rather it directly computes the  
206 expected fitness payoff that would result from each possible choice of  $x_i$ . For a given sensory  
207 state  $x_i$ , there is a posterior probability distribution (given, as with the Truth strategy, by  
208 Bayes’ formula) on the possible world states, as well as a fitness value corresponding to  
209 each world state. The organism weights these fitness values by the posterior probability  
210 distribution, in order to compute the expected fitness that would result from the choice  $x_i$ .  
211 And it picks the one with the highest expected fitness.



212

213 *Figure 3.* The framework within which we define the two resource strategies. We assume a  
 214 fixed perceptual map  $P : W \rightarrow X$  as well as a fixed fitness function  $f : W \rightarrow [0, \infty)$ . Given a  
 215 choice of available territories sensed through the sensory states, say  $x_1, x_2, \dots, x_n$ , the  
 216 organism's goal is to pick one of these, seeking to maximize its fitness payoff.

217

## 218 5. Theorems from Evolutionary Game Theory

219 In an evolutionary game between the two strategies, say  $A$  and  $B$ , the **payoff matrix** is as  
 220 follows:

	<i>against A</i>	<i>against B</i>
<i>A plays</i>	$a$	$b$
<i>B plays</i>	$c$	$d$

222 Here  $a, b, c$ , and  $d$  denote the various payoffs to the row player when playing against the  
 223 column player. E.g.,  $b$  is the payoff to  $A$  when playing  $B$ . We will refer to three main  
 224 theorems from evolutionary game theory relevant to our analysis, as follows.

225 We first consider games with infinite populations. These are investigated by means of a  
 226 deterministic differential equation, called the *replicator equation*, where time is the  
 227 independent variable and the relative population sizes  $x_A, x_B$  are the dependent variables,

228 with  $x_A + x_B = 1$  (Taylor and Jonker, 1978, Hofbauer and Sigmund, 1990, Nowak 2006). In  
229 this context, there are four generic behaviors in the long run:

230 **Theorem 1.** (Nowak 2006) *In a game with an infinite population of two types, A and B, of*  
231 *players, either*

232 (i) **A dominates B** (in the sense that a non-zero proportion of A players will eventually  
233 take over the whole population), if  $a \geq c$  and  $b \geq d$  (with at least one of the  
234 inequalities being strict);

235 (ii) **B dominates A**, if  $a \leq c$  and  $b \leq d$  (with at least one of the inequalities being strict);

236 (iii) **A and B coexist**, if  $a \leq c$  and  $b \geq d$  (with at least one of the inequalities being strict),  
237 at a stable equilibrium given by  $x_A^* = \frac{b-d}{b+c-a-d}$  (and  $x_B^* = 1 - x_A^*$ );

238 (iv) **The system is bistable**, if  $a \geq c$  and  $b \leq d$  (with at least one of the inequalities being  
239 strict) and will tend towards either all A or all B from an unstable equilibrium at the  
240 same value of  $x_A^*$  as above.

241 A fifth, non-generic possibility is that  $a = c$  and  $b = d$ , in which case we have that A and B  
242 are neutral variants of one another: any mixture of them is stable.

243 Games with a *finite* population size  $N$  can be analyzed via a *stochastic*, as against  
244 deterministic, approach. The dynamics are described by a birth-death process, called the  
245 Moran process (Moran 1958). The results are more nuanced than in the infinite population  
246 sized case: there are now *eight* possible equilibrium behaviors, and they are population  
247 dependent, not just payoff dependent.

248 Let  $\rho_{AB}$  denote the *fixation probability* of a single A individual in a population of  $N-1$  B  
249 individuals replacing (i.e., taking over completely) that population. Similarly, let  $\rho_{BA}$   
250 denote the *fixation probability* of a single B individual in a population of  $N-1$  of A individuals

251 replacing (i.e., taking over completely) that population. In the absence of any selection, we  
252 have the situation of *neutral* drift, where the probability of either of these events is just  $\frac{1}{N}$ .  
253 We say that *selection favors A replacing B* if  $\rho_{AB} > \frac{1}{N}$  and that *selection favors B replacing A* if  
254  $\rho_{BA} > \frac{1}{N}$ .

255 By analyzing the probabilities of a single individual of each type interacting with an  
256 individual of either type, or of dying off, we can use the payoff matrix above to compute the  
257 fitness  $F_i$ , when there are  $i$  entities of type  $A$ , and the fitness  $G_i$  of (the  $N-i$  individuals) of  
258 type  $B$ . If we set  $h_i = F_i - G_i$  ( $i = 1, \dots, N$ ), we can see that  $h_1 > 0$  implies that *selection*  
259 *favors A invading B*, while  $h_{N-1} > 0$  implies that *selection favors B invading A*. There are now  
260 sixteen possibilities, depending upon whether selection favors  $A$  replacing  $B$  or not;  $B$   
261 replacing  $A$  or not; whether selection favors  $A$  invading  $B$  or not; and whether selection  
262 favors  $B$  invading  $A$  or not. Of these, eight are ruled out by a theorem of Taylor, Fudenberg,  
263 Sasaki and Nowak (2004). A full description is provided in that paper, along with a number  
264 of theorems detailing the possibilities in terms of the payoff values and population size.  
265 Their Theorem 6, interpreted below as our Theorem 2, is most relevant to our analysis of  
266 evolutionary resource games: it gives conditions under which selection is *independent* of  
267 population size and is reproduced below. Interestingly, for finite populations the  
268 relationship between payoffs  $b$  and  $c$  becomes relevant:

269 **Theorem 2.** *In a game with a finite population of two types of players,  $A$  and  $B$ , if  $b > c$ ,  $a >$   
270  $c$  and  $b > d$ , we have for all  $N$ ,  $h_i > 0 \forall i$  and  $\rho_{AB} > \frac{1}{N} > \rho_{BA}$ : **selection favors A.***

271 Finally, we also consider, within large finite populations, the limit of *weak selection*. In order  
272 to model the strength of selection, a new parameter  $w$  is introduced. This parameter, lying  
273 between 0 and 1, is a measure of the strength of selection: we write the fitness of  $A$  now as

274  $f_i = 1 - w + wF_i$  and the fitness of  $B$  now as  $g_i = 1 - w + wG_i$ . When  $w = 0$ , there is no  
 275 selection: the fitnesses are equal and we have neutral drift. When  $w = 1$ , we have selection  
 276 at full strength. An analysis of the dynamics of the Moran process under weak selection (i.e.,  
 277 in the limit as  $w \rightarrow 0$ ), reveals (following Nowak 2006, equation 7.11) that:

278 **Theorem 3.** *In a game with a finite population of two types of players, A and B, and with weak*  
 279 *selection,  $(a - c) + 2(b - d) > \frac{2(a-c)-(b-d)}{N}$  implies that  $\rho_{AB} > \frac{1}{N}$ . Thus, if  $a > c$  and  $b > d$ ,*  
 280 *for large enough  $N$ , selection favors A.<sup>3</sup>*

281

## 282 6. Evolutionary Resource Games

283 For our situation of two resource strategies, we may define the payoff matrix as follows:

$a$ : to <b><i>Fitness-Only</i></b> when playing against <b><i>Fitness-Only</i></b>	$b$ : to <b><i>Fitness-Only</i></b> when playing against <b><i>Truth</i></b>
$c$ : to <b><i>Truth</i></b> when playing against <b><i>Fitness-Only</i></b>	$d$ : to <b><i>Truth</i></b> when playing against <b><i>Truth</i></b>

284

285 In a game with a very large (effectively infinite) number of players, the ***Fitness-Only***  
 286 resource strategy dominates the ***Truth*** strategy (in the sense that ***Fitness-Only*** will  
 287 eventually drive ***Truth*** to extinction) if the payoffs to ***Fitness-Only*** as first player always  
 288 exceed those of ***Truth*** as first player, regardless of who the second player is, i.e. if  $a \geq c$  and  
 289  $b \geq d$  and at least one of these is a strict inequality. If neither of these inequalities is strict,  
 290 then at the least ***Fitness-Only*** will never be dominated by ***Truth***.

---

<sup>3</sup> The value of  $N$  at which this happens depends upon the payoff matrix, but can be arbitrarily large over the set of all payoff matrices satisfying  $a > c$  and  $b > d$ .

291 Our main claim in this paper is that the *Truth* strategy—attempting to infer to the “true”  
292 state of the world that is most likely correspond to a given sensory state—confers no  
293 evolutionary advantage to an organism. In the next section, we state and prove a theorem—  
294 the “Fitness Beats Truth” theorem—which states that *Fitness-Only* will never be dominated  
295 by *Truth*. Indeed, the *Truth* strategy will generally result in a lower expected-fitness payoff  
296 than the *Fitness-Only* strategy, and is thus likely to go extinct in any evolutionary  
297 competition against the *Fitness-Only* strategy. (The statement of the FBT theorem  
298 articulates the precise way in which this is true.) We begin, first, with a numerical example  
299 that exemplifies this.

### 300 **6.1 Numerical Example of Fitness Beating Truth**

301 We give a simple example to pave the way for the ideas to follow. Suppose there are three  
302 states of the world,  $W = \{w_1, w_2, w_3\}$  and two possible sensory stimulations,  $X = \{x_1, x_2\}$ .  
303 Each world state can give rise to a sensory stimulation according to the information  
304 contained in Table 1. The first two columns give the likelihood values,  $\mathbb{P}(x|w)$ , for each  
305 sensory stimulation, given a particular world state; for instance,  $\mathbb{P}(x_1 | w_2) = 3/4$ . The third  
306 column gives the prior probabilities of the world states. The fourth column shows the  
307 fitness associated with each world state. If we think of the world states as three different  
308 kinds of food that an organism might eat, then these values correspond to the fitness benefit  
309 an organism would get by eating one of the foods. With this analogy,  $w_1$  corresponds to an  
310 extremely healthful food, while  $w_2$  and  $w_3$  correspond to moderately healthful foods, with  
311  $w_2$  being more healthful than  $w_3$  (see Table 1). This setup is the backdrop for a simple game  
312 where observers are presented with two sensory stimulations and forced to choose  
313 between them.

	Likelihood: $x_1$ given $w_j$	Likelihood: $x_2$ given $w_j$	Prior	Fitness
	$\mathbb{P}(x_1 w_j)$	$\mathbb{P}(x_2 w_j)$	$\mathbb{P}(w_j)$	$f(w_j)$
$w_1$	1/4	3/4	1/7	20
$w_2$	3/4	1/4	3/7	4
$w_3$	1/4	3/4	3/7	3

314 *Table 1: Likelihood functions, priors and fitness for our simple example where the **Truth***

315 *observer **minimizes** expected fitness, while **Fitness-only** observer **maximizes** it.*

316 Using Bayes' theorem we have calculated (see Appendix) that for  $x_1$  the **Truth** (i.e. the  
317 maximum-a-posteriori) estimate is  $w_2$ , and that for  $x_2$  this estimate is  $w_3$ . Thus, if a **Truth**  
318 observer is offered a choice between two foods to eat, one that gives it stimulation  $x_1$  and  
319 one that gives it stimulation  $x_2$ , it will perceive that it has been offered a choice between the  
320 foods  $w_2$  and  $w_3$ . Assuming that it has been shaped by natural selection to choose, when  
321 possible, the food with greater fitness, it will always prefer  $w_2$ . So, when offered a choice  
322 between  $x_1$  and  $x_2$ , the **Truth** observer will always choose  $x_1$ , with an expected utility of 5.

323 Now suppose a **Fitness-Only** observer is given the same choice. The **Fitness-Only** observer  
324 is not at all concerned with which "veridical" food these signals most likely correspond to,  
325 but has been shaped by natural selection to only care about which stimulus yields a higher  
326 expected fitness. We have calculated (see Appendix) that the expected fitness of sensory  
327 stimulation  $x_1$  is 5 and the expected utility of stimulation  $x_2$  is 6.6. Thus, when offered a  
328 choice between  $x_1$  and  $x_2$ , the **Fitness-Only** observer will always, maximizing expected  
329 fitness, choose  $x_2$ .



330 The implications of these results are clear. Consider a population of **Truth** observers  
331 competing for resources against a population of **Fitness-Only** observers, both occupying the  
332 niche described by Table 1. Since, in this case, the **Truth** observer's choice minimizes  
333 expected utility and the **Fitness-Only** observer's choice maximizes expected utility, the  
334 **Fitness-Only** population will be expected to drive the population of **Truth** observers to  
335 extinction. Seeing truth can minimize fitness; thereby leading to extinction. This conclusion  
336 is apart from considerations of the extra *energy* required to keep track of truth (see Mark,  
337 Marion and Hoffman 2010 for discussion on energy resources).

338

## 339 **7. Mathematical Background for the Main Theorem**

340 We assume that there is a fixed preliminary map,  $p$ , which associates to each world state  
341  $w \in W$  a sensory state  $x \in X$ . And we assume a fitness map on  $W$  (recall Figure 3). This  
342 places the **Truth** strategy and the **Fitness-only** strategy on a common footing where they  
343 can be set in direct competition against each other within the context of an evolutionary  
344 resource game.

345 We begin with some mathematical definitions and assumptions regarding these spaces and  
346 maps.

347 It will suffice for a basic understanding of the development in what follows, to think  
348 of  $W$  as a finite set (as in the example in 6.1).<sup>4</sup> In general, we take the **world**  $W$  to be a  
349 compact regular Borel space whose collection of measurable events is a  $\sigma$ -algebra, denoted  
350  $\mathcal{B}$ .<sup>5</sup> We assume that  $\langle W, \mathcal{B} \rangle$  comes equipped with an **a priori probability measure**  $\mu$  on

---

<sup>4</sup> in which case all the integral signs below can be replaced by summations.

<sup>5</sup> An example is a closed rectangle in some  $k$ -dimensional Euclidean space, such as the unit interval  $[0, 1]$  in one dimension, or the unit square in two.

351  $\mathcal{B}$ . We will consider only those probability measures  $\mu$  that are absolutely continuous with  
352 respect to the Borel measure on  $\mathcal{B}$ . That is, if we write  $dw$  for the uniform, or Borel,  
353 probability measure on  $W$ , then the a priori measure satisfies  $\mu(dw) = g(w) dw$ . Here  
354  $g: W \rightarrow \mathbb{R}_+$  is some non-negative measurable function, called the **density** of  $\mu$ , satisfying  
355  $\int g(w) dw = 1$ . We will take any such density to be continuous, so that it always achieves  
356 its maximum on the compact set  $W$ . This constitutes the structure of the world: a structure  
357 that applies to most biological and perceptual situations.

358 We assume that a given species interacts with its world, employing a perceptual  
359 mapping that “observes” the world via a measurable map  $p: W \rightarrow X$ . We refer to this as a  
360 *pure perceptual map* because it involves no dispersion: each world state can yield only a  
361 single sensory state  $x$ . We assume that the set of perceptual states  $X$  is a finite set, with the  
362 standard discrete  $\sigma$ -algebra  $\mathcal{X}$ , i.e., its power set (so that *all* subsets of  $X$  are measurable). In  
363 the general case, the perceptual map may have dispersion (or noise), and is mathematically  
364 expressed as a Markovian kernel  $p: W \times \mathcal{X} \rightarrow [0,1]$ . That is, for every element  $w$  in  $W$ , the  
365 kernel  $p$  assigns a probability distribution on  $X$  (hence it assigns a probability value to each  
366 measurable subset of  $X$ ). Because  $X$  is finite and all of its subsets are measurable, here the  
367 kernel may be viewed simply as assigning, for every element  $w$  in  $W$ , a probability value to  
368 each element of  $X$ .

## 369 **7.1 General Perceptual Mappings and Bayesian Inference**

370 We use the letter  $\mathbb{P}$  to indicate any relevant probability. Bayesian inference consists in a  
371 computation of the conditional probability measure  $\mathbb{P}(dw | x)$  on the world, given a  
372 particular perception  $x$  in  $X$ . The **likelihood** function is the probability  $\mathbb{P}(x | w)$  that a  
373 particular world state  $w$  could have given rise to the observed sensory state  $x$ . Then the

374 conditional probability distribution  $\mathbb{P}(dw | x)$  is the *a posteriori* probability distribution in a  
375 (partially) continuous version of Bayes formula:

$$376 \quad \mathbb{P}(dw | x) = \frac{\mathbb{P}(x | w) \mathbb{P}(dw)}{\mathbb{P}(x)}.$$

377 Since  $\mu$ , the prior on  $W$ , has a density  $g$  with respect to the Borel measure  $dw$ , we can recast  
378 this formula in terms of  $g$ : indeed,  $\mathbb{P}(dw | x)$  also has a **conditional density**,  $g(w | x)$ , with  
379 respect to the Borel measure<sup>6</sup> and we obtain

$$380 \quad g(w | x) = \frac{\mathbb{P}(x | w) g(w)}{\int \mathbb{P}(x | w') g(w')}.$$

381 We now define a **maximum a posteriori estimate** for  $x$  in  $X$  to be any  $w_x$  at which this  
382 conditional density is maximized:  $g(w_x | x) = \max\{g(w | x) | w \in W\}$ . At least one such  
383 maximum will exist, since  $g$  is bounded and piecewise continuous; however, there could be  
384 multiple such estimates for each  $x$ .

385 For a given sensory state  $x$ , the only world states that could have given rise to it lie in the  
386 **fiber over**  $x$ , i.e., the set  $p^{-1}\{x\} \subset W$ . So, for a given  $x$ , the mapping  $w \rightarrow \mathbb{P}(x | w)$  takes the  
387 value 1 on the fiber, and is zero everywhere else. This mapping may thus be viewed as the  
388 **indicator function** of this fiber. We denote this indicator function by  $1_{p^{-1}\{x\}}(w)$ .

389 For a pure mapping the **conditional density** is just

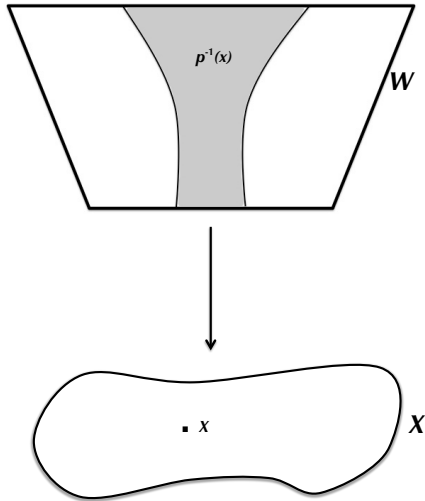
$$390 \quad g(w | x) = \frac{g(w) \cdot 1_{p^{-1}\{x\}}(w)}{\mu(p^{-1}\{x\})},$$

391 where  $\mu(p^{-1}\{x\})$  is the *a priori* measure of the fiber.

---

<sup>6</sup> That is,  $\mathbb{P}(dw | x) = g(w|x)dw$ .

392 In this special case of a pure mapping that has given rise to the perception  $x$ , we can  
 393 diagram the fiber over  $x$  on which this average fitness is computed. This is the shaded  
 394 region in figure 3 below.



395

396 *Figure 4.* The expected fitness of  $x$  is the average, using the **posterior** probability, over the  
 397 fiber  $p^{-1}(x)$ .

## 398 7.2 Expected Fitness

399 Given a **fitness function**  $f: W \rightarrow [0, \infty)$  that assigns a non-negative fitness value to each  
 400 world state, the **expected fitness** of a perception  $x$  is

$$401 \quad F(x) = \int f(w) \mathbb{P}(dw | x) = \int f(w)g(w | x) dw .$$

## 402 7.3 Two Perceptual Strategies.

403 We may build our two perceptual strategies  $P_T, P_F$ , called “*Truth*” and “*Fitness-Only*”  
 404 respectively, as compositions of a “sensory” map  $p: W \rightarrow X$  that recognizes territories and

405 “ordering” maps  $d_T, d_F: X \rightarrow X$ , where  $P_T = d_T \circ p$  and  $P_F = d_F \circ p$ . That is, the map  $d_T$  re-  
406 names the elements of  $X$  by re-ordering them, so that the best one, in terms of its *Bayesian*  
407 *MAP* estimate, is now the first,  $x_1$ , the second best is  $x_2$  etc. The map  $d_F$ , on the other hand,  
408 re-orders the elements of  $X$  so that the best one, in terms of its *expected fitness* estimate, is  
409  $x_1$ , the second best is  $x_2$  etc. The organism picks  $x_1$  if it can,  $x_2$  otherwise.

410 We can now assert our main theorem, in various contexts of evolutionary games: with  
411 infinite populations, finite populations with full selection, and sufficiently large finite  
412 populations with weak selection.

413

## 414 **8. Results**

### 415 **8.1 The “Fitness Beats Truth” Theorem**

416 The following theorem applies to infinite populations, or to large finite populations  
417 including those with weak selection:

418 **Theorem 4:** *Over all possible fitness functions and a priori measures, the probability that the*  
419 ***Fitness-only** perceptual strategy strictly dominates the **Truth** strategy is at least  $(|X| -$*   
420  *$3)/(|X| - 1)$ , where  $|X|$  is the size of the perceptual space. As this size increases, this*  
421 *probability becomes arbitrarily close to 1: in the limit, **Fitness-only** will generically strictly*  
422 *dominate **Truth**, so driving the latter to extinction.*

423 *Proof:* For any given  $x$ , the Bayesian MAP estimate is a world point  $w_x$  (it is the  $w_x$  such that  
424  $g(w_x | x) = \max\{g(w|x) \mid w \in W\}$ ). This point has fitness  $f(w_x)$ ; let  $x_M$  be that  $x$  for which  
425 the corresponding  $f(w_x)$  is maximized. Then this  $x_M$  is, if available, is chosen by **Truth** and  
426  $F(x_M)$ , its expected fitness, is the payoff to **Truth**.

427 On the other hand, the fitness payoff to the **Fitness-only** strategy is, by definition the  
428 maximum expected fitness  $F(x_I)$  over *all* fibers, so clearly,  $F(x_M) \leq F(x_I)$ .

429 As defined earlier, our evolutionary game has as payoffs,  $a$ : to **Fitness-only** when playing  
430 against **Fitness-only**;  $b$ : to **Fitness-only** when playing against **Truth**;  $c$ : to **Truth** when  
431 playing against **Fitness-only**;  $d$ : to **Truth** when playing against **Truth**.

432 We need to estimate the probability that  $a \geq c$  and  $b \geq d$ . We assume that if both strategies  
433 are the same, then each has an even chance of picking its best territory first. Thus if, in any  
434 given play of the game, two competing strategies both take a particular territory as their  
435 most favored one, then each strategy has an even chance of picking that territory and then  
436 the other strategy picks its next-best choice of territory.

437 If **Fitness-only** meets **Fitness-only**, then each has an even chance of choosing its best  
438 territory, say  $x_I$ ; the second to choose then chooses its second best territory, say  $x'_I$ . Since  
439 each player has an equal chance of being first, we have

440 
$$a = [F(x_I) + F(x'_I)]/2.$$

441 If **Truth** meets **Fitness-only**, its choice will be  $x_M$ , as long as this value differs from  $x_I$ . In this  
442 instance, we have  $a > c$ . If, however,  $x_M = x_I$ , half the time **Truth** will choose  $x_M$  and the  
443 other half  $x'_M$ , where  $x'_M$  is the second best of the optimal territories for **Truth**. Hence

444 
$$c = \begin{cases} F(x_M), & \text{if different best territories} \\ \frac{F(x_I) + F(x'_M)}{2}, & \text{if same best territories} \end{cases}$$

445 and since  $F(x_M) \leq F(x_I)$  and  $F(x'_M) \leq F(x'_I)$  we get  $a \geq c$ .

446 What happens when **Fitness-only** meets **Truth**? If **Fitness-only** goes first, the payoff will be  
447  $b = F(x_I)$ . The same is true if **Truth** goes first and the two best territories are different. If,

448 however, the two best territories are the same, then the payoff to **Fitness-only** is its second-  
449 best outcome:

$$450 \quad b = \begin{cases} F(x_I), & \text{if different best territories} \\ F(x'_I), & \text{if same best territories} \end{cases}$$

451 Finally, when **Truth** meets **Truth**, we have that

$$452 \quad d = \frac{[F(x_M) + F(x'_M)]}{2}.$$

453 So it is clear that  $b \geq d$ , as long as the two best territories are different. If they are the same,  
454 this may or may not be true: it depends on the relative size of the average  $d$  and  $F(x'_I)$   
455 (which, in this instance, also lies in between  $F(x'_M)$  and  $F(x_I) = F(x_M)$ ).

456 Now, *a priori*, there is no canonical relation between the functions  $f$  and  $g$ , both of which can  
457 be pretty much arbitrary (in fact,  $f$  need not even be continuous anywhere, and could have  
458 big jumps as well as bands of similar value separated from each other in  $W$ ). Also,  
459 generically the maximum for each strategy will be unique and also the expected fitnesses  
460 for the different territories will all be distinct.

461 Thus, generically,  $F(x_M)$  and  $F(x'_M)$  will be different from and indeed strictly less than  $F(x_I)$   
462 (and also  $F(x'_M) < F(x'_I)$ ). The only impediment to the domination of **Fitness-Only** can  
463 come from the situation where the best territories for both strategies are the same. Let  $X$   
464 have size  $|X| = n$ . There are  $n$  ways the two strategies can output the same territory, out of  
465 the  $n!/[2!(n-2)!]$  ways of pairing territories. Thus, across all possibilities for  $f$  and  $g$ , the  
466 probability that randomly chosen fitness and *a priori* measures would result in choosing the  
467 *same* territory for both strategies, i.e., that  $F(x_M) = F(x_I)$ , will happen with a probability of

$$468 \quad \frac{n}{n!/[2!(n-2)!]} = \frac{2}{n-1}$$

469 Finally, the probability of the two fibers being different is the complement:  $1 - \frac{2}{n-1} = \frac{n-3}{n-1}$ .

470 □

## 471 **8.1 Dynamic Fitness Functions**

472 A possible objection to the applicability of this theorem is that it seems to assume a *static*  
473 fitness function, whereas realistic scenarios may involve *changing*, or even rapidly changing,  
474 fitness functions. However, a close scrutiny of the proof of the theorem reveals that at any  
475 moment, the fitness function *at that time* being the same for both strategies, the relative  
476 payoffs remain in the same *generic* relation as at any other moment. Hence the theorem also  
477 applies to dynamically changing fitness functions.

478

## 479 **9. Discussion**

480 As we noted in the *Introduction*, it is standard in the literature to assume that more accurate  
481 perceptions are fitter perceptions and that, therefore, natural selection drives perception to  
482 increasing veridicality—i.e. to correspond increasingly to the “true” state of the objective  
483 world. This assumption informs the prevalent view that human perception is, for the most  
484 part, veridical.

485 Our main message in this paper has been that, contrary to this prevalent view, attempting to  
486 estimate the “true” state of the world corresponding to a given a sensory state, confers no  
487 evolutionary benefit whatsoever. Rather a strategy that simply seeks to maximize expected-  
488 fitness payoff, with no attempt to estimate the “true” world state, does consistently better  
489 (in the precise sense articulated in the statement of the “Fitness Beats Truth” Theorem).  
490 Indeed, this “Fitness-only” strategy does not estimate any single world state; it simply



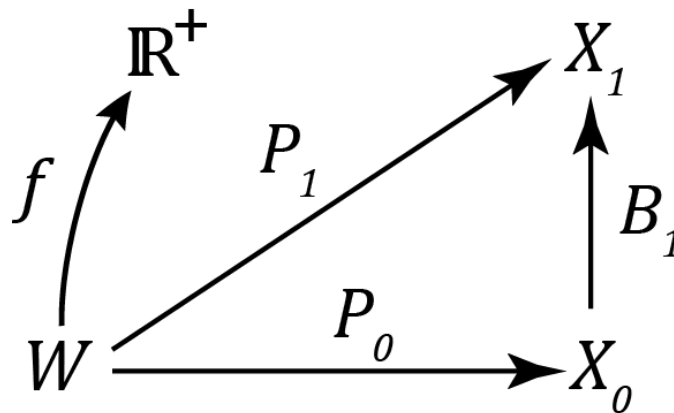
491 averages over all possible world states to compute the expected-fitness payoff  
492 corresponding to any given sensory state (this is analogous to a model-averaging strategy in  
493 model selection). And yet, as the theorem shows, in an evolutionary competition, this  
494 strategy is likely to drive the “truth” strategy to extinction.

495 At first glance, this expected-fitness strategy, based on averaging over all possible world  
496 states, may seem implausible: After all, in our own perceptual experience, we perceive  
497 things to be one particular way; we certainly don’t experience a superposition or “smear”  
498 resulting from averaging over various ways that the world *could* be. While this is  
499 undoubtedly true, one should note that this is a fact about *perceptual experience*, and  
500 provides no support whatsoever for a strategy that involves estimating the “true” state of  
501 the world. In what follows, we sketch out a more complete answer to the seeming  
502 implausibility of averaging, based on our *Interface Theory of Perception* (Hoffman, Singh, &  
503 Prakash, 2015).

504 For the purpose of the current analysis, it was essential to place the two strategies to be  
505 compared—“Truth” and “Fitness-only”—within a common framework involving Bayesian  
506 inference from the space of sensory states,  $X$ , to the world,  $W$  (recall Figure 3). This allowed  
507 us to place the two strategies on the same footing, so they could compete directly against  
508 each other. However, this result strongly supports our belief that the very idea of perception  
509 as probabilistic inference *to states of the objective world* is misguided. Perception is indeed  
510 fruitfully modeled as probabilistic inference, but the inference happens in a space of  
511 perceptual representations, and not in an objective world.

512 These ideas are part of larger theory, the *Interface Theory of Perception*, that we have  
513 described in detail elsewhere (Hoffman, 2009; Hoffman & Prakash, 2014; Hoffman & Singh,  
514 2012; Hoffman, Singh, & Prakash, 2015; see also Koenderink, 2011; 2013; 2014; von

515 Uexküll, 1934). For the purposes of the current discussion, the key point is that the standard  
 516 Bayesian framework for visual perception conflates the interpretation space (or the space  
 517 of perceptual hypotheses from which the visual system much choose) with the objective  
 518 world. This is a mistake; it is essentially the assumption that the language of our perceptual  
 519 representation is the correct language for describing objective reality—rather than simply a  
 520 species-specific interface that has been shaped by natural selection. In our ITP framework,  
 521 the probabilistic inference that results in perceptual experience takes place in a space of  
 522 perceptual representations, say,  $X_I$ , that may have no isomorphic or even homomorphic  
 523 relation whatsoever to  $W$ . The extended framework of this *Computational Evolutionary*  
 524 *Perception* is sketched in Figure 5 (see Hoffman & Singh, 2012; Hoffman, Singh, & Prakash,  
 525 2015; Singh & Hoffman, 2013).



526

527 *Figure 5.* The framework of *Computational Evolutionary Perception* in which perceptual  
 528 inferences take place in a space of representations  $X_I$  that is not isomorphic or  
 529 homomorphic to  $W$ . The more complex representational format of  $X_I$  evolves because it  
 530 permits a higher-capacity channel  $P_1 : W \rightarrow X_1$  for expected fitness, thereby allowing the  
 531 organism to choose and act more effectively in the environment (i.e. in ways that result in  
 532 higher expected-fitness payoffs).

533

534 Thus, the reason we generally perceive a single interpretation is because the probabilistic  
535 inference in the perceptual space  $X_1$  generally results in a unique interpretation. But the  
536 perceptual space  $X_1$  is not the objective world, nor is it homomorphic to it. It is simply a  
537 representational format that has been crafted by natural selection in order to support more  
538 effective interactions with the environment (in the sense of resulting in higher expected-  
539 fitness payoff). In other words, a more complex or higher-dimensional representational  
540 format (e.g. involving 3D representations in  $X_1$ , in place of 2D representations in  $X_0$ ) evolves  
541 because it permits a higher-capacity channel  $P_1 : W \rightarrow X_1$  for expected fitness (see Figure  
542 5). But this does not in any way entail that this representational format somehow more  
543 closely “resembles” the objective world. Evolution can fashion perceptual systems that are,  
544 in this sense, ignorant of the objective world because natural selection depends only on  
545 fitness and not on seeing the “truth.”

546 These considerations strongly undermine the standard assumptions that seeing more  
547 veridically enhances fitness, and that therefore one can expect that human perception is  
548 largely veridical. As human observers, we are prone to imputing structure to the objective  
549 world that is properly part of our own perceptual experience. For example, our *perceived*  
550 world is three-dimensional and populated with objects of various shapes, colors, and  
551 motions, and so we tend to conclude that the *objective world* is as well. But if, as the *Fitness-*  
552 *beats-Truth Theorem* shows, evolutionary pressures do not push perception in the direction  
553 of being increasingly reflective of objective reality, then such imputations have no logical  
554 basis whatsoever.<sup>7</sup>

## 555 **Acknowledgments**

---

<sup>7</sup> See also the *Invention of Space-Time Theorem* in Hoffman, Singh, & Prakash (2015).

556 We thank, Federico Faggin for enlightening discussions. This work has been partially  
557 funded by the Federico and Elvia Faggin Foundation.

558

- 560            1. Adelson E. H., & Pentland A. (1996). *The perception of shading and reflectance*. In: D  
561            Knill and W Richards (Eds.), *Perception as Bayesian Inference*, pp. 409-423 (New  
562            York: Cambridge University Press).
- 563            2. Chemero, A. (2009). *Radical Embodied Cognitive Science*. Cambridge, MA: MIT Press.
- 564            3. Dennett, D. C. (1995). *Darwin's Dangerous Idea*. New York: Simon & Schuster.
- 565            4. Feldman J. (2013) Tuning your priors to the world. *Topics in Cognitive Science*, 5, 13-  
566            34.
- 567            5. Geisler, W. S. and Diehl, R. L. (2003). A Bayesian approach to the evolution of  
568            perceptual and cognitive systems. *Cognitive Science* 27, 379-402.
- 569            6. Hofbauer, J., and Sigmund, K. (1998). *Evolutionary games and population dynamics*.  
570            Cambridge: Cambridge University Press.
- 571            7. Nowak, Martin A. (2006). *Evolutionary Dynamics*. Harvard University Press.
- 572            8. Hoffman, D. D. (2009). The interface theory of perception. In S. Dickinson, M. Tarr, A.  
573            Leonardis, B. Schiele, (Eds.) *Object categorization: computer and human vision*  
574            perspectives, New York: Cambridge University Press, pp. 148–165.
- 575            9. Hoffman, D. D., & Prakash, C. (2014). Objects of consciousness. *Frontiers of*  
576            *Psychology*, 5:577. DOI: 10.3389/fpsyg.2014.00577.
- 577            10. Hoffman, D. D., & Singh, M. (2012). Computational evolutionary perception.  
578            *Perception*, 41, 1073–1091.
- 579            11. Hoffman, D. D., Singh, M., & Mark, J. T. (2013). Does natural selection favor true  
580            perceptions? *Proceedings of the SPIE 8651, Human Vision and Electronic Imaging*  
581            XVIII, 865104. DOI: 10.1117/12.2011609.
- 582            12. Hoffman, D. D., Singh, M., & Prakash, C. (2015). Interface Theory of Perception.  
583            *Psychonomic Bulletin & Review*. DOI 10.3758/s13423-015-0890-8
- 584            13. Kleffner, D., & Ramachandran, V. (1992). On the perception of shape from shading.  
585            *Perception and Psychophysics*, 52, 18–36.
- 586            14. Knill, D. and Richards, W. (1996). *Perception as Bayesian inference*. Cambridge  
587            University Press, Cambridge, UK.
- 588            15. Koenderink, J. J., van Doorn, A., de Ridder, H., Oomes, S. (2010). Visual rays are  
589            parallel. *Perception*, 39, 1163-1171.
- 590            16. Koenderink, J. J. (2011). Vision as a user interface. *Human Vision and Electronic*  
591            *Imaging XVI, Interface theory of perception* 55 *SPIE Vol. 7865*. doi:  
592            10.1117/12.881671.
- 593            17. Koenderink, J. J. (2013). World, environment, umwelt, and inner-world: a biological  
594            perspective on visual awareness, *Human Vision and Electronic Imaging XVIII, SPIE*  
595            *Vol. 8651*. doi:10.1117/12.2011874.
- 596            18. Koenderink, J. J. (2014). The all seeing eye? *Perception*, 43, 1–6.
- 597            19. Li, Y., Sawada, T., Shi, Y., Steinman, R. M., & Pizlo, Z. (2013). Symmetry is the sine qua  
598            non of shape. In S. Dickinson & Z. Pizlo (Eds.), *Shape perception in human and*  
599            *computer vision* (pp. 21–40). London: Springer.
- 600            20. Loomis, J. M. (2014). Three theories for reconciling the linearity of egocentric  
601            distance perception with distortion of shape on the ground plane. *Psychology &*  
602            *Neuroscience*, 7, 3, 245-251.
- 603            21. Mamassian P, Landy M, Maloney L T, 2002 “Bayesian modeling of visual perception”,  
604            in *Probabilistic Models of the Brain: Perception and Neural Function*, Ed. R Rao, B  
605            Olshausen, M Lewicki (Cambridge, MA: MIT Press) pp 13 – 36.

606 22. Marion, B. The Impact of Utility on the Evolution of Perceptions. Dissertation,  
607 University of California, Irvine, 2013.  
608 23. Mark, J., Marion, B. and Hoffman, D. D. (2010). Natural selection and veridical  
609 perception. *Journal of theoretical Biology* 266, 504-515.  
610 24. Marr, D. (1982). *Vision*. San Francisco: Freeman.  
611 25. Moran, P. A. P. 1958. Random processes in genetics. *P. Camb. Philos. Soc.* 54: 60– 71.  
612 1962. *The statistical processes of evolutionary theory*. Oxford: Clarendon Press.  
613 26. Nowak, M. A. (2006). *Evolutionary Dynamics: Exploring the Equations of Life*. Belknap  
614 Harvard University Press, Cambridge, MA.  
615 27. Palmer, S. (1999). *Vision Science: Photons to Phenomenology*. MIT Press, Cambridge,  
616 MA.  
617 28. Pizlo, Z., Li, Y., Sawada, T., & Steinman, R.M. (2014). *Making a machine that sees like*  
618 *us*. New York: Oxford University Press.  
619 29. Shepard, R. (1994). Perceptual-cognitive universals as reflections of the world.  
620 *Psychonomic Bulletin & Review*, 1(1), 2-28.  
621 30. Singh, M., & Hoffman, D. D. (2013). Natural selection and shape perception: shape as  
622 an effective code for fitness. In S. Dickinson and Z. Pizlo (Eds.) *Shape perception in*  
623 *human and computer vision: an interdisciplinary perspective*. New York: Springer,  
624 pp. 171–185.  
625 31. Taylor, C, Fudenberg, D., Sasaki, A. & Nowak, M. Evolutionary Game Dynamics in  
626 Finite Populations. *Bulletin of Mathematical Biology* (2004) **66**, 1621–1644  
627 32. Taylor, P. D., and L. B. Jonker. 1978. “Evolutionary stable strategies and game  
628 dynamics.” *Math. Biosci.* 40: 145– 156.  
629 33. von Uexküll, J. (1934). *A stroll through the worlds of animals and men: A picture book*  
630 *of invisible worlds*. Reprinted in *Instinctive behavior: Development of a modern*  
631 *concept*, C.H. Schiller (1957). New York: Hallmark Press.  
632 34. Yuille, A., & Bülthoff, H. (1996). Bayesian decision theory and psychophysics. In:  
633 *Perception as Bayesian inference*, D. Knill and W. Richards, Eds., New York:  
634 Cambridge University Press.

635

## 636 **Appendix: Calculations for the numerical example in Table 1.**

637 In this appendix we perform the Bayesian and expected fitness calculations using the data  
638 given in Table 1.

639 To compute the **Truth** estimates, we first need the probability of each stimulation  $\mathbb{P}(x_1)$  and  
640  $\mathbb{P}(x_2)$ . These can be computed by marginalizing over the priors in the world as follows:

$$641 \quad \mathbb{P}(x_1) = p(x_1|w_1)\mu(w_1) + p(x_1|w_2)\mu(w_2) + p(x_1|w_3)\mu(w_3) = \frac{1}{4} \cdot \frac{1}{7} + \frac{3}{4} \cdot \frac{3}{7} + \frac{1}{4} \cdot \frac{3}{7} = \frac{13}{28}$$

$$642 \quad \mathbb{P}(x_2) = p(x_2|w_1)\mu(w_1) + p(x_2|w_2)\mu(w_2) + p(x_2|w_3)\mu(w_3) = \frac{3}{4} \cdot \frac{1}{7} + \frac{1}{4} \cdot \frac{3}{7} + \frac{3}{4} \cdot \frac{3}{7} = \frac{15}{28}$$

643 By Bayes' Theorem, the posterior probabilities of the world states, given  $x_1$ , are

644 
$$p(w_1|x_1) = p(x_1|w_1) \cdot \frac{\mu(w_1)}{\mathbb{P}(x_1)} = \frac{1}{4} \cdot \frac{1}{7} / \frac{13}{28} = \frac{1}{13}$$

645 
$$p(w_2|x_1) = p(x_1|w_2) \cdot \frac{\mu(w_2)}{\mathbb{P}(x_1)} = \frac{3}{4} \cdot \frac{3}{7} / \frac{13}{28} = \frac{9}{13}$$

646 
$$p(w_3|x_1) = p(x_1|w_3) \cdot \frac{\mu(w_3)}{\mathbb{P}(x_1)} = \frac{1}{4} \cdot \frac{3}{7} / \frac{13}{28} = \frac{3}{13}$$

647 Thus the maximum *a posteriori*, or **Truth** estimate for stimulus  $x_1$  is  $w_2$ .

648 Posterior probabilities of the world states, given  $s_2$ , are:

649 
$$p(w_1|x_2) = p(x_2|w_1) \cdot \frac{\mu(w_1)}{\mathbb{P}(x_2)} = \frac{3}{4} \cdot \frac{1}{7} / \frac{15}{28} = \frac{1}{5}$$

650 
$$p(w_2|x_2) = p(x_2|w_2) \cdot \frac{\mu(w_2)}{\mathbb{P}(x_2)} = \frac{1}{4} \cdot \frac{3}{7} / \frac{15}{28} = \frac{1}{5}$$

651 
$$p(w_3|x_2) = p(x_2|w_3) \cdot \frac{\mu(w_3)}{\mathbb{P}(x_2)} = \frac{3}{4} \cdot \frac{3}{7} / \frac{15}{28} = \frac{3}{5}$$

652 Thus the maximum *a posteriori*, or **Truth** estimate for stimulus  $x_2$  is  $w_3$ .

653 Finally, the expected-fitness values of the different sensory stimulations  $x_1$  and  $x_2$  are,  
654 respectively:

655 
$$F(x_1) = p(w_1|x_1)f(w_1) + p(w_2|x_1)f(w_2) + p(w_3|x_1)f(w_3) = \frac{1}{13} \cdot 20 + \frac{9}{13} \cdot 4 + \frac{3}{13} \cdot 3 = 5;$$

656 
$$F(x_2) = p(w_1|x_2)f(w_1) + p(w_2|x_2)f(w_2) + p(w_3|x_2)f(w_3) = \frac{1}{5} \cdot 20 + \frac{1}{5} \cdot 4 + \frac{3}{5} \cdot 3 = 6.6.$$

657 Thus  $x_2$  has a larger expected fitness than  $x_1$ .

658

659

660

661

662 **Highlights**

- 663 • We make rigorous mathematical definitions of two perceptual strategies employable  
664 by a given species, for a given action class and within a given environment: **Truth**,  
665 based on Bayesian estimation of assumed objective properties of the world, and  
666 **Fitness**, tuned to an arbitrary fitness function;
- 667 • Under the assumption of universal Darwinism (Dennett, 1995) we subject the two  
668 strategies to an evolutionary game analysis;
- 669 • We conclude that the **Fitness** will generally drive **Truth** to extinction, for generic  
670 fitness functions and priors;
- 671 • The likelihood of **Fitness** dominating **Truth** exceeds  $1/2$  as soon as the sensorium has  
672 more than five elements, and rises monotonically to 1 as the size of the sensorium  
673 grows towards infinity;
- 674 • This theorem holds in the presence of changing fitness functions and for large finite  
675 populations.

676