

The evolutionary puzzle of human cooperation is a red herring

2

3

Abstract

4

5 An influential line of theoretical modeling and analysis is directly motivated by
6 describing human cooperation as an “evolutionary puzzle.” This puzzle paradigm asks how
7 human cooperative behaviors can be understood as the outcome of an evolutionary process.
8 We reveal a critical flaw in the puzzle paradigm for the study of human cooperation: This
9 flaw stems primarily from the fact that puzzle models define cooperation in terms of fitness
10 costs and benefits, which starkly contrasts with the standard evolutionary approach of
11 deriving the fitness costs and benefits of a behavior from a consideration of how the behavior
12 interacts, via life-history, with evolutionary selection. As a consequence, puzzle models exist
13 in a closed theoretical loop, rendering them incapable of making meaningful predictions or
14 providing any insight into the origin or existent patterns of human cooperation. We also
15 discuss why the puzzle paradigm cannot be salvaged by conceiving of the evolutionary
16 process as cultural rather than genetic. We conclude that this approach to human cooperation
17 is a red herring but that valuable ideas for future work can be gained from questioning its
18 tacit assumptions.

18

19 Keywords. altruism, cultural evolution, evolutionary game theory, genetic evolution,
20 human cooperation, phenotypic gambit

21

22

1. Introduction

Our aim of this paper is to offer a critique of a popular research paradigm. We shall refer to it as the *puzzle paradigm* as its use tends to be motivated by an assertion that human cooperation is an evolutionary puzzle. We use four influential papers as examples of the puzzle paradigm. As described in detail in section 2, these papers are characterized by a two-step sequence of mathematical modelling: (1) Human cooperative behavior is modeled as an altruistic strategy in a giving game. The puzzle is that although humans seem to cooperate a great deal, altruistic strategies in a giving game should be unsuccessful in any evolutionary process whereby replication is based on game payoffs. (2) The giving game is then embedded in an iterated social context. This embedding is claimed to resolve the puzzle of human cooperation through a demonstration that successful strategies may involve altruistic moves.

We shall argue that the puzzle paradigm is implicitly employing a research strategy from evolutionary biology called the phenotypic gambit. In section 3 we describe the phenotypic gambit and discuss how it is not justified unless researchers take care to satisfy several criteria for its use. Two examples are given of such careful use of the phenotypic gambit in evolutionary biology. In section 4 we show how puzzle paradigm research compares unfavorably with these examples and discuss the problems that ensue. In section 5 we question the premise that human cooperation should be studied as an evolutionary puzzle. We discuss how questioning this premise opens up more promising avenues for research.

2. The puzzle paradigm

We have selected four influential papers (Boyd et al. 2003; Nowak & Sigmund 2005; Hauert et al. 2007; Boyd et al. 2010) as canonical representatives of the puzzle paradigm. These papers, which appeared in the highest ranked general science journals, study how cooperation in humans may have evolved in tandem with costly punishment mediated

1 through a variety of mechanisms. Each paper begins by asserting or implying that human
2 cooperation is an evolutionary puzzle:

3 “Unlike any other species, humans cooperate with non-kin in large groups. This behavior
4 is puzzling from an evolutionary perspective [...]” (Boyd et al. 2003).

5 “[...] human societies, are organized around altruistic, cooperative interactions. How can
6 natural selection promote unselfish behaviour?” (Nowak & Sigmund 2005).

7 “All human populations seem willing to use costly punishment to varying degrees, and
8 their willingness to punish correlates with the propensity for altruistic contributions. This
9 raises an evolutionary problem [...]” (Hauert et al. 2007).

10 “Humans are a uniquely cooperative species [...] Recent theoretical studies provide an
11 evolutionary explanation for such cooperative behavior [...] There are two important
12 problems with this explanation of human cooperation” (Boyd et al. 2010).

13 These quotes highlight the two distinct but related research questions the puzzle
14 paradigm seeks to answer: Why are humans unique among species with regard to
15 cooperation, and how can cooperative behaviors be selected for by an evolutionary process?
16 The second question is precisely the central question of the branch of theoretical biology that
17 studies the evolution of altruism (e.g., Fletcher & Doebeli 2006).

18 **2.1 The first step: Modeling cooperation as altruism in a giving game**

19 In the puzzle paradigm cooperative behavior is taken as a synonym of altruism and
20 defined as a behavior where an individual pays a personal fitness cost to bestow a fitness
21 benefit upon others. We will refer to a situation that offers the possibility for such an
22 altruistic act as a “giving game”. Thus, in the puzzle paradigm cooperative behavior is

1 modeled as an altruistic act in some version of a giving game as indicated by the following
2 quotes.

3 “There are two behavioral types, contributors and defectors. Contributors incur a cost c to
4 produce a total benefit b that is shared equally among group members. Defectors incur no
5 costs and produce no benefits.” (Boyd et al. 2003)

6 “In the simplest model, the altruistic act consists in conferring a benefit b on the recipient
7 at a cost c to the donor.” (Nowak & Sigmund 2005)

8 “Those who participate can decide whether or not to contribute an investment at a cost c to
9 themselves. All individual contributions are added up and multiplied with a factor $r > 1$.
10 This amount is then divided equally among all participants of the public goods game.”
11 (Hauert et al. 2007)

12 “Cooperation costs the cooperator c and benefits each member of the group b/n ($b > c >$
13 b/n .)” (Boyd et al. 2010)

14 From these quotes it is clear that cooperation is defined in terms of the payoffs in this
15 giving game. Moreover, these payoffs are assumed to be in the currency of fitness. In other
16 words, the strategy that incurs the largest net benefit will increase most in frequency. A
17 “cooperative” strategy of paying a cost to give a benefit to others should therefore be
18 outcompeted by a “non-cooperative” strategy. The emergence and persistence of high levels
19 of cooperation under these conditions should not be possible.

20 **2.2 The second step: Embedding the giving game to resolve the puzzle**

21 The puzzle paradigm’s research goal is to resolve the proposed puzzle by showing that
22 high levels of cooperation can emerge and persist under the additional assumption that the
23 giving game is embedded in an iterated social context, in other words a meta-game. This

1 embedding varies from paper to paper and is rather lengthy to describe. For instance, the
2 embedding of Boyd et al. (2003) amounts to the following four steps: (1) structuring the
3 population into groups, with a certain amount of migration between groups; (2) adding a third
4 strategy type that both cooperates *and* punishes non-cooperators in their group, at certain
5 costs to themselves and the punished parties; (3) letting relative replication of strategies
6 depend on the total payoffs obtained from the giving game and punishments; (4) letting
7 groups meet each other in pairwise competition, with the probability of winning decided by
8 the frequency of cooperation in each group, and with the winning group replacing the losers.
9 Given this embedding it is then demonstrated that, as meta-game strategies replicate,
10 cooperative behavior in the giving game and punishing behavior in the punishment game are
11 both sustained at high levels.

12 The other three puzzle paradigm papers in our sample similarly embed the simple
13 giving game in a complex meta-game, each with some particular features: Nowak and
14 Sigmund (2005) deal with reputation and indirect reciprocity, Hauert et al. (2007) with
15 optional participation, and Boyd et al. (2010) with coordination and economies of scale.

16 **3. The phenotypic gambit**

17 A “gambit” is a chess opening in which a piece is sacrificed in order to gain what is
18 hopefully a better position. The phenotypic gambit was perhaps first made by Fisher (1930)
19 in an attempt to understand the evolution of sex ratios. The term phenotypic gambit was
20 coined in Grafen (1984). We quote part of his definition:

21 “The phenotypic gambit is to examine the evolutionary basis of a character as if the
22 simplest possible genetic system controlled it: as if there were a haploid locus at which
23 each distinct strategy was represented by a distinct allele, as if the payoff rule gave the
24 number of offspring for each allele, and as if enough mutation occurred to allow each
25 strategy to invade.” (Grafen 1984)

1 This quote emphasizes the simplification of genetic inheritance mechanisms, while
2 perhaps deemphasizing the simplification of the complex and context dependent mapping
3 from genotype to fitness via phenotype. We encourage readers to read Grafen's introduction
4 of the term in its entirety, as well as other related accounts of when and why such a gambit is
5 justified (Smith 1978; Dawkins 1982 pp. 30-54). Here we will do our best to paraphrase these
6 accounts.

7 **3.1 The sacrifice of complexity**

8 The basic idea of the phenotypic gambit can be expressed as follows: Instead of starting
9 with a genetic replicator and asking whether its frequency will increase or decrease in a
10 population as a result of evolution, the researcher starts with a phenotype and asks the same
11 question. In order for the frequency of this phenotype to be affected by evolution at all, the
12 variation in the phenotype must be to some extent heritable (i.e., co-vary with underlying
13 alleles at least at one, but possibly many, genetic loci). The details of this genetic heritability
14 may be complicated by various factors such as pleiotropy, epistasis, and linkage
15 disequilibrium. The first thing the phenotypic gambit puts on the chopping block is this
16 genetic complexity. Replacing it is a set of haploid alleles at a single locus, which correspond
17 one-to-one with the set of possible phenotypes or strategies.

18 The question of which phenotypes are actually possible is dependent on the genetic
19 mechanisms by which new phenotypes arise. In general, this phenotypic variation will be
20 constrained by the phylogenetic history of the organism, the ontological process by which the
21 phenotype develops, and the bio-physical processes by which the phenotype manifests. The
22 phenotypic gambit also largely ignores these constraints by leaving the choice of strategies to
23 be considered at the discretion of the researcher. Choosing the set of strategies to be
24 considered is a delicate task, fraught with the potential for implicit assumptions. In addition,

1 not only the set of available strategies but also the manner in which new strategic variants
2 (mutants) derive from pre-existing types radically impacts evolutionary outcomes (e.g. van
3 Veelen et al. 2012).

4 To summarize, the phenotypic gambit sacrifices several layers of complexity at a
5 considerable risk of over-simplification. The gain comes in the form of access to powerful
6 theoretical tools, to be described next.

7 **3.2 The analysis enabled by the sacrifice of complexity**

8 The notion of evolutionarily stable strategies (ESS; Smith & Price 1973) provides a
9 clear and intuitive way of thinking about evolutionary outcomes and making evolutionary
10 predictions. Taylor and Jonker (1978) formally linked the ESS to the equilibria of genetic
11 dynamical systems of the kind described in the above quote from Grafen (1984). In other
12 words, the particular assumptions of the phenotypic gambit empower a theorist to forgo an
13 explicit investigation of genetic dynamics, which may be intractable, and focus instead on a
14 relatively simple ESS analysis. However, even after this simplification there remains the task
15 of deriving an appropriate fitness function.

16 A fitness function takes as inputs the frequencies of each strategy in the current
17 population and gives as output the relative rates of growth or decline for each strategy. The
18 derivation of the fitness function constitutes the core of a phenotypic gambit analysis. If the
19 strategies being considered have been empirically observed, and the fitness of each strategy is
20 known, this function may be fit to data. Typically though, the phenotypic gambit is used to
21 provide an evolutionary explanation for the existent variation, both inter and intra species, in
22 light of the strategic situation created by the interaction of various ecological factors. In such
23 cases empirical fitting of the fitness function may be infeasible. In some cases the behaviors
24 considered are intrinsically linked to individual fitness. The derivation of a fitness function

1 from the given strategy set and life history assumptions is then unambiguous, though possibly
2 mathematically complex. In more challenging cases, strategies have ambiguous impacts on
3 fitness and theoreticians are left with the hard choice of either making further simplifying
4 assumptions, until the derivation of fitness becomes unambiguous (e.g., Higginson et al.
5 2012), or attempting to work with a general broad class of possible fitness functions (e.g.,
6 Pen & Taylor 2005).

7 **3.3 The caveats of the phenotypic gambit**

8 Once derived, the fitness function manifests the essence of an interesting trade-off or
9 strategic social interaction. However, even if the fitness function is ecologically valid there is
10 no guarantee that the implications of this fitness function under the assumptions of the
11 phenotypic gambit will still hold when these assumptions are relaxed. A researcher
12 employing the gambit hopes, but does not know, that the dynamics and equilibria resulting
13 from the simplified strategy set and genetic system are a good indication of the dynamics and
14 equilibria of the actual phenotype, which is determined by a complex genetic system
15 interacting with ontogeny via messy bio-mechanical processes and evolving under significant
16 phylogenetic constraints.

17 Some experts have expressed some doubt of the usefulness of ESS analyses when the
18 actual genetic complexity is greater than the model allows for (e.g., Taylor 1979). There have
19 also been more passionate critics of the phenotypic gambit (Lewontin 1979; Gould &
20 Lewontin 1979; Gould 1982), as represented by the following quote:

21 “An adaptationist programme... based on faith in the power of natural selection...
22 proceeds by breaking an organism into unitary ‘traits’ and proposing an adaptive story for
23 each considered separately... We criticize this approach and attempt to reassert a
24 competing notion that organisms must be analysed as integrated wholes.” (Gould &
25 Lewontin 1979)

1 When the dust settled after this debate the mainstream conclusion was not that
2 organisms must be analyzed as integrated wholes but at least that care needs to be taken when
3 identifying the particular trait that is to be the focus of an adaptationist analysis (Dawkins
4 1982 pp. 33-54; Grafen 1984), though see Andrews et al. (2002). Next we shall present two
5 careful applications of the phenotypic gambit.

6 **3.4 The phenotypic gambit in practice**

7 To provide clear examples to contrast with puzzle paradigm models we discuss two
8 specific applications of the phenotypic gambit. One of these (Hamilton 1967) was written
9 prior to the formalization and entrenchment of the phenotypic gambit and is thus very explicit
10 about the various assumptions made. The second one (Kokko & Jennion 2008) was written
11 recently within a research culture where many of the critical assumptions of the phenotypic
12 gambit have become largely tacit. We shall now discuss how, despite the decades separating
13 these papers, they both share the same set of critical features ensuring that the gambit is not
14 made in vain: The phenotype is well-defined and can reasonably be studied in isolation; a
15 fitness function is derived from interaction between the phenotype and relevant ecological
16 factors; and the theoretical results yield interesting predictions that can be tested in empirical
17 work.

18 *3.4.1 The phenotype is well-defined and can reasonably be studied in isolation*

19 The first critical feature of a successful phenotypic gambit is a clear definition of the
20 phenotype to be explained. For Hamilton (1967) the phenotype of interest is the sex ratio at
21 birth, an unambiguous concept. For Kokko and Jennion (2008) the focal phenotype is
22 parental investment, specifically operationalized as the ratio of time spent caring for existent
23 offspring versus time spent seeking opportunities to create more offspring, again
24 unambiguous.

1 The phenotypes considered by Hamilton (1967) and Kokko and Jennion (2008) seem
2 reasonable to study in isolation. The "decision" of which sex ratio to have is made by every
3 (female or male depending) organism in a species every time it reproduces and this has direct
4 and relatively consistent fitness consequences at all times and places. The same is true of the
5 "decision" between investing in current existent offspring and investing in making more
6 offspring (although environmental changes can have significant impact here and so the theory
7 applies only when the relevant aspects of the environment are sufficiently stable). Barring
8 some deeply engrained phylogenetic or ontological constraints on variation, the fact that
9 these traits have direct and consistent fitness effects for every individual in every generation
10 suggests that the phenotypic gambit will produce reliable results.

11 *3.4.2 A fitness function is derived from interaction between the phenotype and relevant*
12 *ecological factors*

13 The next critical feature of a successful phenotypic gambit is a clear hypothesis about
14 what ecological factors should be considered, and a derivation of a fitness function
15 parameterized by these ecological factors. For instance, Hamilton (1967) hypothesized that
16 sex ratios would be affected by localized mate competition among males coupled with
17 population wide competition among females for breeding sites. The hypothesis of Kokko and
18 Jennion (2008) is that sexual divergence in parental care might be explained by differential
19 mortality rates when providing parental care versus searching or competing for mating
20 opportunities, by differences in parental certainty, and by feedbacks between these two
21 factors mediated through operational sex ratio. The researchers of these behaviors derived
22 fitness functions by making careful arguments about how the relevant ecological factors
23 should interact with strategies to affect the number of surviving offspring under different life
24 histories. Although the specific life history model of Kokko and Jennion (2008) is largely

1 implicit, those familiar with the field can reconstruct the explicit model from which the
2 various fitness functions are derived. For example, they implicitly assume that male and
3 female parental care can evolve independently so that males and females effectively form
4 independent populations.

5 In both papers the primary goal is to explain variation across species, though within
6 species variation is briefly discussed in both. Kokko and Jennion are careful to point out that
7 the factors that explain individual variation within species need not be the same factors that
8 explain between-species variation.

9 *3.4.3 The theoretical results yield interesting predictions that can be tested in empirical work*

10 When a fitness function has been derived, researchers can determine evolutionarily
11 stable strategies (although in 1967 when the paper of Hamilton was written the term ESS had
12 not yet been fully developed; foreshadowing future developments he phrased his arguments
13 in terms of “unbeatable” sex ratios). Under the crucial assumption that the evolutionary
14 process is close to equilibrium, this yields some predictions about the world. The model of
15 Hamilton predicts that the sex ratio of a species will be female biased if females compete
16 globally within a population for breeding sites but males compete locally (and hence often
17 with their siblings) for mating opportunities, and that the degree of this bias will increase with
18 the degree to which male offspring compete and breed with their siblings. These qualitative
19 predictions were compared with empirically observed patterns, allowing Hamilton (1967) to
20 conclude that “in its direction, the effect that Wylie and Jackson have independently reported
21 accords with the theory.” Similarly, Kokko and Jennion (2008) drew connections between
22 previous empirical observations and their model predictions of parental investment under
23 different combinations of the relevant ecological factors. In addition, the authors suggested
24 specific future empirical work based on the models presented.

1 It is a very valuable feature of a theory that it produces predictions (qualitative or
2 quantitative) that can be empirically refuted or verified. Theories with this feature have the
3 possibility of becoming part of a larger body of interlocking theoretical, observational, and
4 experimental work, as has been the case for both of these examples.

5 **4. Comparing the puzzle paradigm to the phenotypic gambit**

6 Thus far we have described the puzzle paradigm approach to studying human
7 cooperation (section 2) and the research strategy in evolutionary biology known as the
8 phenotypic gambit (section 3). We shall now compare the two.

9 Recall that the phenotypic gambit is to examine the evolutionary basis of a character as
10 if the simplest possible genetic system controlled it (Grafen 1984). This is exactly what the
11 puzzle paradigm does: It studies the evolution of cooperative vs. non-cooperative behavior as
12 if these behaviors were controlled by distinct alleles at a haploid locus (although the “as if”
13 part is typically not stated so explicitly). Because of this similarity, it is worth examining how
14 the puzzle paradigm fares with respect to the criteria for successful use of the phenotypic
15 gambit that were outlined in section 3.4.

16 **4.1 The phenotype of the puzzle paradigm is not well-defined and cannot reasonably be** 17 **studied in isolation**

18 Cooperation is the phenotype of interest in the puzzle paradigm, and is clearly defined
19 within this paradigm as paying a cost to bestow a benefit to others. It is crucial to note that
20 this definition is in terms of the payoff produced by the behavior. Thus the phenotype is not a
21 concrete behavior but rather an abstract category of behaviors organized according to their
22 fitness effects. This leads to two problems: The first problem is that the definition of
23 cooperation in terms of payoffs is ambiguous. The second problem is that strategies defined
24 by payoffs cannot reasonably be studied in isolation.

1 *4.1.1 The definition of cooperation is ambiguous*

2 The puzzle paradigm sets out to study the actual phenomenon of human cooperation.
3 However, it then provides a definition of cooperation that is identical to the definition of
4 altruism used by theoretical biologists, provided that the costs and benefits of the giving
5 game are in terms of individual reproductive success. Within the altruism research paradigm,
6 researchers begin by positing the existence of genes with fitness effects equivalent to
7 cooperation and defection in the giving game, and then attempt to model the evolutionary
8 dynamics of these genes within a given reproductive system, without regard for what the
9 altruistic behavior might be that causes these fitness effect or the causal relationship between
10 gene and behavior (e.g., Hamilton 1963). Without doing the hard work of linking such
11 abstract definitions to any actual behavior, we do not know what actual behavior the models
12 apply to. This is fine if the aim of the research is confined to the purely theoretical domain,
13 but not if the aim is to understand an actual phenomenon such as human cooperation.

14 In the four papers we have chosen to represent the puzzle paradigm no concrete
15 behaviors are referred to. Behaviors are only described in abstract terms such as “unselfish”,
16 “altruistic”, and “costly”, which are defined in reference to assumed payoffs. However, other
17 papers by some of the same authors have mentioned behaviors that they consider to be
18 cooperative, such as giving blood (Boyd & Richerson 1998) and giving to charity (Bowles &
19 Gintis 2002). The hard work then remains to show that these concrete behaviors satisfy the
20 abstract definition (e.g., that giving to charity reduces one’s fitness), and that they can
21 reasonably be studied in isolation.

22 This problem is exacerbated when puzzle paradigm research (unlike altruism research
23 in theoretical biology) remains uncommitted as to whether the underlying cause of
24 cooperative behaviors lies in the genetic or cultural domain. Puzzle paradigm models are

1 often claimed to capture both genetic and cultural evolutionary processes. Here are some
2 representative quotes:

3 “Cost and benefit are measured in terms of fitness. Reproduction can be genetic or
4 cultural.” (Nowak 2006).

5 “In evolutionary game theory it is not assumed that players are rational but only that
6 successful strategies spread—by being inherited, for instance, or copied through imitation
7 or learning.” (Nowak & Sigmund 2005).

8 “This score determines the player’s success in the subsequent imitation/reproduction stage,
9 either through replication and displacement or through neighbors imitating and adopting a
10 more successful strategy. Note that the two interpretations reduce to the same dynamics.”
11 (Brandt et al. 2003).

12 “It would be possible to construct an otherwise similar genetic model in which natural
13 selection played the same role that payoff biased imitation plays in the present model”
14 (Boyd et al. 2003).

15 Is it possible that the same model can capture both a genetic process and a cultural
16 process? In theory, yes, but only under extremely strong assumptions: It requires that the
17 cultural process and the genetic process act on the *same set of strategies* and that replication
18 is based on the *same payoffs* in both processes. Under these conditions, certain imitative
19 learning processes and a certain simple genetic replication process lead to identical dynamics
20 (Schuster & Sigmund 1983; Hofbauer & Sigmund 2003). However, in general there is no
21 reason to believe that cultural and genetic evolution act on the same strategies, nor that
22 replication is based on the same payoffs in the two processes. Indeed, some puzzle paradigm
23 researchers have themselves, in other work, pointed out that cultural transmission “can favor

1 different phenotypic variants than would be favored if the trait was genetically transmitted”
2 (Boyd & Richerson 1985, p. 183).

3 We shall return to other issues with the cultural interpretation in section 4.4. Our focus
4 here is that if it is not clear whether models describe genetic or cultural evolution it is even
5 more difficult to say which behaviors the definition of cooperation apply to. It is the
6 replication process, cultural or genetic, that defines what fitness is. Behaviors may have
7 different payoff consequences depending on whether payoffs represent genetic fitness or
8 cultural fitness. Hence, a behavior that does not count as cooperation in one process may
9 count as cooperation in the other process.

10 *4.1.2 The strategies cannot reasonably be study in isolation*

11 In puzzle paradigm models the replicating units are strategies in a meta-game. A meta-
12 game strategy includes, among other things, a rule about what strategy should be played in a
13 particular round of a giving game. For a given puzzle paradigm meta-game we can try to
14 infer an actual phenotype which corresponds to the strategies of that meta-game and then ask
15 if that phenotype can reasonably be studied in isolation. In particular, behaviors in this meta-
16 game must have a negligible correlation with behaviors in other strategic situations of fitness
17 relevance. If this is not the case, the analysis must be extended to an even more complex
18 meta-meta-game with an increasingly complex set of strategies.

19 Recall that giving games and cooperative strategies in such games are defined in terms
20 of payoffs. This makes isolation highly unlikely. After all, the only difference between
21 playing a giving game and other games, say a coordination game, is a small change in the
22 payoffs of the game. Isolation would therefore require that the same genes that have broad-
23 reaching systematic effects across a perceptually diverse set of giving games must cease to
24 have systematic effects across a potentially perceptually similar set of coordination games. In

1 other words, a gene is required to impact the willingness to “cooperate” when the payoff
2 structure is such that there is a net cost to self and a greater net benefit to others (as in a
3 giving game), but *not* impact willingness for the same behavior in circumstances where the
4 payoff structure is such that the relationship of these payoffs is contingent upon the behaviors
5 of other players (as in a coordination game). This seems highly unlikely given the evidence
6 for correlations of behavior across distinct strategic situations (e.g., Yamagishi et al. 2012).

7 **4.2 The fitness function is not derived from interaction between the phenotype and** 8 **relevant ecological factors**

9 The next critical feature of a successful phenotypic gambit is a clear hypothesis
10 concerning which ecological factors might explain variation in the phenotype of interest.
11 Section 3.4.2 demonstrated how successful applications of the phenotypic gambit focus on
12 concrete and straightforward life history factors. Puzzle paradigm papers offer up a variety of
13 more abstract hypotheses to the effect that if the giving game is iterated within a meta-game,
14 and the possibility of targeted punishment contingent upon behavior in the giving game is
15 also a possibility within the meta-game, then cooperation can emerge. Importantly,
16 punishment too is defined in terms of payoffs, as a behavior which reduces the payoff of the
17 punished party, typically at a cost to the punisher. In the case of costly punishment, simply
18 iterating rounds of the giving game followed by rounds of costly targeted punishment does
19 not allow for the emergence and maintenance of cooperation. Our canonical puzzle paradigm
20 papers go on to hypothesize that some additional factors may explain the evolution of
21 cooperation. These factors include inter-group conflict, which creates the possibility for
22 group selection (Boyd et al. 2003); reputation, which creates the possibility for indirect
23 reciprocity (Nowak & Sigmund 2005); optional participation in cooperative endeavors
24 (Hauert et al. 2007); economies of scale in either punishment or the giving game (Boyd et al.

1 2010); and the requisite communication and coordination of these various features. On a
2 general level, the hypotheses that human cooperative behavior is affected by these additional
3 factors are reasonable. However, the problem remains that in these models the core
4 ecological factor – punishment – suffers from being abstractly defined in terms of payoffs
5 instead of actual behavior.

6 To derive a fitness function, payoffs in the meta-game are simply calculated from the
7 payoffs of the simple giving game and the payoffs of punishment. Herein lies a critical
8 difference to the successful applications of the phenotypic gambit in section 3.4. Those
9 applications of the phenotypic gambit used simplifying assumptions about life history to
10 derive the fitness payoffs for concrete behaviors. No such connection is made in the puzzle
11 paradigm, because the model behaviors are already defined in terms of payoffs. Simply put,
12 the definition of cooperative behavior in terms of assumed payoffs creates a closed
13 theoretical loop which prevents the connection of actual behaviors and fitness.

14 **4.3 Puzzle paradigm models do not yield interesting predictions that can be tested in** 15 **empirical work**

16 With the fitness function in hand a researcher employing the phenotypic gambit is
17 empowered to make some predictions concerning the variation in the phenotype of interest
18 with respect to other potentially relevant ecological factors. For instance, recall that
19 Hamilton's (1967) analysis yielded the prediction that female bias in the sex ratio at birth will
20 increase with the degree to which male offspring compete and breed with their siblings. We
21 have seen no such exciting and directly testable prediction about the world come out of the
22 puzzle paradigm. The closest to predictions that our four sample papers get are general
23 statements such as:

1 “people should be less inclined to pay fixed than variable punishment costs” (Boyd et al.
2 2003)

3 “if the joint enterprise is optional, cooperation backed by punishment is more likely than if
4 the joint enterprise is obligatory.” (Hauert et al. 2007)

5 “the model predicts that only some individuals will engage in punishment.” (Boyd et al.
6 2010)

7 The lack of exciting predictions is to be expected given the closed theoretical loop of
8 the puzzle paradigm (section 4.2). The theory never links concrete ecological factors to
9 concrete behaviors.

10 **4.4 Shifting to cultural evolution does not salvage the puzzle paradigm**

11 In section 4.1.1 we discussed that puzzle paradigm papers tend to be ambiguous as to
12 whether the underlying replication process driving cooperative behavior is cultural or genetic.
13 We then went on to largely ignore cultural imitation processes, critiquing the puzzle
14 paradigm primarily under a genetic interpretation. As a final point we want to emphasize that
15 the puzzle paradigm is not salvaged by a shift to a cultural interpretation. There are two
16 senses in which such a simple shift is unjustified. To begin with, it would be wrong to
17 entirely abandon the genetic level. There are many interesting questions concerning the
18 genetic evolutionary history of human behavioral plasticity itself, as well as of the various
19 more general cognitive capacities that enable culture and cooperation (Heyes 2012;
20 Tomasello et al. 2005).

21 An even stronger reason not to shift the evolutionary puzzle of human cooperation to
22 cultural evolution is that the shift is not valid. Evolution at the cultural level is different from
23 evolution at the genetic level with the fundamental concepts of the latter, such as replication
24 and fitness, failing to carry over to culture in any well-defined way (Gabora 2011; Henrich et

1 al. 2008). Whereas genetic fitness can be roughly measured by looking at individuals and
2 counting their surviving descendants some time later, there is no corresponding way to
3 measure cultural fitness by looking at an individual and counting how many others have
4 copied her behavior. Although the term “payoff biased imitation” and the payoff based
5 replication it implies may sound a reasonable thing, its application to puzzle paradigm models
6 of human cooperation requires several strong assumptions to be met. One requirement is that
7 people must be able to observe other individuals’ strategies in the focal repeated game.
8 However, it is generally not possible to infer complex strategies in repeated games from
9 observed behavior, because the same observed behavior can be generated by many different
10 strategies. Thus, already this fundamental requirement is unlikely to be met. The following
11 sections discuss two additional assumptions that are necessary but unlikely to hold.

12 *4.4.1 Imitation of someone else’s strategy must be the dominant mode of behavioral change*

13 If the strategies in puzzle paradigm models are to be interpreted as cultural replicators,
14 they are learned and possibly relearned many times over a lifetime. But learning is not
15 necessarily imitation. Adoption of a behavior can occur through many other pathways.
16 Learning in any one situation may spill over to other situations; e.g., it is conceivable that
17 someone who goes through a vaccination program might become more inclined to give blood
18 because of the contextual similarity between the two situations. An individual may also learn
19 general rules such as “always be considerate of others”, which would clearly affect behavior
20 in many situations. Thus, a behavioral change in one situation may be brought about by
21 cultural learning in several strategically distinct situations, and at different levels of
22 generality. This stands in stark contrast with the assumption that behavior in a giving game
23 situation changes mainly due to individuals adopting a new meta-game strategy in a unique
24 meta-game.

1 Or consider the impact on giving behavior of the various identities of individuals, such
2 as their personal identity, familial identity, and community identity. Aaker and Akutsu (2009)
3 discuss how different identities may range from specific to broad, may be activated in
4 different contexts, may interact with emotions, may shift over the life-span, may be
5 reinforced or change due to one's own choices to give or not, and may generally color the
6 way in which one sees the world. See also Weber et al. (2004). It is difficult to see how the
7 impact of context-activated identities on giving can be reconciled with the notion of imitation
8 of meta-game strategies.

9 *4.4.2 People must be able to observe other individuals' total payoffs in the focal repeated*
10 *game and preferentially imitate high-payoff strategies*

11 So far we have argued that it is unlikely that meta-game strategies can be observed so
12 that they can be imitated. Moreover, we have argued that even if they could be observed, it is
13 unlikely that behavioral change is mainly driven by imitation of such strategies. Our last
14 point is that even if meta-game strategies are imitated, it is unlikely that people base such
15 imitation on accurate observations of cultural fitness payoffs that individuals accrue.

16 The raw payoffs of behaviors come in a variety of currencies (e.g., resources, effort,
17 pain and pleasure, and reputation). Preferences for these currencies vary. To the extent that
18 individuals at all make decisions by adding up these incommensurable terms to a single total
19 payoff, they will surely do so in idiosyncratic and highly subjective ways. These subjective
20 total payoffs are not observable to others. To the extent that some component of the total
21 payoff is observable (e.g., an individual's wealth or level of happiness), this payoff
22 component will be subject to contributions from other parts of the individual's life as well. It
23 seems quite impossible for others to correctly estimate what proportion of an individual's

1 wealth or happiness, say, should be attributed to the individual's behavior in the focal
2 repeated game.

3 **5. The evolutionary puzzle of human cooperation is a red herring**

4 In section 4 we used a comparison with the phenotypic gambit in evolutionary biology
5 to examine the puzzle paradigm for studying human cooperation. The comparison highlighted
6 several weaknesses in puzzle paradigm research. Here we shall argue that very premise of the
7 paradigm – that human cooperation is an evolutionary puzzle – is questionable.

8 To start with, it is questionable which specific puzzle the paradigm is meant to address.
9 On the one hand, all the papers we have quoted from include various statements indicating
10 that the primary variation of interest is inter-species variation with humans as outliers (e.g.,
11 “Unlike any other species, humans cooperate with non-kin in large groups”, Boyd et al.
12 2003). On the other hand, other statements made in the same papers explicitly refer to within-
13 species variation instead (see quotes in section 4.3). Thus, despite the references to human
14 uniqueness we believe puzzle paradigm research means to address the within-species puzzle
15 of why altruistic behavior has not been selected against. It is questionable to frame this
16 question as an evolutionary puzzle of human cooperation to be solved by the identification of
17 the correct meta-game. To do so requires a number of tacit assumptions:

- 18 1. The expression of altruistic behavior is the output of identifiable strategies that
19 serve as replicators in an evolutionary process. (Otherwise it cannot be selected
20 against.)
- 21 2. For these strategies, and the altruistic behaviors they result in, there are measurable
22 fitness payoffs such that the altruistic behaviors have negative payoff but the
23 strategy as a whole has not. (Otherwise there is no puzzle and no solution,
24 respectively.)

- 1 3. The evolutionary process with respect to these strategies has reached an
2 equilibrium. (Otherwise selection may just be ongoing.)
- 3 4. The kind of cooperation at which humans excel amounts, in fact, to altruistic
4 behavior. (Otherwise the puzzle is not about human cooperation.)

5

6 With respect to the last point, it is certainly not clear that human cooperation is
7 characterized by altruistic behavior. In a book devoted to human cooperation, Argyle (1991)
8 explicitly rejects the notion of cooperation as a game theoretic strategy and instead defines it
9 as “acting together in a coordinated way at work, leisure or in social relationships, in the
10 pursuit of shared goals, the enjoyment of the joint activity, or simply furthering relationship”.

11 Assuming the goal is restricted to understand human altruism and not cooperation in this
12 wider sense, we still believe it is more fruitful to question the first three assumptions than to
13 take them for granted. Questioning these assumptions points to a number of intriguing
14 empirical research questions. These have the potential to result in genuine knowledge and
15 inform future attempts at theoretical modeling.

16 **5.1 Can behavior be captured as the output of replicating strategies?**

17 At the heart of game theoretic models lies the assumption of a certain set of strategies.
18 As illustrated by puzzle paradigm models, strategies may be simple (“in situation A do X”) or
19 of greater complexity (e.g., “in situation A do X if C holds” or “in situation A do X and in
20 situation B do Y”). It is a genuine challenge to assess if this concept of strategy can capture
21 how people actually behave and, if so, what strategies are actually used. There are two main
22 reasons why this is challenging. First, observations of behavior may be consistent with any
23 number of different strategies (e.g., an observation of behavior X in situation A would be
24 consistent with any of the above three examples of strategies). Second, there is no

1 unambiguous way to categorize situations or behaviors in the first place. In other words,
2 given that no two situations or behaviors are ever perfectly identical, it may not be clear
3 which should be taken as instances of A and X. Somehow these difficulties must be faced if a
4 given game theory model is ever to convince those who do not a priori believe in the validity
5 of its strategy set. The notion of strategies as units of replication suggests one way in which
6 their complexity in terms of scope should be investigated, namely, analysis of how clusters of
7 behaviors tend to vary together between individuals and change together within individuals.
8 Such clusters are candidates for units of replication.

9 **5.2 What are the payoffs?**

10 What is the actual payoff, to self and others, of various strategies? This question has
11 several layers of complexity. One problem is that it is not clear what type of payoff to
12 measure. It is possible, though difficult, to measure both genetically relevant payoff (number
13 of surviving children or grandchildren) and various other payoff currencies that people may
14 care about, such as health, wealth, and life satisfaction. A more daunting problem is how to
15 assess the link between these payoffs and the replication of specific strategies (assuming
16 these have been identified, see section 5.1). This is comparable to the notoriously difficult
17 problem of accurately assessing health risks and benefits of behaviors such as consumption of
18 certain food or drink. Perhaps research methods from that field can be adapted to behaviors
19 related to cooperation.

20 **5.3 Are altruistic behaviors at equilibrium?**

21 According to the most recent official Swedish statistics
22 (<https://geblod.nu/fakta/blodgivningsstatistik>, accessed 4th July 2016), the number of blood
23 donors in Sweden has declined every year between 2010 and 2015. The total decline in just
24 five years is almost 12%, despite the population growing by almost 5% in the same period.

1 These data suggest that blood donations in Sweden are not at equilibrium. Before searching
2 for equilibria in models, it seems worthwhile to study actual empirical trends in behavior and
3 search for their causes.

4 **5.4 Conclusion**

5 In this essay we have criticized an approach to human cooperative behavior where
6 evolutionary explanations are proposed and studied with mathematical models. There have
7 been previous critiques of misapplication of evolutionary theory to human cooperation (e.g.,
8 El Mouden et al. 2012; Scott-Phillips et al. 2011; West et al. 2007; West et al. 2011). We
9 believe that these critiques do not go far enough. Ours is more in the spirit of Arnold's (2013)
10 recent critique of simulations of the repeated prisoner's dilemma to study the evolution of
11 cooperation in that it questions the scientific value of an entire paradigm.

12 Let us emphasize that we, the authors, love both mathematics and evolutionary theory.
13 We believe that mathematical modeling of evolutionary processes is sometimes an extremely
14 useful research tool—but only sometimes. It is unlikely to be useful if the modeling paradigm
15 is chosen without careful consideration of whether it fits the phenomenon one wants to study.
16 It is fine to study models of payoff based replication of giving game strategies for their own
17 sake. But students of human behavior should choose that particular modeling paradigm only
18 if it is able to capture what is really going on in human behavior.

19 **References**

- 20 1. Aaker, J.L. & Akutsu, S. (2009) Why do people give? The role of identity in giving.
21 *Journal of Consumer Psychology* 19:267–270.
- 22 2. Andrews, P., Gangestad, S. & Matthews, D. (2002) Adaptionism – how to carry out
23 an adaptionist program. *Behavioral and Brain Sciences* 25:489–553.
- 24 3. Argyle, M. (1991) *Cooperation: The Basis of Sociability*. Routledge.

- 1 4. Arnold, E. (2013). Simulation models of the evolution of cooperation as proofs of
2 logical possibilities. How useful are they?. *Ethics & Politics*, 2, 101-138.
- 3 5. Bowles, S. & Gintis, H. (2002) Behavioural science: homo reciprocans. *Nature* 415
4 (6868):125–127.
- 5 6. Boyd, R., Gintis, H. & Bowles, S. (2010) Coordinated punishment of defectors
6 sustains cooperation and can proliferate when rare. *Science* 328 (5978):617–620.
- 7 7. Boyd, R., Gintis, H., Bowles, S. & Richerson, P. (2003) The evolution of altruistic
8 punishment. *Proceedings of the National Academy of Sciences* 100 (6):3531–3535.
- 9 8. Boyd, R. & Richerson, P. (1988) The evolution of reciprocity in sizable groups.
10 *Journal of Theoretical Biology* 132 (3):337–356.
- 11 9. Brandt, H., Hauert, C. & Sigmund, K. (2003) Punishment and reputation in spatial
12 public goods games. *Proceedings of the Royal Society of London. Series B:*
13 *Biological Sciences* 270 (1519):1099–1104.
- 14 10. Dawkins, R. (1999) *The extended phenotype: The long reach of the gene*. Oxford
15 University Press, USA.
- 16 11. El Mouden, C., Burton-Chellew, M., Gardner, A. & West, S.A. (2012) What do
17 humans maximise? In: *Evolution and Rationality*, ed. S. Okasha & K. Binmore, [23-
18 49]. Cambridge University Press.
- 19 12. Fisher, R., 1930. *The Theory of Natural Selection*. Oxford University Press, London.
- 20 13. Fletcher, J. A., & Doebeli, M. (2006). How altruism evolves: assortment and
21 synergy. *Journal of Evolutionary Biology*, 19(5): 1389-1393.
- 22 14. Gabora, L. (2011) Five clarifications about cultural evolution. *Journal of Cognition*
23 *and Culture* 11:61– 83.
- 24 15. Gardner A. (2013) Ultimate explanations concern the adaptive rationale for organism
25 design. *Biology and Philosophy* 28:787–791.

- 1 16. Gould, S. (1982) Darwinism and the expansion of evolutionary theory. *Science* 216
2 (4544):380–387.
- 3 17. Gould, S. & Lewontin, R. (1979) The spandrels of san marco and the panglossian
4 paradigm: a critique of the adaptationist programme. *Proceedings of the Royal
5 Society of London Series B. Biological Sciences* 205 (1161): 581–598.
- 6 18. Grafen, A., 1984. Natural selection, kin selection and group selection. In: *Behavioural
7 ecology: an evolutionary approach 2*, eds. J. R. Krebs & N. B. Davies, [62-84].
8 Cambridge University Press.
- 9 19. Hamilton, W.D. (1963) The evolution of altruistic behavior. *The American Naturalist*
10 97 (896): 354–356
- 11 20. Hamilton, W. D. (1967) Extraordinary sex ratios. *Science* 156 (3774): 477–488.
- 12 21. Hauert, C., Traulsen, A., Brandt, H., Nowak, M. & Sigmund, K. (2007) Via freedom
13 to coercion: the emergence of costly punishment. *Science* 316 (5833): 1905–1907.
- 14 22. Henrich, J., Boyd, R., & Richerson, P. J. (2008) Five misunderstandings about
15 cultural evolution. *Human Nature*, 19 (2):119–137.
- 16 23. Heyes, C. (2012) New thinking: the evolution of human cognition. *Philosophical
17 Transactions of the Royal Society B: Biological Sciences* 367 (1599):2091–2096.
- 18 24. Higginson, A., McNamara, J. & Houston, A. (2012) The starvation-predation trade-
19 off predicts trends in body size, muscularity, and adiposity between and within taxa.
20 *The American Naturalist*, 179 (3): 338–350.
- 21 25. Hofbauer, J. & Sigmund, K. (2003) Evolutionary game dynamics. *Bulletin of the
22 American Mathematical Society* 40 (4): 479–519.
- 23 26. Kokko, H. & Jennions, M. D. (2008). Parental investment, sexual selection and sex
24 ratios. *Journal of Evolutionary Biology* 21 (4): 919–948.

- 1 27. Lewontin, R. (1979) Sociobiology as an adaptationist program. *Behavioral Science* 24
2 (1): 5–14.
- 3 28. Nowak, M., 2006. Five rules for the evolution of cooperation. *Science* 314 (5805):
4 1560–1563.
- 5 29. Nowak, M. A. & Sigmund, K. (2005) Evolution of indirect reciprocity. *Nature* 437
6 (7063): 1291–1298.
- 7 30. Pen, I., & Taylor, P. D. (2005) Modelling information exchange in worker–queen
8 conflict over sex allocation. *Proceedings of the Royal Society B: Biological Sciences*,
9 272 (1579): 2403–2408.
- 10 31. Schuster, P. & Sigmund, K. (1983) Replicator dynamics. *Journal of Theoretical*
11 *Biology* 100 (3), 533–538.
- 12 32. Scott-Phillips, T., Dickins, T. & West, S. (2011) Evolutionary theory and the
13 ultimate–proximate distinction in the human behavioral sciences. *Perspectives on*
14 *Psychological Science* 6 (1), 38–47.
- 15 33. Smith, J. (1978) Optimization theory in evolution. *Annual Review of Ecology and*
16 *Systematics* 9, 31–56.
- 17 34. Smith, J. & Price, G. (1973) The logic of animal conflict. *Nature* 246: 15–18.
- 18 35. Taylor, P. & Jonker, L. (1978) Evolutionary stable strategies and game dynamics.
19 *Mathematical Biosciences* 40 (1), 145–156.
- 20 36. Taylor, P. D. (1979) Evolutionarily stable strategies with two types of player. *Journal*
21 *of Applied Probability* 16: 76–83.
- 22 37. Tomasello, M., Carpenter, M., Call, J., Behne, T. & Moll, H. (2005) Understanding
23 and sharing intentions: The origins of cultural cognition. *Behavioral and Brain*
24 *Sciences*, 28 (5): 675–690.

- 1 38. van Veelen, M., García, J., Rand, D. G. & Nowak, M. A. (2012) Direct reciprocity in
2 structured populations. *Proceedings of the National Academy of Sciences* 109
3 (25):9929–9934.
- 4 39. Weber, J.M., Kopelman, S. & Messick, D.M. (2004) A conceptual review of decision
5 making in social dilemmas: Applying a logic of appropriateness. *Personality and*
6 *Social Psychology Review* 8:281–307.
- 7 40. West, S., El Mouden, C. & Gardner, A. (2011) Sixteen common misconceptions
8 41. about the evolution of cooperation in humans. *Evolution and Human Behavior* 32
9 (4):231–262.
- 10 42. West, S., Griffin, A. & Gardner, A. (2007) Social semantics: altruism, cooperation,
11 mutualism, strong reciprocity and group selection. *Journal of Evolutionary Biology* 20
12 (2):415–432.
- 13 43. Yamagishi, T., Horita, Y., Mifune, N., Hashimoto, H., Li, Y., Shinada, M., Miura, A.,
14 Inukai, K., Takagishi, H. & Simunovic, D. (2012). Rejection of unfair offers in the
15 ultimatum game is no evidence of strong reciprocity. *Proceedings of the National*
16 *Academy of Sciences* 109 (50):20364–20368.