

Evolution of cooperation and skew under imperfect information

Erol Akçay^{a,1}, Adam Meirowitz^b, Kristopher W. Ramsay^b, and Simon A. Levin^{a,1}

Departments of ^aEcology and Evolutionary Biology and ^bPolitics, Princeton University, Princeton, NJ 08544

Contributed by Simon A. Levin, July 27, 2012 (sent for review March 26, 2012)

The evolution of cooperation in nature and human societies depends crucially on how the benefits from cooperation are divided and whether individuals have complete information about their payoffs. We tackle these questions by adopting a methodology from economics called mechanism design. Focusing on reproductive skew as a case study, we show that full cooperation may not be achievable due to private information over individuals' outside options, regardless of the details of the specific biological or social interaction. Further, we consider how the structure of the interaction can evolve to promote the maximum amount of cooperation in the face of the informational constraints. Our results point to a distinct avenue for investigating how cooperation can evolve when the division of benefits is flexible and individuals have private information.

other-regarding preferences | social evolution | incentive compatibility | reproductive transactions | cheap-talk bargaining

Cooperative interactions drive much of the ecological, evolutionary, and social dynamics of organisms ranging from soil bacteria to primates, including—and especially—humans. Whereas much theory focuses on various mechanisms that promote cooperative behaviors (1–6), some fundamental questions remain unresolved. Among them is how the benefits of cooperation are to be divided among cooperating agents. Most theoretical work conceives of cooperation as a binary affair with payoffs to individuals from each outcome set a priori. However, frequently, the surplus from cooperation, whether it is the kill of a cooperatively hunting group or the reproductive output of a breeding group, can be partitioned among individuals in different ways; and how this division is achieved affects how likely individuals are to cooperate. Further, most research on biological cooperation focuses implicitly or explicitly on situations where individuals make decisions under perfect information of their and others' payoffs. However, private information, where some individuals have access to information and others do not, is a feature of many biological and social interactions. Although private information has been studied in a few specific contexts before, including mate choice (7, 8), parental care (9, 10), and animal conflicts (11), the role of private information in the evolution of cooperation in general remains understudied.

We introduce a distinct approach to biology to study how cooperation can be maintained when the division of benefits is flexible and individuals have private information. This approach, called mechanism design (12) and borrowed from economics, inverts the standard methodology of game-theoretic modeling. Instead of specifying a particular game and analyzing its equilibria, we analyze the properties of equilibrium outcomes in a large class of games and also ask what the consequences of different game structures are for the fitness of different individuals and the group's reproductive output.

As a case study, we use a problem of central importance to behavioral ecology and social evolution: the partitioning of reproduction, or reproductive skew, within a breeding group. A large body of work in behavioral ecology aims to understand the evolution of reproductive skew (13) as a function of demographic, individual, and ecological variables. However, patterns of reproductive skew remain contradictory: A recent review (14) concludes that whereas theory explains between-species patterns with some success, within-species patterns of skew often

do not conform to theoretical predictions. We suggest that these failures occur because existing theory (14–16) assumes that reproductive skew evolves under perfect information about all relevant variables (17). In reality, however, individuals might be expected to have private information about themselves or the environment, which as we show dramatically affects both the scope of cooperation and the division of the benefits when cooperating. A related problem is that the proliferation of models in skew theory, driven in part by the empirical difficulties, has resulted in a situation where many contradictory patterns can be predicted, depending on the details of the model (14). Together with a systematic theory of which models apply in different settings, this could be a desirable property, but there is currently no such theory; hence the abundance of models fails to generate the clarity that theory is supposed to provide. Our approach avoids this problem by obtaining results independent of the precise game structure for a large class of games and also provides a first step in asking how the transactions game itself might evolve.

Our basic setup is a twist on the canonical reproductive skew model. Consider two individuals, labeled 1 and 2, who have the option of forming a group and breeding together or breeding alone. Label their expected success when breeding alone—their outside options—as o_1 and o_2 , respectively. We depart from the canonical model in assuming that these options are not observed directly by both individuals; individual 1 only “knows,” i.e., can condition its behavior on, o_1 , but cannot condition on o_2 , and vice versa. The outside options are distributed according to some probability distribution, and hence natural selection will lead to optimal strategies according to their expected fitness consequences. This assumption differs from that in previous models, where the outside options are commonly known; hence individual 1's strategy can be conditional on the lowest share 2 will accept, and vice versa. This is not possible in our setting and thus optimal strategies will typically miss some mutually beneficial opportunities for cooperation. If the individuals form the group, they can obtain a joint breeding success of Ω . This joint reproductive success is to be divided between 1 and 2 through some sort of game; our goal is to study which groups and divisions of reproduction are compatible with evolutionary stability of strategies in any kind of game and what game structures can implement the optimal outcome. For most of the following, it is more convenient to work with the potential losses and gains from group formation, defined as $g_1 = \Omega - o_1$ (the reproduction individual 1 gives up by entering the group) and $g_2 = \Omega - o_2$ (the reproduction individual 2 could potentially gain by entering the group). We sometimes call g_1 and g_2 players' “types.” Suppose that $g_1 \in [a_1, b_1]$ and $g_2 \in [a_2, b_2]$ are distributed according to $f_1(g_1)$ and $f_2(g_2)$, with cumulative distributions $F_1(\cdot)$ and $F_2(\cdot)$, respectively. We assume that these distributions are attributes of the environmental variation, i.e., do not change with the strategies of individuals or the game structure. Hence,

Author contributions: E.A., A.M., K.W.R., and S.A.L. designed research; E.A., A.M., and K.W.R. performed research; E.A., A.M., K.W.R., and S.A.L. contributed new reagents/analytic tools; and E.A., A.M., K.W.R., and S.A.L. wrote the paper.

The authors declare no conflict of interest.

¹To whom correspondence may be addressed. E-mail: eakcay@princeton.edu or slevin@princeton.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1212925109/-DCSupplemental.

natural selection can optimize individual strategies with regard to their expected fitness under these distributions. (In the rationality-based language of economic game theory, these distributions would have been “common knowledge” in a Bayesian game; see *SI Text* for more on the correspondence between evolutionary stability and Bayesian Nash equilibria.) Finally, individuals are related to each other with a coefficient of relatedness, r , so that they gain indirect fitness benefits from each other’s reproduction. The relatedness in a group is a function of demographic parameters (18), including the reproductive skew in the population (19); as a simplifying assumption, we take r as constant and to be the appropriately scaled relatedness that is consistent with the equilibrium level of skew in our model.

Overall, our setting is analogous to the bilateral trade problem with private information in economics (20), where the trade corresponds to formation of groups, and the payment corresponds to the allocation of the reproduction within the group. The important analytical difference is that previous work in economics assumes null relatedness, which as we show has important consequences. As such, our results for $r > 0$ also contribute to the economics literature.

If individuals’ actions could be conditioned on both outside options o_1 and o_2 , then it would be mutually beneficial for groups to form (with a division that makes both players better off) when $o_1 + o_2 < \Omega$, which translates to $g_2 > g_1$. With private information, however, individuals cannot compute this condition. They could potentially signal their outside options, but such signaling will not, in general, be evolutionarily stable. An individual with higher outside option can demand a higher share of the reproduction, and therefore individuals have an incentive to behave as if their outside options were higher. We show that this incentive problem can preclude cooperation even in cases where it is mutually beneficial, unless individuals are highly related to each other. Furthermore, the possibility of costly signaling cannot always remedy this problem.

We begin by specifying (i) an information structure for an interaction (i.e., which variables individual strategies can be conditioned on), (ii) a set of feasible games (i.e., mappings from combinations of actions to different payoffs), (iii) an equilibrium concept (e.g., evolutionary stability), and (iv) an objective function (e.g., maximizing the probability of cooperation). We obtain two kinds of results: First, we find the properties of evolutionarily stable outcomes in any possible game. Second, we describe a family of games that include the most salient transactional models and describe equilibrium play and its fitness consequences in this family of games as a first step to consider how the structure of the social interaction can evolve.

Our first kind of result is made possible by a celebrated theorem from economics, the revelation principle (21), which allows us to focus on a special class of simple games—direct mechanisms—that can represent all possible equilibria in all possible games. A direct mechanism is the simplest possible game structure, in which individuals simply send a message reporting their information to an (imagined) central arbiter and get assigned a payoff on the basis of the messages received by the arbiter. (A “mechanism” simply refers to a game or, more precisely, a mapping from combinations of strategies to the vector of payoffs to the players.) The function that determines the players’ payoffs (assumed to be known to the players before sending their messages) determines whether players will find it optimal to report their information truthfully instead of misrepresenting it. If truth telling is optimal, then the mechanism is called incentive compatible. The revelation principle (21) states that all Bayesian Nash equilibria [a necessary condition for evolutionarily stable strategies (*SI Text*)] to any game of imperfect information can be represented by incentive-compatible direct mechanisms. Thus, determining whether an outcome is possible in a direct mechanism tells us whether that outcome could ever be the result of equilibrium behavior in any evolutionary game. (See more on the revelation principle in *SI Text*.)

Denoting the players’ reports of their outside options by θ_1 and θ_2 , respectively, we characterize a direct mechanism with two

functions of these reports. The first function, $p(\theta_1, \theta_2)$ denotes the probability of group formation [$p(\cdot, \cdot)$ could be binary or continuous], and the second function, $x(\theta_1, \theta_2)$ gives the share of reproduction allocated to individual 1 (and taken from individual 2, and hence $x(\cdot, \cdot)$ can be viewed as a “payment”). For this section we are interested in the most general class of games, so we do not necessarily require that the payment is made only when the group forms (in the section *Nondirect Mechanisms and Incomplete Control*, we investigate a class of games that do have this more realistic property). Hence, $x(\theta_1, \theta_2)$ is the expected payment to individual 1 when the reports are θ_1, θ_2 , regardless of whether the group forms or not. Given a direct mechanism with functions (p, x) , individual 1’s expected change in inclusive fitness when it has a loss of g_1 and reports θ_1 is given by

$$W_1(g_1, \theta_1) = \int_{a_2}^{b_2} ((1-r)x(\theta_1, g_2) - (g_1 - rg_2)p(\theta_1, g_2))f_2(g_2)dg_2. \quad [1]$$

[It should be noted that inclusive fitness calculations may in general be dependent on the frequency of different genotypes in the population when there are nonadditive fitness effects (22). However, with weak selection or small-effect mutants, additivity can be approximately restored and inclusive fitness calculations become approximately accurate. Here, we make use of this approximation, as our interest is to compute the optimal conditional strategies given the strategic interaction and not the dynamics of a particular set of genotypes.] Similarly, individual 2’s expected change in inclusive fitness is given by

$$W_2(g_2, \theta_2) = \int_{a_1}^{b_1} ((g_2 - rg_1)p(g_1, \theta_2) - (1-r)x(g_1, \theta_2))f_1(g_1)dg_1. \quad [2]$$

A mechanism is said to be incentive compatible (IC) if $W_1(g_1, g_1) \geq W_1(g_1, \theta_1)$ for all $\theta_1 \in [a_1, b_1]$ and $W_2(g_2, g_2) \geq W_2(g_2, \theta_2)$ for all $\theta_2 \in [a_2, b_2]$. Furthermore, we require that participation in the game (i.e., sending the reports and accepting the outcome of the arbiter) is voluntary: Individuals can opt out of the interaction altogether and breed alone (thus obtaining their outside option) if their expected gains from the interaction are negative. Hence, in addition to IC, we require from our mechanism that $W_1(g_1, g_1) \geq 0$ and $W_2(g_2, g_2) \geq 0$ for all g_1 and g_2 . Note that this condition applies after individuals know their own outside option, but before they have learned their partner’s, so we call it the interim participation constraint (IPC). In *SI Text*, we provide the necessary and sufficient conditions for a mechanism to be both IC and IPC.

Our first major result concerns whether any game exists that ensures cooperation in all cases that are mutually beneficial (i.e., whenever $g_2 > g_1$). The Myerson–Satterthwaite theorem answers this question in the negative: If there are any pairs of individuals for whom cooperation is not mutually beneficial (i.e., the distributions of g_1 and g_2 overlap), there are no mechanisms that guarantee that all groups that are mutually beneficial will form. However, as we show below, this result is potentially changed when individuals are related and maximize their inclusive fitness (alternatively, one can think of the game being played by two agents with other-regarding preferences where the level of other regard is parameterized by r). In particular, assuming that $b_1 \geq b_2$ (see *SI Text* for $b_2 > b_1$), full cooperation becomes possible when

$$r \int_{a_1}^{a_1+b_2-a_2} F_1(t)(1-F_2(t)) dt \geq \int_{a_2}^{b_2} F_1(t)(1-F_2(t)) dt. \quad [3]$$

Thus, with high enough relatedness between the individuals, a game exists that ensures the individuals will form a group whenever it is mutually beneficial, a condition we term full cooperation,

or efficiency. If relatedness is below that threshold, some cooperation will still take place, but some pairs of individuals that would be better off in a group will nevertheless not cooperate at equilibrium, because of the incentives to misrepresent their private information. Note that the level of relatedness needed to ensure full efficiency can be very high, even exceeding 1. In the special case where both 1 and 2's outside options come from the same distribution, the integrals on both sides of Eq. 3 are identical to each other, and full cooperation becomes possible only when $r = 1$, i.e., with clonal groups.

In biology the study of costly signaling has become a common approach to dealing with informational inefficiencies (7–9). In the context of mate choice and parent–offspring relations, costly signaling can indeed restore efficiency [e.g., choosing the optimal mate (8) or optimally provisioning for offspring (23–25)], so one might be tempted to claim that allowing individuals to signal their outside options before the interaction might allow them to circumvent the inefficiency due to private information. This intuition turns out to be only partially correct: When signal costs depend only on the message (θ_i) and not on individuals' true type (g_i) (this is the case in models of parent–offspring signaling), we can show (SI Text) that all outcomes that can be implemented in any game with potential for costly signaling can also be implemented in a game without costly signaling, along with some that cannot be implemented with costly signals. Hence, costly signaling where costs depend only on the signal value cannot remedy the inefficiencies identified above. The reason is that the impossibility of sustaining inefficiency does not result from an inability to incentivize truth telling per se. Rather, the reason is that the payments required for truth telling become so high that they make some individuals better off not participating in the game at all, violating the IPC. Signaling with costs only exacerbates this problem, because the signal costs (unlike payments) are wasted and do not supply an inclusive fitness benefit. The same applies to games where parties can invest in a costly competitive behavior to grab a portion of reproduction for themselves (e.g., compromise models) (16, 26).

When signals can depend on the true type of the individual (e.g., individuals with high type pay a relatively lower cost for the same signal) (8, 27), the effect of costly signaling depends on the exact nature of the cost function. It is easy to show functions always exist that allow truthful revelation at little or no cost in equilibrium (28), effectively removing the private information problem and restoring full efficiency. However, if the signaling mechanism inevitably leads to costs at the equilibrium path (as is commonly presumed), the participation constraint might again preclude full efficiency. More generally, mechanism design can be used to learn about the types of cost function that can implement efficient outcomes. However, any sharper prediction about whether efficient cooperation can be maintained would require careful accounting of biological constraints specific to the interaction. Likewise, if the costly actions of one individual actually provide a benefit to the other, it is easy to show that full efficiency can be restored provided that the benefits are high enough. However, such a solution arguably redefines the problem, because it implies that cooperative benefits are not all accounted for in the original problem statement.

Given that full cooperation may not be possible for some values of relatedness between individuals, we next consider the properties of the interaction that maximizes the aggregate expected gains to the individuals, i.e., the expected value of $W_1(g_1, g_1) + W_2(g_2, g_2)$. This question can be answered by maximizing the expected probability that groups form,

$$\bar{p} = \int_{a_2}^{b_2} \int_{a_1}^{b_1} p(g_2, g_1) f_2(g_2) f_1(g_1) dg_1 dg_2, \quad [4]$$

subject to the constraints $g_2 \geq g_1$ and inequality Eq. S19 in SI Text. We can show that the mechanism maximizing the gains from cooperation prescribes that groups should form whenever

$$g_2 - g_1 \geq \max\left(\frac{1}{1+r} \left[\frac{1-F_2(g_2)}{f_2(g_2)} + \frac{F_1(g_1)}{f_1(g_1)} - \frac{1}{\lambda} \right], 0\right), \quad [5]$$

where $\lambda > 0$ is a Lagrange multiplier, which is determined by substituting the function Eq. 5 into the constraint Eq. S19 and satisfying that constraint with equality. The solution Eq. 5 can be interpreted as follows: In the space of gains (Fig. 1), the *Upper Left* region (shaded) consists of cases where the group is efficient and forms in equilibrium, and the *Lower Right* region (below the dashed diagonal) consists of cases where the group is not efficient and does not form. The slice of g_1 - g_2 pairs in between the shaded region and the dashed diagonal corresponds to cases where the group does not form even though it is mutually beneficial. The intuition for why groups do not form in this region is that each individual has relatively high outside options and therefore “demands” too high a share of the group production to consent to cooperation. The size of the slice between the shaded region and the diagonal in Fig. 1 captures the extent to which the incentive problem causes inefficiencies and is determined by the distributions and r .

An Example

To apply the general results shown above, consider the scenario where a dominant and a subordinate individual interact to determine their share of the group reproduction. Normalize the group's productivity to $\Omega = 1$ and assume that the outside options o_1 and o_2 are uniformly distributed over $[w_d, 1]$ and $[0, w_s]$, respectively. In other words, w_d is the lowest possible reproduction that a dominant can expect to achieve on its own, and w_s is the highest possible reproduction for a subordinate. To focus on a particular interesting set of cases we assume $w_s \leq 0.5 \leq w_d$ and $1 - w_s \leq w_d$. With these assumptions, the optimal mechanism maximizing the total benefit while being IC and IPC (i.e., the solution Eq. 5) prescribes that groups will form when

$$g_2 - g_1 \geq \frac{(1-r)(1-w_d)}{2(2+r)}. \quad [6]$$

Thus, unless the two individuals are clones of each other with $r = 1$, some combinations of dominants and subordinates will not be able to form a group even though there is a degree of skew that allows both individuals to benefit from the group. This result illustrates the essence of the inefficiency that is created by private information.

If groups do form, the share of the group reproduction going to the dominant (given by the payment function) is predicted to be $x(g_1, g_2) = 1/3g_1 + 1/3g_2 + (1 + w_s)/6$. Thus, the dominant's share of the reproduction (and hence, the reproductive skew) increases with the outside option of the dominant and decreases with the outside option of the subordinate (as $g_2 \equiv \Omega - o_2$). Among all of the groups that form, the mean payment (i.e., skew) is predicted to be $\bar{x} = (1 + w_d)/2$, which interestingly does not depend on the relatedness between individuals. Similarly, we can compute the variance of skew among all of the groups that form and find $\text{var}(x) = (1 - w_d)^2(1 + r)^2/(24(2 + r)^2)$, which is increasing in r , but decreasing in w_d .

Nondirect Mechanisms and Incomplete Control

The above analysis tells us what outcomes can be evolutionarily stable in any possible game and characterizes the equilibrium properties of the game structures that maximize the total benefit (a similar analysis can be done for any other objective function). It does not, however, answer the questions of how such an outcome might be implemented in an actual game that does not involve the imaginary arbiter or whether natural selection (or cultural evolution or learning) will actually lead to such a game structure. The latter problem is very much understudied; to our knowledge, very few studies have dealt with how the payoff

structure of a social interaction might evolve (29, 30). This section focuses on a family of simple negotiation games and shows that the optimal mechanism with uniform distributions can be implemented by a member of this family. Although our analysis does not provide a complete answer to the question of how the game might evolve, it suggests possible ways this question can be addressed.

In the context of reproductive transactions, a focal and contentious question about game structures has been who controls the division of the reproduction (16). In the so-called concession models, the dominant concedes to the subordinate the minimum reproduction that is required for the subordinate to prefer being in the group, whereas in the restraint model, the roles are reversed. Between these two extremes, both individuals would have a say, with the final division somewhere in between the two individuals' offers, termed incomplete control by both individuals (31).

Here, we present a model that combines the concession and restraint models and extends them to situations with private information. The basic informational environment remains as above with uniform distributions, but instead of focusing on direct mechanisms, we now consider the following class of games: Individuals simultaneously declare offers $\theta_i(g_i)$; θ_1 is the minimum amount of reproduction that individual 1 (the dominant) requests to assent to group formation, whereas θ_2 is the maximum amount of reproduction individual 2 (the subordinate) is willing to "pay" to the dominant. If $\theta_2 \geq \theta_1$, the group forms, and the dominant gets a share of reproduction $k\theta_2 + (1-k)\theta_1$, where k is between 0 and 1. Each value of k defines a particular game: With $k = 0$, the payment to 1 is its own offer, which corresponds to a restraint model (because the dominant is getting the minimum it requires). Similarly, $k = 1$ corresponds to a concession model, and $0 < k < 1$ to cases where neither side has complete control over the division. This setup is similar to two-person bargaining under incomplete information (32), again with the difference of nonzero relatedness that creates interdependent preferences. Note that the offer signals θ_i are not costly; hence this is a model of "cheap-talk" bargaining where individuals are free to "bluff" if they choose to.

First, we consider the equilibrium of the game for a given k . In particular, we are interested in separating equilibria, where individuals' offers θ_i are continuous and increasing functions of their private information g_i . Using the first-order conditions for the optimal offer strategies, we arrive at a set of coupled differential equations,

$$\begin{aligned} (\theta_1^{-1}(\theta_2(g_2)) - rg_2 - (1-r)\theta_2(g_2))f_2'(g_2) \\ = -(1-r)(1-k)(1-F_2(g_2))\theta_2'(g_2) \end{aligned} \quad [7]$$

$$(\theta_2^{-1}(\theta_1(g_1)) - rg_1 - (1-r)\theta_1(g_1))f_1'(g_1) = (1-r)kF_1(g_1)\theta_1'(g_1), \quad [8]$$

where $\theta_i^{-1}(\cdot)$ denotes the inverse of the offer function of individual i . To illustrate what these equations entail, we assume g_i are uniformly distributed as in the previous example. In that case, the equilibrium offer strategies θ_i are linear in g_i , with the slope and intercept being functions of w_d , k , and r (see *SI Text* for full expressions). Regardless of k and w_d , for $r < 1$, the slopes of both offer functions are less than 1, and the intercepts are nonnegative. Fig. 2 illustrates the general nature of the offer functions. Importantly, each individual i will "shade" or "mark up" its offer (i.e., bid below or above g_i , respectively), with the level and direction of this shading dependent on the relatedness.

The expected total gain in group output from cooperation is then given by

$$\frac{(1-w_d)^2(2-k(1-r)+r)(1+k+r(2-k))}{2w_s(2+r)^3} \quad [9]$$

Likewise, the expected inclusive fitness gains for a dominant and a subordinate are given by

$$\bar{W}_d = \frac{(1-w_d)^2(2-k(1-r)+r)(1+k+r(2-k))}{6w_s(2+r)^3} \quad [10]$$

$$\bar{W}_s = \frac{(1-w_d)^2(2-k(1-r)+r)(1+k+r(2-k))^2}{6w_s(2+r)^3}, \quad [11]$$

respectively. Eqs. 10 and 11 imply that the expected inclusive fitness of the dominant is decreasing in k , whereas the expected inclusive fitness of the subordinate is increasing. On the other hand, Eq. 9 implies that $k = 1/2$; i.e., a "split-the-difference" rule maximizes the total expected gain from cooperation. Such a rule corresponds to an incomplete control model where neither individual is able to impose his/her preferred division (θ_i) on the other. Moreover, equilibrium behavior under the splitting-the-difference rule implies that groups will form when

$$g_2 - g_1 > \frac{(1-r)(1+w_d)}{2(2+r)}, \quad [12]$$

which is the exact same condition as in Eq. 6. In other words, the evolutionarily stable strategies in the split-the-difference game yield the maximum gains from cooperation among all possible games. Fig. 1 depicts the region of g_1 and g_2 where groups are predicted to form and not form and the predicted share of the dominant in the groups that form. Dominants and subordinates with high outside options (higher g_1 and lower g_2 , respectively) are

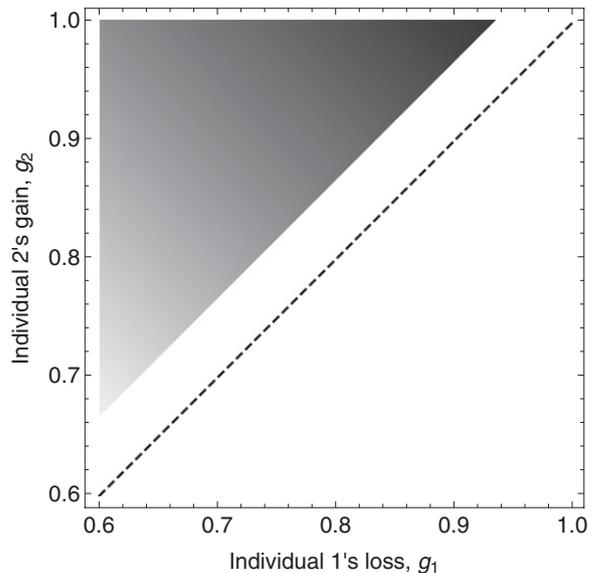


Fig. 1. The predicted outcome in the bargaining mechanism with $k = 1/2$ see as a function of g_1 and g_2 . The shaded area denotes the region in the g_1 - g_2 space where groups are predicted to form. The dashed diagonal separates the regions where groups are mutually beneficial (above the diagonal) and not (below the diagonal); the region between the diagonal and the shaded region represents groups that are mutually beneficial but cannot form at equilibrium. In the region where the group is predicted to form, darker shading indicates a higher share of reproduction for the dominant, i.e., higher skew. Parameters are $w_d = 0.6$, $w_s = 0.5$, $r = 0.25$.

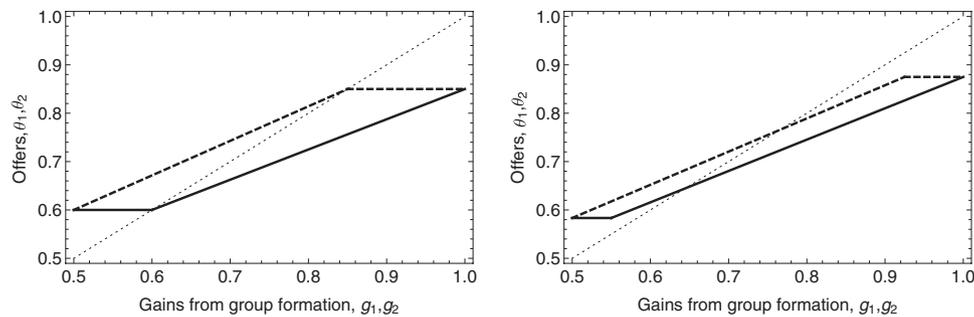


Fig. 2. The offer functions $\theta_1(g_1)$ (dashed lines) and $\theta_2(g_2)$ (solid lines). The dotted line is the 45° diagonal, corresponding to the offers being equal to the gains g_i (i.e., truthful revelation). (Left) $r = 0$; (Right) $r = 0.4$. For low values of g_2 , the subordinate cannot make an offer that the dominant will accept, and hence starts out with a constant offer strategy that is given by the offer that is accepted by the dominant with the lowest g_1 . Likewise, for high values of g_1 , no subordinate can afford to pay enough to the dominant, and hence dominants with high enough g_1 also follow a constant offer strategy. Note that when $r = 0$, no individual ever offers below its gain or loss from group formation, whereas with positive r , such offers can be part of the equilibrium strategy when either individual has high outside options (corresponding to high g_1 and low g_2), due to the inclusive fitness effect. Parameters for both panels: $w_s = w_d = 0.5$, $k = 0.6$.

predicted to get relatively higher shares of reproduction compared with the same role individuals with lower outside options. When one looks at the mean personal fitness of dominants and subordinates, however, one can see that for a large region of the parameter space, solitary individuals of either role do better than their counterparts within groups (Fig. 3). This result is due to a self-selection effect: Only individuals with relatively high outside options are predicted to remain solitary at equilibrium. Hence, even though group formation confers a net benefit to those individuals that form groups, those who choose to remain solitary are still better off on the average. This simple selection effect might explain findings where the solitary fitness of individuals is higher than the share of reproduction they are getting in the group (14, 33).

Although the split-the-difference rule maximizes the total gains from group formation, this fact does not automatically mean that this bargaining game will be the end result of evolution, even if the social system is constrained to the family of games parameterized by k . A full treatment of how evolution can change the mechanisms is beyond the scope of the current paper; however, we can note a few basic predictions from our results. In particular, Eqs. 10 and 11 imply that alleles that alter the social game to favor the subordinate's offer (i.e., higher k) will be selected for in subordinates (individuals 2) and alleles that effectively lower k will be favored in dominants (individuals 1). Where the exact balance will depend on the demography of the species, including the relatedness between individuals, which itself is a function of the skew resulting from the social game (19). These factors can lead to complex feedbacks between the evolution of the social game and the demographic properties of a species; elucidating

these feedbacks remains an open question. Another important caveat here is that Eqs. 10 and 11 assume that each individual is playing an optimal strategy given the bargaining rule k . Hence, any statement based on these fitness functions is predicated on the bargaining rule evolving much more slowly than individuals' strategies in a given social structure (or alternatively, that optimal strategies are learned behaviorally) (30). Relaxing this assumption and allowing individual's strategies to be "mismatched" to the game they are playing is likely to change evolutionary dynamics substantially.

Conclusions

For reproductive skew theory, our analyses show that the addition of private information to reproductive transactions theory can change the predictions from the models significantly. Our main result highlights a previously unrecognized constraint to the evolution of cooperation: When the distributions of individuals' gains overlap, so that there is uncertainty over whether cooperation would be mutually beneficial, it may not be possible to guarantee cooperation in all cases where it is in fact mutually beneficial. Rather than being an idiosyncrasy of a particular game setup, this result holds in any evolutionarily stable equilibrium of any game. In a sense, the possibility that cooperation might not be beneficial to both "poisons the well" and as a result, some mutually beneficial cases are forsaken. Moreover, the nature of the inefficiency is one-sided. Groups that should form do not, but groups that should not form will not mistakenly form in the optimal equilibria.

Empirically, our model draws the distinction between patterns at the between-species (represented by species means) and within-species levels (represented by distributions within populations).

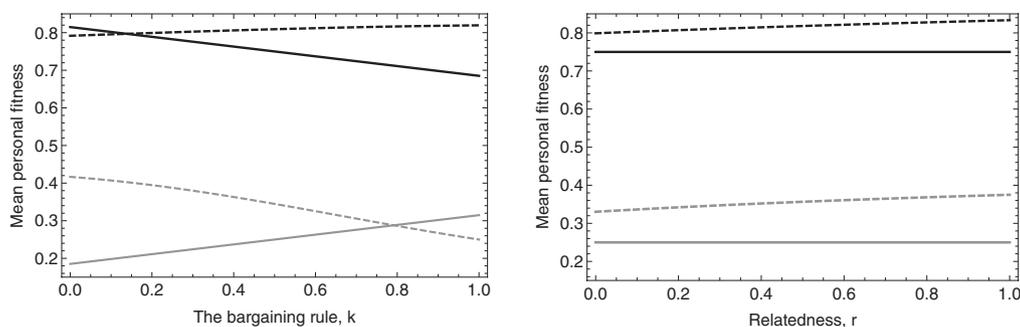


Fig. 3. The mean fitness of individuals who form a group vs. stay solitary under the cheap-talk bargaining game. The solid curves are for mean fitness within group, and dashed curves are for solitary individuals; black represents the dominants and gray the subordinates. Parameters are $w_d = w_s = 0.5$, $r = 0.25$ (Left), and $k = 0.5$ (Right).

An appreciation of this distinction might help explain why extant theory does well at explaining the former, but not the latter. However, we also see that predictions at both levels can depend on how individuals' outside options are distributed. A cautionary note here is that the estimation of these distributions can be tricky unless important selection effects are addressed. The naturally observed distribution of outside options will not represent the underlying distribution, but will be biased toward those with high outside options. This fact might resolve the apparent paradox that in some species, the within-group fitnesses of individuals appear too low to explain group stability (33). A recent resolution to this paradox has been to account for benefits later in life of individuals (34, 35); our model suggests that such benefits might not be needed to explain this pattern.

One of the prominent criticisms of reproductive skew theory is that the multitude of models of reproductive skew makes it possible to make dramatically different predictions depending on the details of the model, without an overarching theory of which model should be used in each case. In particular, pure transactional models (e.g., ref. 15) allow a range of outcomes that are consistent with group stability, with the precise level of skew being dependent on who is assumed to be in control. In contrast, compromise models (26) do predict a unique skew, but those predictions are critically sensitive to the details of the model (such as the functional shape of the reproductive share resulting from competitive effort) (14). Our approach avoids both of these issues and allows us to derive general results that are independent of the structure of the game and—given the distribution of outside options—to provide sharp predictions for the expected level of skew. Moreover, considering a specific class of games that include the concession and restraint models, we show that for uniform distribution of the outside options, the highest amount of cooperation is maintained by a mechanism where both the dominant and the subordinate have partial control over the division of reproduction. Our results thus represent a first step in elucidating the larger question of how natural selection shapes the structure of the social interaction leading to the division of reproduction.

For the evolution of social behavior and cooperation in general, our study illustrates how the methods of mechanism design can be used to study evolution of both individual behaviors and the social interaction under uncertainty and private information. Here, we studied the effects of private information about outside options; different kinds of uncertainty (such as over fighting ability) will be important in different contexts, such as parental care (9) and agonistic interactions (11). We believe that mechanism design will be a powerful tool in the evolutionary biologists' toolkit because it can be used to obtain results about evolutionary stability that are independent of particular assumptions about game structures. The general approach of mechanism design also provides a framework for studying how the social structure of a species—in addition to individuals' strategies in a given social game—evolves. Addressing this question fully will require extending existing mechanism design theory to a dynamic evolutionary setting.

Finally, our results also have some significance for economics, where our analysis can be interpreted as applying to bargaining over trades between individuals with other-regarding preferences (36, 37). Ever since Myerson and Satterthwaite (20), economists have accepted that private information leads to unavoidable inefficiencies. Our results show that other regard mitigates and, in some cases, can entirely counteract the inefficiencies in trade. One implication of this result is that in human history, groups with other-regarding agents would have an easier time taking advantage of the opportunities offered by trade. This effect would produce another route for the evolution of other regard, as such groups would have an advantage in (cultural or genetic) between-group selection. Hence, our results suggest that other-regarding preferences might have facilitated the emergence of, and coevolved with, trade and economic activity in human history.

ACKNOWLEDGMENTS. We thank seminar participants at Princeton University and University of California at Irvine for discussion and comments on the work and Christina Riehl, Jeremy Van Cleve, Dustin Rubenstein, and two reviewers for valuable feedback on the manuscript. This study was supported by National Science Foundation Grant EF-1137894.

- Trivers RL (1971) The evolution of reciprocal altruism. *Q Rev Biol* 46:35–57.
- Axelrod R, Hamilton WD (1981) The evolution of cooperation. *Science* 211:1390–1396.
- Sachs JL, Mueller UG, Wilcox TP, Bull JJ (2004) The evolution of cooperation. *Q Rev Biol* 79:135–160.
- Lehmann L, Keller L (2006) The evolution of cooperation and altruism—a general framework and a classification of models. *J Evol Biol* 19:1365–1376.
- Levin SA (2009) *Games, Groups, and the Global Good (Springer Series in Game Theory)*, ed Levin SA (Springer, Berlin), pp 143–153.
- Levin SA (2010) Crossing scales, crossing disciplines: Collective motion and collective action in the Global Commons. *Philos Trans R Soc B* 365:13–18.
- Zahavi A (1975) Mate selection: A selection for a handicap. *J Theor Biol* 53:205–214.
- Grafen A (1990) Biological signals as handicaps. *J Theor Biol* 144:517–546.
- Godfray HCJ (1991) Signalling of need by offspring to their parents. *Nature* 352:328–330.
- McNamara JM, Gasson C, Houston AI (1999) Incorporating rules for responding into evolutionary games. *Nature* 401:368–371.
- Parker GA, Rubenstein DI (1981) Role assessment, reserve strategy and acquisition of information in asymmetric animal conflicts. *Anim Behav* 29:221–240.
- McCarty NA, Meirowitz A (2007) *Political Game Theory: An Introduction* (Cambridge Univ Press, Cambridge, UK).
- Vehrencamp SL (1983) Optimal degree of skew in cooperative societies. *Am Zool* 23:327–335.
- Nonacs P, Hager R (2011) The past, present and future of reproductive skew theory and experiments. *Biol Rev Camb Philos Soc* 86:271–298.
- Reeve HK, Ratnieks FLW (1993) *Queen Number and Sociality in Insects*, ed Keller L (Oxford Univ Press, Oxford), pp 45–85.
- Johnstone RA (2000) Models of reproductive skew: A review and synthesis. *Ethology* 106:5–26.
- Kokko H (2003) Are reproductive skew models evolutionarily stable? *Proc R Soc Lond B Biol Sci* 270:265–270.
- Frank S (1998) *Foundations of Social Evolution* (Princeton Univ Press, Princeton).
- Johnstone RA (2008) Kin selection, local competition, and reproductive skew. *Evolution* 62:2592–2599.
- Myerson R, Satterthwaite M (1983) Efficient mechanisms for bilateral trading. *J Econ Theory* 29:265–281.
- Myerson R (1979) Incentive compatibility and the bargaining problem. *Econometrica* 47:61–74.
- Queller DC (1992) Quantitative genetics, inclusive fitness, and group selection. *Am Nat* 139:540–558.
- Nöldeke G, Samuelson L (1999) How costly is the honest signaling of need? *J Theor Biol* 197:527–539.
- Roughgarden J, Song Z (2013) *Human Nature, Early Experience and the Environment of Evolutionary Adaptedness*, eds Narvaez D, Panksepp J, Schore A, Gleason T (Oxford Univ Press, New York).
- Akçay E (2012) Incentives in the family ii: Behavioral dynamics and the evolution of non-costly signaling. *J Theor Biol* 294:9–18.
- Reeve HK, Emlen ST, Keller L (1998) Reproductive sharing in animal societies: Reproductive incentives or incomplete control by dominant breeders? *Behav Ecol* 9:267–278.
- Spence M (1973) Job market signaling. *Q J Econ* 87:355–374.
- Bergstrom CT, Lachmann M (1998) Signaling among relatives. III. Talk is cheap. *Proc Natl Acad Sci USA* 95:5100–5105.
- Worden L, Levin SA (2007) Evolutionary escape from the prisoner's dilemma. *J Theor Biol* 245:411–422.
- Akçay E, Roughgarden J (2011) The evolution of payoff matrices: Providing incentives to cooperate. *Proc Biol Sci* 278:2198–2206.
- Clutton-Brock TH (1998) Reproductive skew, concessions and limited control. *Trends Ecol Evol* 13:288–292.
- Chatterjee K, Samuelson W (1983) Bargaining under incomplete information. *Oper Res* 31:835–851.
- Nonacs P, Liebert A, Starks P (2006) Transactional skew and assured fitness return models fail to predict patterns of cooperation in wasps. *Am Nat* 167:467–480.
- Field J, Cronin A, Bridge C (2006) Future fitness and helping in social queues. *Nature* 441:214–217.
- Leadbeater E, Carruthers J, Green J, Rosser N, Field J (2011) Nest inheritance is the missing source of direct fitness in a primitively eusocial insect. *Science* 333:874–876.
- Fehr E, Gächter S (2000) Fairness and retaliation: The economics of reciprocity. *J Econ Perspect* 14(3):159–181.
- Sobel J (2005) Interdependent preferences and reciprocity. *J Econ Lit* 43:392–436.