

# 2

## Explaining Variation in Cooperative Behavior

### Perspectives from the Economics Literature

Friederike Mengel and Joël van der Weele

#### Abstract

Casual observation and controlled experiments show that humans display great heterogeneity in their tendency to exploit others or invest in mutual cooperation. This chapter reviews models in the economics literature that can explain the coexistence of free riders (exploiters) and cooperators (investors). A distinction is made between models of full and bounded rationality. Although some models provide tentative explanations, there is a large gap between the empirical and theoretical literature, and there has been little effort to integrate long- and short-run models.

#### Introduction

Human societies thrive when individuals invest in cooperative relationships and common interest. At the same time, however, this situation offers opportunities for individuals to exploit the investments of others. Economists have studied the choice between investment and exploitation in the context of so-called social dilemma games, which juxtapose two kinds of actions. Individuals can “cooperate” with others by investing individual resources in actions that benefit the group. Alternatively, they can exploit, “defect,” or “free ride” by choosing a strategy that benefits the individual but implies a material cost to the group.<sup>1</sup> This dichotomy offers a fruitful way to model many social and economic interactions that humans engage in on a regular basis.

---

<sup>1</sup> For the remainder of this chapter, we will use the terms “cooperation” and “defection” strategies in the dilemma situation that is the focus of the economics literature, and that fall within the broader concepts of “investment” and “exploitation.” We will use the words “free ride,” “defect,” and “cheat” interchangeably.

Casual observation and experimental evidence show that there is great heterogeneity in people's tendency to cooperate or defect in social dilemma games (Ledyard 1995; Fischbacher et al. 2001; Camerer 2003). Explaining the origins and the stability of such heterogeneity is therefore an important task for social scientists. In this overview, we discuss whether economic models can help explain the origins and stability of heterogeneity. To do so, we distinguish between heterogeneity on different levels:

1. Long-term processes of natural selection, cultural transmission, and learning generate diversity in people's capabilities, knowledge, and preferences.
2. Institutional environments determine how such heterogeneity is expressed in behavior, depending on the incentives that such institutions provide.

Methodologically, both types of models are examples of game theoretic analysis, but the underlying assumptions are quite different. Long-run models rely on the assumption that strategies spread through social learning, cultural transmission, or the creation of offspring. These processes are often boundedly rational in nature. By contrast, short-run models typically assume rational or optimizing responses to the incentives provided by different institutions. This distinction is reflected in our discussion.

We begin with a discussion of how, for a given distribution of preferences, institutions that punish defection can induce heterogeneity in behavior. Thereafter we turn to models of bounded rationality and discuss how learning, imitation, and evolution may contribute to equilibria with heterogeneity in preferences or strategies. Throughout, we focus on stylized social dilemma situations, whereas Oldekop and Hajjar (this volume) focus on the importance of contextual factors.

The distinction between long-run processes of preference formation and short-run reactions to incentives is not always clear-cut. Institutions that incentivize certain behaviors in the short term can also cause long-term changes in preferences and behavior (Bowles 1998). Conversely, long-term movements in preferences will change the kind of institutions that will be necessary to promote cooperation. Where available, we discuss a literature that addresses this two-way relationship and argue for the need for more comprehensive theories.

### **Punishment Institutions**

Social dilemma games (e.g., prisoner's dilemma, public goods game) exemplify the trade-off between exploitation and investment that is the focus of this volume. In such dilemmas, individuals choose between a personally costly "cooperative" action that benefits other group members and "free riding" or

“defection,” which maximizes the material payoffs to the individual at the expense of the group.

All individuals in a group are better off when they can achieve cooperation from all members, compared to a situation where everyone defects. To limit defection, societies have developed a diverse set of institutions to punish defectors. However, punishment may not deter all defections, thus allowing both cooperative and free-riding behavior to coexist. Here we review evidence for the coexistence of both behaviors in institutional contexts that have been analyzed in the economic literature. While we focus on the effect that institutions have on the behavior of optimizing agents, we also highlight how institutions could affect long-term processes of preference formation.

### **Centralized Punishment**

Punishment by a central authority has been a major topic in the philosophical literature since at least Thomas Hobbes. In the seventeenth century, Hobbes argued how a government or Leviathan could improve the situation of noncooperating individuals in a “state of nature” and be sustained as part of a social contract between individuals. An economic take on this idea is the economic model of crime, formalized by Becker (1968). The starting point in this model is a society of agents who have the opportunity to take an action (crime) that brings personal benefits but hurts others in the society. Typically individuals will differ in the personal cost and benefits of crime, due to differences in moral convictions, wealth, and other personal circumstances. Potential criminals rationally weigh the costs and benefits of the crime. A central authority can influence the calculations of these individuals and improve efficiency by raising the cost of crime, for example, through sanctions or imprisonment. Thus, the model of crime extends “the economist’s usual analysis of choice” (Becker 1968:170), by analyzing crime as a good and punishments as that good’s price.

Becker showed that it is not optimal to deter all crime when deterrence is costly and the costs and benefits of crime differ between individuals. Rather, crime should be deterred only up to the point where the marginal costs of enforcement plus the marginal benefits to the criminal equal the marginal cost of crime to the victim. A pragmatic, optimizing authority will allow some crime to occur, either because it is too costly to eradicate or because it is relatively harmless to the victim, or both. Thus, the optimal policy that flows from this model results in the coexistence of both compliant and criminal behavior.

There are many subtle forms of exploitation or free-riding behavior that are technically not “crimes” punishable by law. Becker’s model is very general and can be applied to these subtle forms of defection. Similarly, there is flexibility in the incentives and the authorities that can be considered. For example, the management of a firm can discourage shirking behavior by paying part of a salary in bonuses for high performance, or a football coach can bench players who perform poorly.

The central implication of this theory is that crime should fall with both higher penalties and a higher probability of getting caught. There is a decade-old debate about the empirical validity of this “deterrence hypothesis.” Although there is some evidence that the probability of getting caught matters, there is no consensus as to whether higher penalties reliably deter crime. The effects of deterrent policies appear to be highly dependent on the social context (van der Weele 2012a).

In the long term, centralized punishment can also affect the coexistence of cooperative and selfish “types” who differ in their preferences for cooperation. Consider, for example, interactions involving bilateral exchange, where one party can cheat the other party (Huck 1998). Trading partners cannot distinguish between cooperative and defector types, so the latter will take advantage to cheat on their contracts. In this situation, probabilistic detection of cheating associated with penalties favors cooperative types who are more likely to comply. As a consequence, cooperative types become more prevalent in the population, causing the optimal size of sanctions to decrease over the long term.

By contrast, when the type of the interaction party is observable at least with some probability, the imposition of sanctions can favor defectors. In the absence of sanctions, cooperative types would never interact with defectors, who would suffer low payoffs and make up only a small share of the population. Bohnet et al. (2001) showed that the presence of enforcement weakens this kind of ostracism and leads to an increase in the population share of defectors, as long as penalties are relatively low. Thus, cooperation or trustworthiness is crowded out with weak enforcement, but crowded in with strong enforcement. Indeed, in their experiment, Bohnet et al. (2001) found that weak penalties lead to an increased prevalence of cheating.

In summary, coexistence of cooperators and free riders exists naturally in a world with costly enforcement and heterogeneity in the costs and benefits of crime. The long-run effects of centralized enforcement depend on the size of the penalties and the available information about the interaction partner.

### **Decentralized Punishment and Social Norms**

Many forms of punishments are not carried out by a central authority, but rather by a community of peers.<sup>2</sup> Peers have an important advantage over a central authority in that they will often have more information about the nature of transgressions and their perpetrators. Peer punishment can take the form of

---

<sup>2</sup> The conceptual distinction between centralized and peer-sanctioning schemes does not mean that the two are independent of each other. For example, Huck and Kosfeld (2007) argue that when an authority reduces sanctions to discourage crime, this may also lead to the abandonment of neighborhood watch groups engaged in peer surveillance. Van der Weele (2012b) shows that when a government has superior information about the number of free riders in society, it may refrain from setting severe penalties to avoid signaling to citizens that being a free rider is the social norm.

tacit or open disapproval, withholding of cooperation, ostracism, or even physical attack against the perpetrator. Peer punishment, however, is often costly to the punisher, and the enforcement of cooperative social norms becomes itself a public good, often referred to as a “second-order public good problem.”

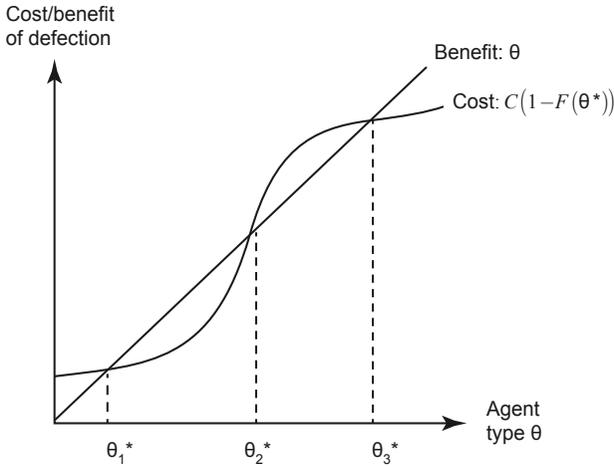
The second-order public good problem is easier to solve than the first-order problem, as a large part of the population seems willing to punish transgressions at a potential cost to themselves. Fehr and Gaechter (2000), as well as many follow-up studies, demonstrated this experimentally in the context of a public goods game, a multiperson version of the famous prisoner’s dilemma. The authors augment this game with a punishment phase, in which participants can take away earnings of noncontributors at a cost to their own experimental earnings. Many participants do indeed choose to punish which, despite the initial destruction of resources, leads to higher cooperation rates and higher efficiency over time (Gaechter et al. 2008).

The cost and effectiveness of peer punishment is likely to depend on the ratio of defectors and cooperators. When cooperators are relatively numerous, punishment resources can be concentrated on a smaller sample of defectors, leading to higher (probability of) punishment. At the same time, an increased expectation of punishment is likely to lead to a higher cooperation rate.

We show the consequences of such punishment complementarities for the coexistence of free-riding and cooperative behavior in a canonical model. Variations of this model have been applied by different authors to different instances of defection such as tax evasion (Lindbeck et al. 1999; Traxler 2010). Suppose that there is a community of a countable infinite number of agents indexed  $i = 1, 2, \dots$ . Each agent can choose to behave cooperatively or to defect. Payoffs from cooperation are zero; payoffs of defection are given by  $\theta_i - C(n)$ , where  $\theta_i$  is the “type” of agent  $i$  that encapsulates all psychological and material payoffs from defection specific to that agent. Agents differ with respect to their  $\theta_i$  and are distributed over the type space according to some distribution  $F(\theta)$  with full support on  $[0, \theta_{\max}]$ . The second term,  $C(n)$ , represents the social cost of punishment to the defector. This cost depends negatively on the fraction  $n$  of the population that defects, reflecting the assumption that the more defectors there are, the more punishment will be diluted.<sup>3</sup>

A rational individual thus defects if and only if  $\theta_i \geq C(n)$ . If we define by  $\theta^*$  the type that is just indifferent between defecting or not, then all types  $\theta > \theta^*$  will defect. Since  $n = 1 - F(\theta^*)$  is the fraction of defectors, it follows that  $\theta^*$  is defined implicitly by the equation  $\theta^* = C(1 - F(\theta^*))$ . Under suitable conditions on the functions  $C(n)$  and  $F(\theta)$ , this model yields a situation with multiple equilibria, as illustrated in Figure 2.1.

<sup>3</sup> To keep things simple, we have assumed here that the payoffs of defection (and the effectiveness of punishment) do not depend on  $n$ . If they do, one can get similar qualitative results, as long as these payoffs of defection fall less quickly with  $n$  than its costs.



**Figure 2.1** Multiple equilibria in a model with heterogeneous benefits of defection. The x-axis shows the agent type  $\theta$ ; the y-axis shows the costs  $C(1 - F(\theta^*))$  and benefits ( $\theta$ ) from defection. Equilibrium is found where the two lines intersect, and costs and benefits are equal.

Equilibria characterized by high levels of defection and low levels of peer enforcement (e.g.,  $\theta_1^*$  in Figure 2.1) coexist with equilibria featuring low levels of defection and high levels of enforcement (e.g.,  $\theta_3^*$  in Figure 2.1). Depending on the shape of the functions  $C(n)$  and  $F(\theta)$ , there may be a large number of such equilibria. Note that the equilibrium associated with  $\theta_2^*$  in Figure 2.1 is unstable. That is, suppose that type  $\theta_2^*$  deviates from this equilibrium and chooses to cooperate. This raises the cost of defections for the remaining agents, making it optimal for some types  $\theta > \theta^*$  to cooperate as well, and causes the equilibrium to unravel. Conversely, if some types just below  $\theta^*$  defect, it becomes optimal for “lower” types to do so as well. Thus,  $\theta_2^*$  is better thought of as a “tipping point” rather than an equilibrium, on either side of which a different stable equilibrium becomes an attractor.

Common terminology associates the (equilibrium) level of  $n$  with a “social norm,” because it measures the degree to which cooperation and defection are normal actions in the population. Glaeser et al. (1996) argued that variation in social norms can explain the lion’s share of the empirical variation in crime. In their intercountry comparison, using data from 1980, the homicide rate in the United States was about 150 times higher than in Japan. Looking at differences between U.S. cities, Glaeser et al. found that the crime rate in Atlantic City, New Jersey, was about 400 times higher than the nearby city of Ridgewood Village. On an even more detailed intra-city level, the crime rate in the 1st Precinct of New York City was about 10 times greater than in the 123rd Precinct. Examining the crime statistics in New York in more detail, the authors show that at most 30% of these differences can be accounted for

by observable differences between different locations (e.g., levels of income, schooling, female-headed households, arrest rates). The rest, they argue, is due to peer effects or social norms.

In summary, complementarities in the effectiveness of peer punishment imply the existence of multiple equilibria or social norms. One can think of these multiple equilibria as “parallel societies” that may exist in different places or at different times. When there is multiplicity, theory cannot predict *ex ante* how a given distribution of preferences translates into cooperative behavior, but each equilibrium may itself involve a stable coexistence of cooperation and free-riding behavior.

### Exclusion and Sorting

A particular form of peer punishment is ostracism or exclusion. Forming and terminating relationships can be part of the strategy set, but separation can also arise as an equilibrium phenomenon even if sorting operates via other mechanisms. Our focus in this section is on the latter type of models; an example of the former is taken up later.

Kosfeld et al. (2009) derived a good example of a mechanism that relies on sorting through their analysis and experimentally testing a coalition formation model within the context of public good provision. They modeled institution formation as a three-stage game. In the first stage, each player decides whether to participate in an organization that, once implemented, exerts punishment on individual members who do not contribute their full endowment to the public good. The organization is costly, and only players who are members of the organization can be punished. In the second stage, players learn how many of the other players are willing to participate. The organization is implemented if and only if all players willing to participate agree to its actual formation. In the final stage, the public goods game is played. Theoretically, two types of equilibria can be sustained: one in which at least a minimal amount of players establish an institution and contribute to the public good; the other where no institution is established. Under the first type of equilibria, coexistence can be established.

Models of coalition formation have also been used to explain cartel formation (Green and Porter 1984). The incentives in these cases are very similar, but most would view the successful establishment of a collusive institution as less beneficial, because of the harm inflicted on third parties, often consumers.

Given the intuitive appeal and wide range of applications of such models, it is not surprising that these ideas have been tested in a variety of experiments. Kosfeld et al. (2009) tested their mechanism in the lab and found that participants are unwilling to support institutions where some individuals can free ride. The vast majority of outcomes found in the lab feature institutions where everyone is part of the institution and contributing, or where no one is. Guererck et al. (2006) found evidence of the effectiveness of sorting in a slightly different

setting. In their experiments, participants endogenously choose to select into a punishment mechanism. The main difference to Kosfeld et al. (2009) is that those who select into the punishment mechanism do not interact in the public goods game with those who did not select into the mechanism. While sorting into the punishment institutions can sustain cooperation in the lab in both studies, it is full cooperation, not coexistence, that is sustained.

In summary, some models have been able to establish coexistence via sorting or exclusion mechanisms. Notably, however, when tested in the lab, coexistence has not been established. Whether these lab results (typically obtained in small groups of 3–5 participants) extend to much larger groups remains to be seen.

### **Reputation Motives**

In our discussion of punishment, we have implicitly assumed that defectors can be identified, at least with some probability. However, the information that is available about a person's past behavior depends itself on institutions. For example, one of the crucial obstacles facing the development of online commerce platforms, such as eBay, is the development of reputation formation mechanisms that allow customers to identify fraudulent sellers.

Researchers have modeled the effect of reputational motives on cooperation in various ways. Kreps et al. (1982), for example, investigated the power of reputational motives for cooperation in the finitely repeated prisoner's dilemma. Under common knowledge that all agents are selfish, theory predicts that cooperation will unravel: since cooperation is an irrational move in the last round of play, promises based on cooperation in later rounds have no credibility. Kreps et al. (1982), however, posit that players believe that, with some probability, their opponent may be a "tit-for-tat" player (cooperate on first move and then match the strategy opponent used on last play). If such beliefs are sufficiently high, they show that it may be rational to pretend to be a tit-for-tat type, at least until the end of the game draws near, so as to exploit the other player's conditional willingness to cooperate. As a consequence, all players may act cooperatively, even though none were of a tit-for-tat type.

Economic experiments that test this reasoning show that there is substantial heterogeneity in cooperative behavior. For example, Andreoni and Miller (1993) conducted a prisoner's dilemma game where the same partners interacted repeatedly for a finite number of rounds. They observed a cooperation rate of about 60% in early rounds, which deteriorated toward the end of the game. By contrast, when reputation formation was not possible, the cooperation rate was much lower at about 20%. In a similar experiment carried out by Bolton et al. (2004), sellers and buyers interacted repeatedly, and sellers had the possibility to defraud the buyers. Bolton et al. included a feedback condition: partners rotated but buyers could observe the behavior of the buyer in previous rounds. Although the conditions where participants played with the

same partner in each round yielded almost full (about 90%) trustworthiness from sellers, trust in the feedback conditions hovered around 70% until the last few rounds of the experiment.

The model by Kreps et al. (1982) does not explain why some people behave like cooperative partners while others do not, but it provides a starting point. One plausible explanation is that participants have different initial beliefs about the likelihood of facing a cooperator, although it is unclear exactly where these beliefs originate, as the experimental conditions were the same for all. Since the theory allows for both defection and cooperation as an equilibrium, it may also be that different experimental subjects play different equilibria, without managing to converge on a single equilibrium. Finally, some forms of learning may rationalize these results (see next section).

In the long run, reputation motives can lead to the evolution of strategies of “indirect reciprocity.” In the image scoring game introduced by Nowak and Sigmund (1998), two players are randomly drawn in each round from a population: one player is randomly selected to be the donor and the other the receiver. The donor can decide to “keep,” yielding a payoff of  $c$  to the donor and 0 to the receiver, or the donor can “give,” yielding a payoff of  $b > c$  to the receiver and 0 to the donor. The donor’s decisions result in an “image score” that is visible to the partner: it is 1 if the donor gave at the last opportunity, and 0 otherwise.

Nowak and Sigmund (1998) considered the strategies employed by universal defectors (who never give), universal altruists (who always give), and “discriminators” (who only give to a partner with a sufficiently high image score). Such discriminators, in fact, practice a form of indirect reciprocity; that is, they cooperate in the hope that the resulting image will induce cooperation from a future interaction partner. Nowak and Sigmund demonstrated that in this model, discriminators can successfully invade a population of universal defectors when the likelihood of knowing the partner’s image score  $q$  exceeds the ratio  $c/b$ . However, universal altruists can invade discriminators, who will occasionally resort to punishment at a cost to themselves. Altruists can be invaded, in turn, by defectors, causing a never-ending evolutionary cycle.<sup>4</sup>

In summary, in finitely repeated dilemma games with reputation formation, a mixture of cooperation and defection can be observed. Economic theory suggests that this is due to differences in beliefs about the behavior of other agents, but it does not explain the origins of such beliefs. Evolutionary models link reputation formation to indirect reciprocity, but these models are hard to test. Further work is thus needed to determine how the increased importance of reputations, in a world with social media and an increasing amount of available information, impacts cooperative behavior.

<sup>4</sup> Note, however, that Lotem et al. (1999) showed that if there is a steady exogenous supply of defectors, the discriminator’s advantage over altruists remains, and the frequency of discriminators stabilizes.

## **Bounded Rationality**

In this section, we focus on explanations under which boundedly rational agents learn to adopt certain behaviors or social norms via processes of social learning or cultural transmission. Boundedly rational agents use heuristics or learning rules which may ultimately make them fail to appreciate possible private gains. Thus, under this class of explanations, punishment is not required to sustain cooperation and/or exploitation.

### **Learning**

Models of learning specify intuitive updating and choice rules under which agents do not always rationally incorporate all available information. Learning rules differ widely, according to their degree of sophistication, and range from very simple reinforcement type rules to forward-looking and optimizing agents. Here we do not attempt to give a complete overview of the learning literature (for a good overview see Dhimi 2016, Chapter V). Instead we focus on select models that can explain the coexistence of cooperators and free riders.

#### *Reinforcement, Endogenous Aspirations, Forward-Looking Learners*

In behavioral psychology, reinforcement increases the frequency of a certain behavior whenever that behavior is followed by a stimulus that is appetitive or rewarding. In economics, reinforcement learning refers to the fact that an agent's propensity to choose any given action will be proportional to past payoffs obtained with that action (Roth and Erev 1995). Payoffs are thus the economic agents' antecedent stimulus and motivator in these models. Standard (Erev–Roth type) reinforcement learning models are approximated by the evolutionary replicator dynamics and can only support Nash equilibria (and hence defection in the prisoner's dilemma) as stable states (Börgers and Sarin 1997). Since Nash equilibrium is a coarsening of evolutionary stability (Maynard-Smith and Price 1973), this also means that cooperative outcomes are not evolutionarily stable.

Some variants of reinforcement learning with endogenous aspirations—where agents are satisfied when their aspiration is met and aspirations adjust to recent payoff experiences (Simon 1956)—have been shown to support cooperation (Karandikar et al. 1998). In the following scenario, studied by Karandikar et al. (1998), players are randomly matched to play a  $2 \times 2$  prisoner's dilemma game. Each player has an aspiration at each date and takes an action. The action is switched at the subsequent period only if the achieved payoff falls below the aspiration level, with a probability that depends on the shortfall. Aspirations are updated in each period, depending on the divergence of achieved payoffs from aspirations in the previous period. Karandikar et al. (1998) showed that if the speed of updating aspiration levels is sufficiently

slow, then the outcome, in the long run, must involve both players cooperating most of the time. While there is no coexistence of cooperators and defectors in a stable state, players may (and occasionally do) profit by deviating from cooperative behavior. The dynamics of the process, however, ultimately leads back to mutual cooperation.

In another class of models, where it is assumed that agents are forward looking (anticipating future path of play) but still adaptive (learning from past experience), learners are much more sophisticated. With forward-looking agents, cooperation can be sustained in finitely repeated interactions because agents learn that histories involving defection are more often followed by defection than histories involving cooperation. Forward-looking agents learn that defection can sour relationships. Such forms of learning can explain why there is more cooperation when reputation formation is possible (Andreoni and Miller 1993; see also above discussion on reputation motives).

### *Imitation, Networks, and Exclusion*

Let us now turn to social learning models, where agents copy behaviors observed in others, and revisit exclusion as a mechanism to sustain cooperation. This time, though, sorting is not a choice by optimizing agents; it arises endogenously via a process of imitation learning by boundedly rational agents.

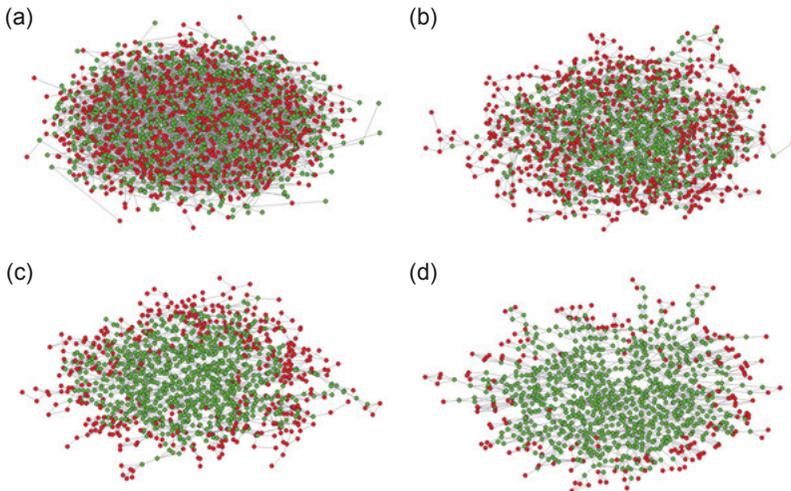
Similar in spirit to the group/kin selection literature in biology, a number of models have been proposed that rely, to some extent, on excluding defectors from beneficial interactions with cooperators. These models exploit the idea that if, for whichever reason, cooperators were to interact more frequently with other cooperators (and defectors with defectors), then cooperators could achieve higher evolutionary fitness, precisely because joint cooperation Pareto-dominates joint defection; that is, all players are better off under joint cooperation compared to joint defection. This idea has appeal because it is grounded in evolutionary considerations and does not rely on the introduction of “types” (or tweaks to payoffs) that would not survive the test of evolutionary fitness.

Imitation learning, where agents imitate successful behavior of others, has been able to produce stable states with coexistence of cooperators and defectors. In a seminal paper, Eshel et al. (1998) modeled agents located on a circle network imitating the successful actions of their neighbors. Agents in this model copied an action if and only if it yielded higher average payoffs in their neighborhood. This is very different from blind copying, customarily used in the behavioral ecology literature (see also Dubois et al., this volume). Under the assumptions made in Eshel et al. (1998) there can be clusters of cooperators coexisting with clusters of defectors in the network. Cooperators in the middle of a cluster or defectors in the middle of a cluster of defectors will not change their behavior via imitation learning simply because everyone they observe (and hence could imitate) chooses the same action as themselves. What about the cooperators and defectors at the fringes of these clusters? The defectors are

getting a better payoff than the cooperators they observe (who they exploit) and thus will not be tempted into imitating. The cooperators also will not switch because, in these equilibria, they observe other cooperators (those in the center of a cluster) who are very well off. These equilibria, however, have proven to be very fragile. In subsequent literature it has been shown that they are not obtained in other networks and that allowing agents to use information from agents further away (even if only second-order neighbors) destroys the result even in the circle (Goyal 2007; Mengel 2009).

In a model of endogenous network formation, the possibility of coexistence of free riders and cooperators has been demonstrated by Fosco and Mengel (2011). With endogenous networks, people imitate actions *and* link choices of successful others; this leads to a coevolution of the network with choices in the prisoner's dilemma. In absorbing states of this coevolutionary process, the shortest path between any two cooperators never involves a defector. The reason is that any two cooperators separated by a defector will want to sever ties with the defector and establish instead mutual links. This leads networks to form dynamically, as illustrated in Figure 2.2: cooperators end up occupying central positions in the network, with free riders in more peripheral positions (Figure 2d).

Free riders will not imitate cooperation, since the only cooperators they observe are linked to many defectors and make poor payoffs. Most cooperators do not observe any defectors. Those that do observe defectors are by and large linked to other defectors, thus making lower payoffs than the cooperators



**Figure 2.2** Coexistence of cooperators (green nodes) and free riders (red nodes) in a model of endogenous network formation. Initially (a) cooperators and free riders are randomly allocated on the network. (b)–(d) As agents start to form and sever links, free riders are pushed increasingly toward the periphery of the social network.

they observe in the periphery. Here a key element of bounded rationality is that these “bridging cooperators” (who are linked to both defectors and cooperators) assess the benefits of cooperation by comparing themselves to people in very dissimilar situations.

There is no clear empirical evidence on whether defectors or cooperators tend to be more central in social networks, though some studies have found that altruists (measured by giving in the dictator game) tend to be more central (Branas-Garza et al. 2010). In a study of smoking behavior of 12,067 people assessed between 1971 and 2003, Christakis and Fowler (2008) found that smokers moved increasingly to the periphery of their social networks. They interpret this to reflect a societal change in perception, where over time, smoking came to be seen as antisocial. These models may also provide some insight into the interactions of hunter-gatherer groups, which were mostly bilateral in kinship, as opposed to the lineage- or village-based grouping of horticulturalists and agriculturalists. Conversely, allowing for kinship structure in some of these models might generate new insights.

In summary, while various classes of individual learning models can only sustain Nash equilibria (and thus defection), models using endogenous aspirations as well as models that rely on limited forward-looking players can sustain some degree of cooperation. None of these, however, has been shown to sustain stable coexistence. By contrast, some social learning models, especially those which rely on imitation, have produced “proper coexistence” of cooperators and defectors who interact with each other in a social network.

### **Categorization**

An important aspect, which economic theories of learning have only started to incorporate recently, is that agents’ experiences and learning are not grounded in a single game alone: they develop across a great variety of situations. Theories of categorization and learning across games model how agents’ learning and behavior in one situation will be affected by others. Categorization occurs to economize on reasoning cost or to make faster decisions. Allport (1954) famously noted that “the human mind must think with the aid of categories. We cannot possibly avoid this.” In the prisoner’s dilemma, agents may cooperate because to do so seems optimal *on average* across a broad category of different situations.

There is only a limited literature available on categorization and cooperation in economics. However, recent advances have been made in understanding when categorization may occur and in showcasing examples of categorization affecting behavior in situations that resemble social dilemmas.

Noting that humans compete in “Machiavellian tournaments,” Samuelson (2001) modeled situations where agents put increasing amounts of effort into the tournament and lump together other decision situations in coarse categories as a consequence. Using this approach he showed that fair splits in ultimatum

games can be sustained in equilibrium as players will lump ultimatum games together with (longer) alternating bargaining games, where fair splits are equilibrium outcomes. More generally, bunching can be advantageous in an evolutionary sense whenever reasoning or other costs make it prohibitive to devise strategies for each game separately. In such settings cooperation can be sustained in games where the incentives to defect are not “too big,” because these games are bunched together in evolutionary equilibrium with others where cooperation is a Nash equilibrium.

Another strand of literature has focused on coarse beliefs. Jehiel (2005) developed a concept called analogy-based expectations equilibrium wherein people form the same beliefs across different situations bunched together in the same equivalence class of games. Irrespective of how beliefs are formed, there are no beliefs that can rationalize choosing a dominated strategy, such as cooperation in the prisoner’s dilemma. Early experimental evidence supports this point, demonstrating how categorization can affect equilibrium play in a range of games, with the exception of those that have dominant strategies (Grimm and Mengel 2012).

This relatively young research area could be a promising avenue for further research, possibly studied in conjunction with identity concerns. People interact in a great variety of situations and how we partition and view the world is part of our culture and social identity. Existing studies have shown that taking these factors into account can rationalize seemingly irrational behavior in experimental games that are closely related to social dilemmas (Jehiel 2005; Mengel 2012). Exploring more deeply such identity-driven motives in conjunction with coarse (culturally determined) reasoning could lead to novel results on coexistence.

### **Cultural Transmission and Evolution of Preferences**

Culture is transmitted within and across generations. The cultural transmission literature tries to model this process, often focusing on the subtle interactions between genetic and cultural evolution thought to occur at differential speeds. Models of cultural transmission distinguish between three types of transmission (Cavalli-Sforza and Feldman 1981; Boyd and Richerson 2005): (a) horizontal transmission, including the formation of social norms and peer pressure, (b) vertical transmission, including parent’s socialization efforts, and (c) oblique transmission through, for example, media, teachers, or other role models.

One of the most interesting approaches in this area focuses on the endogeneity of preferences and stems from the work of Bisin et al. (2004; see also Bowles 1998). Building on earlier work by Cavalli-Sforza and Feldman (1981), Bisin et al. focused on vertical transmission, assuming that parents who socialize their children are motivated by altruism (i.e., a non-selfish concern for their children’s well-being). However, altruism is not perfect in the sense that parents evaluate their children’s payoffs with their own preferences (rather

than their children's). Parents invest some effort into socializing their children: if they succeed, the child will adopt the parents' preferences; if they fail, horizontal transmission will kick in and the child will become socialized according to the preferences of the majority of the child's peers. Altruistic preferences survive as a minority preference, because minorities have higher incentives to socialize their offspring to their own preferences than majorities do. Thus, this is one of the few true coexistence results in the literature. Empirical evidence on how altruistic behavior is affected by the minority or majority status of a group remains scarce.

In models of horizontal transmission, conditionally cooperative preferences have been shown to coevolve with matching structure. Here, unlike in the models discussed earlier (see discussion on exclusion and sorting), the mechanism is not so much one of exclusion, where the fact that some agents are excluded from beneficial interactions with cooperators leads to coexistence. Instead, the strength of moral concerns itself changes with the matching structure (Mengel 2008). The idea is that people feel guilty about free riding, for example avoiding taxes, if these moral concerns are shared by many others in their social environment. If, by contrast, most people do not feel guilty about free riding, then feelings of guilt will be subdued (see also above discussion on exclusion and sorting). The conditionality of the moral concerns gives rise to conditionally cooperative behavior (Fischbacher et al. 2001). Since conditional cooperators adapt their behavior with the share of free riders in their environment, they cannot be easily exploited. This underlies the evolutionary fitness of cooperative behavior. While the previous mechanism works without observing other's preferences prior to interaction, a simpler mechanism that does require observability of preferences was noted by Bester and Gueth (1998). They note that if preferences can be observed, then conditional cooperation can be evolutionarily stable whenever there is enough "strategic complementarity" in decision making. Conditional cooperators cooperate whenever they observe others with cooperative preferences; otherwise, they free ride.

To summarize, mechanisms of cultural transmission have produced coexistence with free riders and cooperators interacting. The results rest on either of two mechanisms: vertical transmission with endogenous socialization efforts by parents, or horizontal transmission with endogenous norm strengths. Currently there are no convincing empirical tests of either mechanism, thus leaving a possible avenue for further research.

## **Discussion and Conclusion**

We have surveyed the economics literature in search of explanations for the coexistence of free riders and cooperators in social dilemma situations. In doing so, we distinguished between theories which take the heterogeneity of agents as given and focus on how institutions guide the behavior of rational

agents, and theories which try to explain the origin of preferences in processes of bounded rationality.

When it comes to punishment institutions, some naturally allow the expression of heterogeneous preferences into heterogeneous behavior. For example, when enforcement costs are positive, it is inefficient to deter all crime, leading to coexistence of cooperative and free-riding behavior. Complementarities in peer punishment lead to multiple equilibria, which can be interpreted as “parallel societies” governed by different social norms. Another set of theories rely on exclusion mechanisms, where cooperators manage to generate surplus among themselves, despite the existence of free riders.

In terms of models of bounded rationality, a few theories have been successful in explaining “proper coexistence,” where defectors and cooperators interact in stable proportions. Again, some of these models rely on the ability to exclude defectors and marginalize them to the fringes of the network (Fosco and Mengel 2011). Another promising avenue is cultural transmission mechanisms, where minority groups sustain themselves through intense socialization efforts by parents (Bisin et al. 2004).

These models provide starting points to think about coexistence. However, there are several reasons why they do not represent a unified body of research. First, there is a substantial gap between theory and empirical work. Theories of learning or cultural transmission predict long-run dynamics that are hard to isolate from other factors. As far as we know, no empirical validation exists for such models. Some empirical research has demonstrated cultural transmission (Dohmen et al. 2012), but testing of theoretical work is limited. Theories of the effects of institutions are easier to test (e.g., in the laboratory), but such tests have so far not produced valid explanations of coexistence. For example, theories of reputation formation predict either universal cooperation or noncooperation, but not the intermediate outcomes observed in actual experiments. Conversely, models of coalition formation, such as by Kosfeld et al. (2009), predict coexistence but the experiments find mostly full cooperation.

Second, this literature has not accounted for evolutionary foundations. While investigating economic institutions, economists simply assume heterogeneity in preferences for the payoffs of others, feelings of warm glow, altruism, or reciprocity. These “social preference” models have remained ad hoc, and there has been little discussion of their evolutionary origins (although see Alger and Weibull 2013). Just as social preference models lack evolutionary foundations, models of the dynamics of learning and evolution have paid little attention to the effect of legal or cultural institutions, although interesting leads are emerging that evaluate the effects of institutions on preferences for cooperation (Bowles and Polania-Reyes 2012). Clearly, future research needs to examine the two-way interactions between institutions and the development of individual tastes. Such models could help to explain the emergence and effects of institutional variations, like religion and government, and to evaluate the long-term effects of economic and political shocks.

Finally, there is ample evidence that social and moral preferences are more complex than assumed in the most popular economic models. The same individuals often behave quite differently in different dilemma situations (e.g., Blanco et al. 2011), and there are important spillover effects between behavior in different environments (e.g., Grimm and Mengel 2012). Thus, it is not clear how stable traits are over time and across different contexts. Models based on the assumption that people intrinsically favor altruistic behavior are unable to explain evidence that behavior that is “too altruistic” is viewed negatively and sometimes even punished (Herrmann et al. 2008), the fact that the framing of a decision problem often matters (Andreoni 1995), or that people seem to care about expectations of others (Charness and Duwenberg 2006). Some attempts have been made to represent the complex psychology of social and moral concerns in game theoretical models, often involving more parameters and more complex equilibrium concepts (Bénabou and Tirole 2006). Without evolutionary foundations or other ways to pin down parameters, such models may lack parsimony and predictive power.

In the process of writing this review, we found that explaining coexistence of free riders and cooperators has not been a goal of the economics literature. The literature summarized in this chapter focuses on explaining a particular behavior or trait, such as altruism or indirect reciprocity, or the impact of some institution, like climate coalition or a punishment scheme. Coexistence sometimes emerges as a byproduct of these endeavors, but it was never the direct object of inquiry.

To make progress, this will have to change. The economics literature could get inspiration by comparing models of social dilemmas in economics with behavioral ecology models of producers and scroungers. While both are models of investing and exploitation, coexistence is an equilibrium outcome in the latter, but not in the former, where the only Nash equilibrium (and hence also the only evolutionary stable state) involves free riding (see Burton-Chellew et al., this volume). More generally, we hope that the discussions brought forth in this volume will help both disciplines to exploit their common interest and help economists invest in new explanations for coexistence.

### **Acknowledgments**

We thank Sam Brown, Luc-Alain Giraldeau, Philipp Heeb, Kiryl Khalmetski, Michael Kosfeld, Julia Lupp, Fred Thomas, Björn Vollan, Bruce Winterhalder, and Arnon Lotem as well as participants at the Ernst Strüngmann Forum for valuable comments on an earlier version of this paper.