# CESIFO WORKING PAPERS

8285 2020

May 2020

### Belief-Dependent Motivations and Psychological Game Theory

Pierpaolo Battigalli, Martin Dufwenberg



#### **Impressum:**

**CESifo Working Papers** 

ISSN 2364-1428 (electronic version)

Publisher and distributor: Munich Society for the Promotion of Economic Research - CESifo

GmbH

The international platform of Ludwigs-Maximilians University's Center for Economic Studies and the ifo Institute

Poschingerstr. 5, 81679 Munich, Germany

Telephone +49 (0)89 2180-2740, Telefax +49 (0)89 2180-17845, email office@cesifo.de

Editor: Clemens Fuest

https://www.cesifo.org/en/wp

An electronic version of the paper may be downloaded

from the SSRN website: <a href="https://www.SSRN.com">www.SSRN.com</a>from the RePEc website: <a href="https://www.RePEc.org">www.RePEc.org</a>

· from the CESifo website: <a href="https://www.cesifo.org/en/wp">https://www.cesifo.org/en/wp</a>

## Belief-Dependent Motivations and Psychological Game Theory

#### **Abstract**

The mathematical framework of psychological game theory is useful for describing many forms of motivation where preferences depend directly on own or others' beliefs. It allows for incorporating, e.g., emotions, reciprocity, image concerns, and self-esteem in economic analysis. We explain how and why, discussing basic theory, experiments, applied work, and methodology.

JEL-Codes: C720, D910.

Keywords: psychological game theory, belief-dependent motivation, reciprocity, emotions, image concerns, self-esteem.

Pierpaolo Battigalli Bocconi University Milan / Italy pierpaolo.battigalli@unibocconi.it Martin Dufwenberg
University of Arizona
Tucson / AZ / USA
martind@eller.arizona.edu

May 2, 2020

We have benefited from many stimulating discussions (over the years) with our coauthors of the articles cited below. For their comments and advice we thank the Editor Steven Durlauf and several referees, as well as Chiara Aina, Lina Andersson, Geir Asheim, Enzo Di Pasquale, Francesco Fabbri, Price Fishback, Amanda Friedenberg, Joe Halpern, Tom Hwang, Kiryl Khalmetski, Senran Lin, Rachel Mannahan, Elena Manzoni, Paola Moscariello, Giulio Principi, Alexander Sebald, Joel Sobel, and Jin Sohn. Financial support of ERC (grant 324219) is gratefully acknowledged.

#### 1 Introduction

Economists argue that a rich variety of human motivations shape outcomes in important ways. Some categories (e.g., profit-maximization, altruism, inequity aversion, maximin preferences, or warm glow) can be readily handled using standard tools, most notably classical game theory. However, many other important sentiments which involve what we will call "belief-dependent motivation" defy standard analysis. A broader mathematical framework called "psychological game theory" (PGT), pioneered by Geanakoplos, Pearce & Stacchetti (1989) (GP&S) and further developed by Battigalli & Dufwenberg (2009) (B&D), provides adequate modeling tools. We explain how and why.

Let us provide early exposure to key ideas via three examples that relate to the following prominent categories of belief-dependent motivations:

- emotions, including guilt, disappointment, elation, regret, joy, frustration, anger, anxiety, suspense, shame, and fear;
- reciprocity, or the inclination to respond to kindness with kindness and to be unkind to whoever is unkind;
- image concerns, e.g., when someone wants others to believe that he is smart, altruistic, or honest.

Example 1: Guilt & tipping This example involves the emotions category. Psychologists Baumeister, Stillwell & Heatherton (1994) argue that "the prototypical cause of guilt would be the infliction of harm, loss, or distress on a relationship partner" and that if "people feel guilt for hurting their partners ... and for failing to live up to their expectations, they will alter their behavior (to avoid guilt) in ways that seem likely to maintain and strengthen the relationship" (see p. 245; cf. Tangney 1995). That outlook is reflected in the following vignette:

Tipper feels guilty if she lets others down. When she travels to foreign countries, and takes a cab from the airport, this influences

<sup>&</sup>lt;sup>1</sup>See also Gilboa & Schmeidler (1988) who in another pioneering contribution on "information-dependent games" anticipated some of the themes that GP&S and B&D developed in more depth. For now dated introductions to the older PGT literature, see Attanasi & Nagel (2008) and Dufwenberg (2008).

the gratuity she gives. Tipper gives exactly what she believes the driver expects to get, to avoid the pang of guilt that would plague her if she gave less.

To model this, consider game form  $G_1$  where Tipper (player 2) chooses tip  $t \in \{0, 1, ..., M\}$  and M > 0 is the amount of money in her wallet. The driver (player 1) is not active, and there is no future interaction. His (material) payoff is t. Choice t thus pins down an end-node. Tipper's material payoff is (M-t). However, her utility equals  $(M-t)-\theta_2\cdot [\tau-t]^+$ , where  $\tau$  is 1's expectation of t (which Tipper can only guess) and  $\theta_2 \geq 0$  is a sensitivity parameter measuring how much Tipper suffers when she lets 1 down. (Also,  $[\tau - t]^+ := \max\{\tau - t, 0\}$ .) In words, Tipper's utility equals own money minus a pang of guilt which is proportional to how much less the driver gets than what he expects to get. Tipper's behavior in the vignette is expected utility maximizing if  $\theta_2 > 1$ . The presence of  $\tau$  in her utility makes it beliefdependent, leading to what is called a "psychological game" (p-game) here given by  $G_1^*$ . The key characteristic is that utility at an end-node depends on beliefs, in this case that 2's utility at t depends on 1's beliefs (via  $\tau$ ). Had we had a traditional game, utilities would be defined on end-nodes independently of beliefs and Tipper's best choice would be independent of her guess of  $\tau$ .

 $[G_1 \text{ and } G_1^*]$ 

Terminology Example 1 illustrates key concepts we rely on throughout: A game form specifies the structure of a strategic situation (the "rules of the game"): the players, how they can choose, and the material consequences (typically money) of players' actions. We reserve the term payoffs for material consequences. Unless players are expected payoff maximizers, payoffs do not represent preferences over end-nodes.<sup>2</sup> These are instead given by utility functions (or utilities). Whereas in traditional game theory utilities are defined on end-nodes only, in PGT they also depend on features of beliefs about behavior, like  $\tau$  in Example 1, and higher-order beliefs. Such beliefs are determined by the strategic analysis, i.e., they are endogenous. We use the term motivations to distinguish conceptually different parts that may affect utility. Tipper is affected by two motivations: own money and guilt. We use the term belief-dependent to describe whether or not a motivation or utility depends on beliefs. Tipper's first motivation (own money) is

<sup>&</sup>lt;sup>2</sup>Or, more generally, probability distributions over nodes.

not belief-dependent while the second one (guilt) is, which implies that Tipper's utility is belief-dependent. We call **psychological game**, or **p-game** for short, the model obtained by appending belief-dependent utilities to the given game form.

Example 2: Reciprocity in the battle-of-the-sexes The idea that people wish to be kind towards those they perceive to be kind, and unkind towards those deemed unkind, is age-old.<sup>3</sup> Early academic discussions can be found in anthropology (Mauss 1954), social psychology (Goranson & Berkowitz 1966), biology (Trivers 1971), and economics where the pioneer is Akerlof (1982), who analyzed "gift-exchange" in labor markets. Akerlof had the intuition that reciprocity would imply a monotone wage-effort relationship (at least up to the level of a "fair wage"), and he posited that such a relationship exists. However, he did not engage in mathematical psychology and formal description of the underlying affective processes. Rabin (1993) realized that such an approach could bring about a generally applicable model, which he developed. Our second example is taken from him.

 $[G_2]$ 

Consider game form  $G_2$ . If the players were motivated solely by their material payoffs we would have a traditional game, with two equilibria: (opera, opera) and (boxing, boxing). These strategy profiles remain equilibria in Rabin's model, where players' utilities are affected by reciprocity and therefore (opera, boxing) is an additional equilibrium. We describe the underlying intuition rather than the mathematical details: The players are "unkind" to each other, in the sense that given equilibrium expectations they minimize each other's material payoffs (to be 0 rather than 1). Each player is sacrificing some own material payoff in the process (getting 0 rather than  $\frac{1}{2}$ ), but the desire to reciprocate the perceived unkindness of the co-player is strong enough to make it worthwhile.

Section 2 explains in more detail why modeling reciprocity involves PGT. The reason is that kindness depends on beliefs. Here we quote Rabin (p.

 $<sup>^3</sup>$ Fehr & Gächter (2000, p. 159) reproduce a  $13^{th}$  century quote from the Edda that conveys the spirit: "A man ought to be a friend to his friend and repay gift with gift. People should meet smiles with smiles and lies with treachery." Dufwenberg, Smith & Van Essen (2013, Section III) give more examples, from popular culture, business, and experiments. Sobel (2005) provides a broad critical discussion.

1285), who compares the (boxing, boxing) and (opera, boxing) equilibria, highlighting a non-standard aspect of his model:

In the natural sense, both of the equilibria ... are strict: each player strictly prefers to play his strategy given the equilibrium. In the equilibrium (boxing, boxing), player 1 strictly prefers playing boxing to opera. In the equilibrium (opera, boxing) player 1 strictly prefers opera to boxing. No matter what payoffs are chosen, these statements would be contradictory if payoffs depended solely on the actions taken.

**Example 3: Status & conformity** Our next example illustrates an image concern as modeled by Bernheim (1994). A special case goes as follows: Agents in a population are uniformly distributed on T = [0, 2], where  $t \in T$  is an agent's "type" of preference for "brightness of clothing". Each agent simultaneously chooses a (garment) color  $c \in T$  to wear. All agents observe these choices, and form beliefs about the type of the chooser conditional on the choice. Let  $t_c$  denote the expected type of an agent who chose c. An agent of type t's utility equals  $-(t-c)^2 - (1-t_c)^2$ . In words, he suffers quadratic losses of (i) letting his chosen color deviate from the one he favors, and of (ii) status by being *perceived* as having an expected type that deviates from 1 (the "fashion standard").

Focusing on the case where agents' types are private information, Bernheim analyzes this model as a signaling game. He looks for equilibria where the "sender" (the only active player) maximizes expected utility given the way beliefs are formed, while beliefs are formed consistently with Bayes' rule given how choices depend on types. He argues that a plausible class of equilibria involve pooling at c = 1. Under our parameterization, such pooling can be universal, if out-of-equilibrium inferences—which cannot be pinned down by Bayes' rule—satisfy (e.g.) that  $t_c \in \{0, 2\}$ , for all  $c \neq 1$ .

The example is interesting to us for two reasons. First, aspect (ii), mentioned above, makes utility belief-dependent,<sup>5</sup> and so creates a p-game  $G_3^*$ . Second, consider a modified version of  $G_3^*$ , call it  $G_3^{**}$ , where agents' types

<sup>&</sup>lt;sup>4</sup>If an agent of type t stays with the proposed equilibrium he gets utility  $-(t-1)^2 - 0$ . If he deviates, the best way to do so would be to choose c = t, in which case he would get utility  $-0 - (1-0)^2$  (or  $-0 - (1-2)^2$ ), hence he cannot gain by deviating.

<sup>&</sup>lt;sup>5</sup>Note that  $t_c$  is a feature of an endogenous belief, because it is derived from an initial belief about types and choices by conditioning on the observed choice.

are observed ex post. Obviously, an agent of type t will be believed to have type t, regardless of his choice c. In the unique equilibrium, he will rely on a dominant strategy: c = t, so the prediction will differ from Bernheim's. A striking observation, from a game-theoretic point of view, is that the difference between  $G_3^*$  and  $G_3^{**}$  concerns information across end-nodes. It is imperfect in  $G_3^*$  but perfect in  $G_3^{**}$ . In traditional game theory, information across end-nodes never affects predictions, and is therefore not even specified. That this property does not extend to p-games shows that information at end-nodes (and more generally the information players have when they are inactive) should be carefully specified.

What's next? The above examples illustrate different belief-dependent motivations that lead to belief-based utility and p-games. They are suggestive of how one needs to move beyond traditional game theory for modeling purposes, and how the framework of PGT can do this. PGT provides a useful intellectual umbrella, under which many trends in psychology and economics can be understood, related, and synthesized via appropriate p-games. Awareness of and interest in PGT is on the rise, yet far from universal. In this article, we explain what PGT is and what motivations can be modeled, highlighting a variety of idiosyncratic features. We present old insights and speculate about new ones that PGT may deliver. We discuss basic theory, experimental tests, and applied work. Although we cite a lot of papers, our primary goal is to highlight the structure and potential of various forms of work involving PGT. Our style is semi-formal, presenting some notions verbally rather than mathematically. Readers who wish to dig deeper should compare with relevant passages of GP&S, B&D, and other articles. This includes, in particular, the recent methodological article by Battigalli, Corrao & Dufwenberg (2019) (BC&D), a text we frequently draw connection to.

Our discussion mostly consists of showing how to functionally represent the belief-dependent motivations in question and highlighting the ensuing best-reply behavior. Sometimes we analyze strategic reasoning either by iterated elimination of non-best replies, or informally applying an equilibrium concept. For a broader discussion of solution concepts see B&D and BC&D, as well as our brief critical remarks in Section 9.

In writing this text we faced a pedagogical challenge. On one hand, notions of belief-dependent motivation are exciting and intuitive, although likely novel and unusual to many. With this, we wish to allow our read-

ers to develop a feel quickly. On the other hand, serious study involving PGT requires that one develops a precise mathematical framework. We also wish to explain how to do that. Had this been a textbook, we would have presented the formal framework before exploring all the notions of belief-dependent motivation that can be analyzed. However, some of our referees convinced us that such an approach might scare off potentially interested readers, who might stop reading before they get to the sections where all the exciting forms of motivation are described. So, we switched the order. In the next four sections (2-5) we elucidate the palette of sentiments that PGT can explore (starting with the three prominent categories of motivations mentioned above). This is doable, because we mainly focus on simple examples for which we can provide needed technical side-comments as we go. The heavy-lifting is postponed until Section 6, which presents the formal framework underlying the analysis of the preceding sections. We encourage readers who like formal analysis to read Section 6 before the preceding ones.

Sections 7-10, finally, concern PGT-related experiments, applications, critical comments on methodology, and wrap-up remarks.

**Abbreviations** Those defined above—PGT, p-games, GP&S, B&D, and BC&D—will be used throughout!

#### 2 Reciprocity

Rabin's model of kindness-based reciprocity pioneered using PGT to explore the general implications of a particular motivation. He focuses on simultaneous-move game forms, as we illustrated via  $G_2$ . But, as Rabin himself points out (p. 1296)—from the perspective of applied economics—it is important to also consider extensive game forms with a non-trivial dynamic structure. Dufwenberg & Kirchsteiger (2004) took on that task,<sup>6</sup> and we sketch their approach. Game form  $G_4$  (akin to their  $\Gamma_1$ ) is useful for introducing main ideas:

 $[G_4]$ 

<sup>&</sup>lt;sup>6</sup>The main difference between Rabin's and Dufwenberg & Kirchsteiger's approaches concerns which class of game forms is considered, but there are other differences too. See Dufwenberg & Kirchsteiger (2004, Section 5; 2019).

A crucial building block of the analysis concerns player i's kindness to j, denoted  $\kappa_{ij}(\cdot)$ .<sup>7</sup> It is the difference between the payoff (i.e., the material/monetary reward) i believes j gets (given i's choice) and a comparison payoff C that i computes as follows: C is the average of the minimum and the maximum payoff that i believe j could get, for other choices of i.<sup>8</sup> In  $G_4$ , if 1 believes there is probability p that 2 would choose take, we get

$$\kappa_{12}(stay, p) = 5 - \frac{1}{2} \cdot [5 + (p \cdot 9 + (1 - p) \cdot 1)] = 2 - 4 \cdot p,$$

$$\kappa_{12}(reach, p) = p \cdot 9 + (1 - p) \cdot 1 - \frac{1}{2} \cdot [5 + (p \cdot 9 + (1 - p) \cdot 1)] = 4 \cdot p - 2,$$

$$\kappa_{21}(take) = 1 - \frac{1}{2} \cdot [1 + 9] = -4, \text{ and}$$

$$\kappa_{21}(give) = 9 - \frac{1}{2} \cdot [1 + 9] = 4.$$

Note that i's kindness to j has the dimension of the (expected, material) payoff of j, it ranges from negative to positive, and it may depend on i's beliefs (as it does for 1 in  $G_4$ ). Player i is taken to maximize (the expectation of) a utility that depends on actions and beliefs according to a functional form of the following kind:

$$u_i(\cdot) = \pi_i(\cdot) + \theta_i \cdot \kappa_{ij}(\cdot) \cdot \kappa_{ji}(\cdot), \tag{2}$$

where  $\pi_i(\cdot)$  is i's (material) payoff function and parameter  $\theta_i \geq 0$  reflects i's reciprocity sensitivity. The desire to reciprocate kindness, as intuitively described in the Introduction, is captured via "sign-matching;"  $\theta_i \kappa_{ij}(\cdot) \kappa_{ji}(\cdot)$  is positive only if the signs of  $\kappa_{ij}(\cdot)$  and  $\kappa_{ji}(\cdot)$  match. To illustrate in  $G_4$ : if  $\theta_2$  is high enough, 2 wants to "surprise" 1, i.e., 2 chooses take if  $p < \frac{1}{2}$  and give if  $p > \frac{1}{2}$ .

We make several PGT-related observations:

(i) Player 2 chooses between end-nodes. So, in traditional game theory, her optimal choice would be independent of beliefs. This is not the case with

<sup>&</sup>lt;sup>7</sup>Here, and in other expressions below, the dot symbol  $(\cdot)$  represents on one or more variables, such as chosen actions (terminal history reached) and beliefs.

<sup>&</sup>lt;sup>8</sup>This definition neglects an important aspect that is commented on below under the heading "Dealing with 'bombs."

<sup>&</sup>lt;sup>9</sup>Player *i* cannot know *j*'s beliefs and must form beliefs about  $\kappa_{ji}(\cdot)$ , denoted  $\lambda_{iji}(\cdot)$  by Dufwenberg & Kirchsteiger, who plug  $\lambda_{iji}(\cdot)$  into  $u_i$ . Our formulation, (2), conformant with Section 6 below, relies on first-order beliefs only, but has equivalent implications.

reciprocity. In  $G_4$ , 2's optimal choice depends on p, 1's belief. This illustrates that  $G_4$ , when played by agents motivated by reciprocity, is a p-game.

- (ii) Relatedly, backward induction can not be used to find 2's optimal choice independently of beliefs. Player 2 must consult her beliefs about p to compute his expected-utility-maximizing action.
- (iii) In traditional game theory, finite perfect-information games have equilibria (justifiable by backward induction) where players rely on degenerate, deterministic plans (intended choices). This is not the case in  $G_4$ , for high values of  $\theta_2$ . We have not defined equilibrium here, but suppose we have a notion of equilibrium that requires 1 to correctly anticipate 2's plan (and that plans are carried out), and for 2 to anticipate that 1 will do so." (Dufwenberg & Kirchsteiger's equilibrium has that property.) If 2 plans to choose take, and 1 anticipates that 2 plans to choose take, then p = 1. But, if 1 anticipates that, then (as explained above) 2's best response would be give, not take. An analogous argument rules out an equilibrium where 2 plans to choose give.

Our next example, the Mini-ultimatum Game form  $G_5$ , gives further insights regarding reciprocity, and will be used for later comparisons as well:

 $[G_5]$ 

Reasoning as before (with p now 1's belief about reject),  $\kappa_{12}(greedy, p)$  is strictly negative for all p.<sup>10</sup> If  $\theta_2$  is large enough, the utility maximizing plan for 2 is reject. Suppose this case is at hand. What should 1 do? If  $\theta_1 = 0$ , meaning that 1 is selfish, then 1 would choose fair (since 5 > 0). If instead  $\theta_1$  is large (enough), then there are two possibilities. The first one is that 1 chooses fair. To see why, suppose that (at the root, i.e., before the start of play) 1 believes that 2 believes that 1 plans to choose fair. Then 1 believes that 2 believes that 2 is not (as evaluated at the root) affecting 1's payoff. That is, at the root, it holds that  $\kappa_{21}(\cdot) = 0$ , implying that, to maximize his utility, 1 should act as if selfish and choose fair (since 5 > 0). The second, very different, possibility is that 1 chooses greedy, despite the anticipation that 2 will choose reject. This is a "street fight" outcome, with negative reciprocity manifesting along the path of play. To get the intuition, suppose 1 believes that 2 believes (at the root) that 1 is going to choose *qreedy*. Then 1 believes that 2 is planning to generate a payoff of 0 rather than 9 for player 1. In this case, 2 would be unkind to 1. Since  $\theta_1$  is large, 1 reciprocates (in

anticipation!) choosing *greedy* thereby generating a payoff of 0 rather than 5 for player 2.

The analysis here reflects a key feature of the approach, namely that players' kindness is re-evaluated at each history. For example, 2's kindness to 1 at the root may be zero (if 2 believes 1 plans to choose fair) and yet 2's kindness after 1 chooses greedy would at that time not be zero.<sup>11</sup>

**Dealing with "bombs"** The account of reciprocity theory just given glosses over a subtle issue which we now flag. To illustrate, let  $G_5^X$  be a modification of  $G_5$  such that player 1 has a third choice at the root—X—which explodes a bomb, leaving each player with a material payoff of -100.

Recall how we defined i's kindness to j as the difference between the payoff i believes j gets and the average of the minimum and maximum payoff that i believes j could get.  $G_5^X$  can illustrate how, in some game forms, absurd implications follow unless the calculation of "the minimum payoff i believes j could get" is modified to not consider choices that hurt both i and j. In  $G_5$  we concluded that 1's kindness when choosing greedy was negative  $(\kappa_{12}(greedy, p) = -2 - \frac{p}{2}$ , as noted in a footnote). Reasoning analogously, in  $G_5^X$  we would get that 1's kindness of choice greedy is instead positive. Arguably, this is implausible; while hurting everyone would surely be unkind, not doing so shouldn't automatically render other choices kind. Arguably (for a given p) the kindness of choice greedy should rather be the same in  $G_5^X$  as in  $G_5$ .

Dufwenberg & Kirchsteiger (2004), as well as Rabin, propose kindness definitions that achieve such an objective, by calculating "the minimum payoff j could get" without regard to so-called "inefficient strategies" that hurt both i and j. Their approaches, while to a large degree similar in spirit, differ in details. The (somewhat contentious) issues involved are too subtle to warrant coverage here. We refer to Dufwenberg & Kirchsteiger (2019) for a detailed discussion, including a response to a related critique by Isoni & Sugden (2019).

Related literature Dufwenberg & Kirchsteiger (2004) limit attention to certain game forms without chance moves, a restriction Sebald (2010) drops,

<sup>&</sup>lt;sup>11</sup>Our account has been sketchy; see van Damme et al. (2014; Section 6, by Dufwenberg & Kirchsteiger) for a fuller analysis of a large class of Ultimatum Game forms.

<sup>&</sup>lt;sup>12</sup>More precisely,  $\kappa_{12}(greedy, p) = (1-p) \cdot 1 - (\frac{1}{2} \cdot 5 + \frac{1}{2} \cdot (-100)) = 48.5 - p.$ 

which allows him to address broader notions of "attribution" and "procedural concerns". Sohn & Wu (2019) analyze situations where players are uncertain about each other's reciprocity sensitivities. Jiang & Wu (2019) discuss alternatives to the belief-revision rules of Dufwenberg & Kirchsteiger (2004). Dufwenberg, Smith & Van Essen (2013) modify the theory to focus on "vengeance;" players reciprocate negative but not positive kindness (achieved by replacing  $\kappa_{ji}(\cdot)$  in (2) by  $[\kappa_{ji}(\cdot)]^-$ ). All these authors hew close to Rabin. Alternative approaches are proposed by Falk & Fischbacher (2006) who combine reciprocity motives with preferences for fair distributions, <sup>13</sup> and Çelen, Schotter & Blanco (2017) who model i's reciprocation to j based on how i would have behaved had he been in j's position.

As PGT-based models gain popularity they will be increasingly used to do applied economics. Most such work to date is based on reciprocity theory (and in particular Dufwenberg & Kirchsteiger's 2004 model). Topics explored include wage setting, voting, framing effects, hold-up, bargaining, gift exchange, insolvency in banking, mechanism design, trade disputes, public goods, randomized control trials, memoranda of understanding, climate negotiations, communication, and performance-based contracts.<sup>14</sup>

#### 3 Emotions

For a long time, neither psychologists nor economists paid much attention to emotions and how they shape behavior. We recommend Keltner & Lerner's (2010) handbook chapter which explains how while "founding figures in psychology" (in particular Charles Darwin and William James) paid significant attention to emotions, during most of the  $20^{th}$  century and "the heyday of behaviorism ... emotions resided ... outside the purview of observable mea-

<sup>&</sup>lt;sup>13</sup>So do Rabin (1993, p. 1298) as well as Charness & Rabin (2002) in appendix-versions of their social preference models. These models and the references in the main text are PGT-based. Levine (1998), Cox, Friedman & Sadiraj (2008), and Gul & Pesendorfer (2016) present reciprocity-related ideas which are not kindness-based and do not use PGT.

<sup>&</sup>lt;sup>14</sup>See Dufwenberg & Kirchsteiger (2000), Hahn (2009), Dufwenberg, Gächter & Hennig-Schmidt (2011), Dufwenberg et al. (2013), van Damme et al. (2014; Section 6), Netzer & Schmutzler (2014), Dufwenberg & Rietzke (2016), Bierbrauer & Netzer (2016), Bierbrauer, Ockenfels, Pollak & Rückert (2017) Conconi, DeRemer, Kirchsteiger, Trimarchi & Zanardi (2017), Dufwenberg & Patel (2017), Jang, Patel & Dufwenberg (2018), Kozlovskaya & Nicolò (2019), Aldashev, Kirchsteiger & Sebald (2017), Nyborg (2018), Le Quement & Patel (2018), and Livio & De Chiara (2019).

surement" and were considered "undeserving of scientific inquiry" (p. 317).<sup>15</sup> Furthermore, Elster (1996, 1998) forcefully argues that economists by and large have neglected to study the emotions, despite that the topic is potentially of great importance. In his 1996 text he goes so far as to note that "all human satisfaction comes in the form of emotional experiences" (p. 1368). He argues that by failing to recognize such an important source of utility economists are potentially failing to get a correct grip on how decisions are formed.

That view is corroborated by the development in more recent psychology. According to Keltner & Lerner, not only has (since 1980) "a robust science of emotion ... emerged" (p. 317), but it has indicated that a large variety of emotions, each one in distinct ways, impacts well-being and behavior. The causalities are complex and hardly fully understood, but a key idea that is often stressed involves what since Lerner & Keltner (2000, 2001) has been called "appraisal-tendency." Lerner, Li, Valdesolo & Kassam (2014) discuss in particular the implications for decision making and how "appraisal tendencies are goal-directed processes through which emotions exert effects on judgments and decisions" (p. 479). The themes they bring up include how emotions affect content and depth of thought, goal activation, and interpersonal assessments. <sup>16</sup>

Reading these psychological discussions is highly inspiring, and we encourage economists to do so. Yet, at times, getting a full grip can be frustrating as the concepts and connections tend to be, not only overwhelmingly plentiful, but also informal. We suspect and hope that some complementary clarity can be brought to the table by invoking analytical methods. PGT provides an adequate set of tools. In his previous article in this *Journal*, Elster (1998) argued that emotions "are triggered by beliefs" (p. 49) and that they can have important economic consequences. How "can emotions help us explain behavior for which good explanations seem to be lacking?" he asked (p. 48). While he lamented economists' dearth of attention to the issue, PGT has subsequently been put to such use, and there is more to do.

 $<sup>^{15}</sup>$ Keltner & Lerner quote Skinner (1948) who wrote that emotions "are the fictional causes to which we ascribe behavior" and "useless and bad for our peace of mind and out blood pressure."

<sup>&</sup>lt;sup>16</sup>See also Keltner & Lerner's Table 9.3 and the related discussion of attention, certainty, control coping, pleasantness, responsibility, legitimacy, and anticipated effort.

<sup>&</sup>lt;sup>17</sup>See also Chang & Smith (2015) who elaborate on this theme.

#### 3.1 Guilt

Among the emotions, guilt has been explored the most using PGT.<sup>18</sup> Motivated by work in psychology (e.g., Baumeister *et al.* and Tangney, cited in the Introduction), Battigalli & Dufwenberg (2007) develop a model allowing exploration of how (two versions of) guilt shapes strategic interaction in a general class of extensive game forms. While most follow-up work has been experimental (see Section 7), a few applied theory papers explored how guilt influences marriage & divorce, corruption, deception, framing, tax evasion, public goods, embezzlement, and expert advice.<sup>19</sup>

We provide (BC&D's account of) Battigalli & Dufwenberg's (2007) notion of "simple guilt." Player i experiences guilt when he believes that the payoff j gets  $(\pi_j(\cdot))$  is lower than the payoff j initially expected given j's first-order beliefs  $\alpha_j$ . This expectation is denoted  $\mathbb{E}[\pi_j; \alpha_j]$ , and it depends on j's beliefs about (own and others') actions.<sup>20</sup> Specifically,  $i \neq j$  maximizes (the expectation of) a utility of the form

$$u_i(z,\alpha_i) = \pi_i(z) - \theta_i \cdot \left[ \mathbb{E}[\pi_i;\alpha_i] - \pi_i(z) \right]^+, \tag{3}$$

where z is the sequence of chosen actions (**terminal history**, **path**, or **end-node**). Again,  $\theta_i \geq 0$  is a sensitivity parameter. As seen in the Introduction, Tipper's behavior in  $G_1$  is captured if  $\theta_2 > 1$ . We now discuss also a trust game form  $G_6$ .<sup>21</sup> Assume that  $\theta_1 = 0$  and  $\theta_2 > 0$  to get the p-game  $G_6^*$ , displayed alongside, where  $\hat{\pi}_1 = \mathbb{E}[\pi_1; \alpha_1] \in [0, 10]$  denotes 1's expected payoff.<sup>22</sup>

$$[G_6 \text{ and } G_6^*]$$

<sup>&</sup>lt;sup>18</sup>Reciprocity, which we do not count as an emotion, has been explored even more than guilt. See Azar (2019) for a statistical analysis of the bibliometric impact of PGT-based reciprocity and guilt theory.

<sup>&</sup>lt;sup>19</sup>See Dufwenberg (2002), Balafoutas (2011), Battigalli, Charness & Dufwenberg (2013), Dufwenberg & Nordblom (2018), Dufwenberg *et al.* (2011), Patel & Smith (2019), Attanasi, Rimbaud & Villeval (2019), and Khalmetski (2019).

 $<sup>^{20}</sup>$ The authors actually assume that i suffers only to the extent that he causes j to get a lower payoff than j initially expected. Stating that precisely leads to a more complicated utility than the one seen here. However, best responses are identical, so we opt for the simpler version here.

<sup>&</sup>lt;sup>21</sup>Compare, e.g., Huang & Wu (1994), Dufwenberg (2002), Dufwenberg & Gneezy (2000), Bacharach, Guerra & Zizzo (2004), and Charness & Dufwenberg (2006).

<sup>&</sup>lt;sup>22</sup>Note that  $[\hat{\pi}_1 - \pi_1 (trust, grab)]^+ = \hat{\pi}_1$  because  $\hat{\pi}_1 \ge 0 = \pi_1 (trust, grab)$ .

 $G_6^*$  is a p-game, because of the presence of  $\hat{\pi}_1$ , an expectation derived from 1's beliefs. One may think of 2's utility as reflecting a form of "state-dependent" preference, i.e., what 2 would prefer if he knew  $\hat{\pi}_1$ . To maximize her utility, 2 must consult her beliefs about  $\hat{\pi}_1$ .<sup>23</sup>

In some strategic settings, powerful predictions may obtain if players reason about each other's reasoning. This may be relevant in p-games, and the emotion of guilt, as modeled in  $G_6^*$ , can illustrate this in a stark way: If  $10 > 14 - \theta_2 \hat{\pi}_1$ , then 2 prefers share over grab, and vice versa. No matter how high  $\theta_2$  is, if  $\hat{\pi}_1$  is low enough 2 prefers grab over share. Nevertheless, 2 may reason that if 1 chose trust then  $\hat{\pi}_1 \geq 5$ , since otherwise 1 would not be rational. If  $\theta_2 > \frac{4}{5}$ , player 2 will then prefer share over grab, and if 1 believes that 2 will reason that way, he should choose trust.<sup>24</sup>

As argued by Charness & Dufwenberg (2006), simple guilt can help explain why communication can foster trust and cooperation. Suppose  $G_6/G_6^*$  is augmented with a pre-play communication opportunity and that 2 promises 1 to choose share. If 1 believes this, and if 2 believes that 1 believes this, then simple guilt makes 2 live up to her promise. A promise by 2 feeds a self-fulfilling circle of beliefs about beliefs that share will be chosen. We emphasize that guilt, per se, does not imply such a positive effect of communication (nor does it rule out a negative effect), but it is consistent with it.<sup>25</sup>

Let us finally discuss the following three guilt-related distinct topics:

Counterfactual emotions In  $G_6^*$ , if 2 chooses *share* to avoid guilt, then 2 will (along the realized path) *not* experience guilt. Nevertheless, guilt has shaped the outcome. This illustrates a more general phenomenon: An

 $<sup>^{23}</sup>$ Early work on guilt (e.g., Dufwenberg 2002) plugged that second-order belief (rather than  $\hat{\pi}_1$ ) into  $u_2$ . As explained by B&D, the two approaches are equivalent. We prefer our chosen one. The shape of 2's utility is kept simpler with only first-order belief in its domain (see Section 6 where we only consider such beliefs).

<sup>&</sup>lt;sup>24</sup>Dufwenberg (2002) calls this line of reasoning "psychological forward induction." See B&D, BC&D, and Battigalli *et al.* (2020) for more discussion and formalization via extensive-form rationalizability.

 $<sup>^{25}</sup>$ In other game forms, one may argue that if a vulnerable party, say player i, were afraid that a guilt averse player j would take an action that could hurt i, then i might wish to tell j either that he had "high expectations" or that (for given expectations) the loss due to the hurtful action would be large. These are other ideas that link guilt aversion and communication, which have been explored by Cardella (2016) and Caria & Fafchamps (2019).

emotion (it could also be, e.g., disappointment or regret, as we'll see in coming sections) need not actually realize in order to affect economic outcomes.

This observation marks a difference, to a degree, between what is the natural focus of economists and psychologists. For economists it is obvious that a counterfactual emotional experience is important, if it influences behavior and who gets what. Psychologists' discussions, by contrast, tend to focus on the impact of guilt when it actually occurs. The quote from Baumeister *et al.*, regarding guilt, which we included in the Introduction, is exceptional.

Expecting too much? Battigalli & Dufwenberg's (2007) model does not distinguish whether or not a belief by j is "reasonable," as regards whether or not guilt of i can be triggered. This assumption was made in order to keep things simple, and it could be unrealistic. For example, in  $G_1$ , if M is large and the driver expected Tipper to give away all she has then she might plausibly find the driver obnoxious, and enjoy giving nothing! Balafoutas & Fornwanger (2017) and Danilov, Khalmetski & Sliwka (2019) discuss such "limits of guilt".

Guilt vs. reciprocity With reference back to Section 2, the following points of comparison are noteworthy. First, in  $G_6$ , the incorporation guilt or reciprocity has opposite connections between q and 2's preference. Under simple guilt (i.e., in  $G_6^*$ ), the higher is q the more inclined 2 will be to choose share (see (iii)). However, the higher is q the less kind is 1 (reasoning as in Section 2), so if 2 were motivated by reciprocity a higher q would spell less inclination to choose share.<sup>26</sup>

Second, under simple guilt, a single utility function, that depends on initial payoff expectations and on which end-node is reached, can be applied at each history where a player moves. By contrast, to capture reciprocity motivation one must describe and re-evaluate each player's kindness at each history.<sup>27</sup>

<sup>&</sup>lt;sup>26</sup>For more on this, see Attanasi, Battigalli & Nagel (2013).

<sup>&</sup>lt;sup>27</sup>Herein lies *two* differences: First, a new utility function is needed for each history; see Dufwenberg & Kirchsteiger (2004) for more on this feature, which we have not illustrated very clearly since players moved once in the games we considered. Second, since kindness depends on (foregone) choice options, game-form details matter in a way that lacks counterparts with simple guilt. See BC&D for a detailed discussion of this distinction, concerning "game-form free" vs. "game-form dependent" preferences.

Third, recall our above remark regarding how, in  $G_6/G_6^*$ , if guilt may make 2 chooses *share* then 2 will *not* experience guilt. By contrast, if 2 were instead motivated by reciprocity, her belief-dependent motivation might be felt as she chooses *share*; at that time she perceives 1 as kind (in inverse proportion to q) which influences her utility as she chooses.

#### 3.2 Disappointment

Dufwenberg (2008) gives the following example which illustrates a critical role of prior expectations:

I just failed to win a million dollars, and I am not at all disappointed, which however I clearly would be if I were playing poker and knew I would win a million dollars unless my opponent got lucky drawing to an inside straight, and then he hit his card.

Belief-dependent disappointment was first modeled by Bell (1985) and Loomes & Sugden (1986). More recent work by Kőszegi & Rabin (2006, 2007, 2009) and also Shalev (2000) is technically closely related, but since it is differently motivated we write about it under the separate heading of "Belief-dependent loss aversion" in Section 5 below. Gill & Prowse (2011) argue that disappointment may help explain behavior in tournaments for "promotions; bonuses; professional partnerships; elected positions; social status; and sporting trophies" (p. 495).

The needed modeling machinery was in part present already in the part on guilt of Section 3. Factor  $[\mathbb{E}[\pi_j; \alpha_j] - \pi_j(z)]^+$ , seen in eq. (3), captures j's disappointment, although in (3) it was used for the purpose of modeling i's guilt.<sup>28</sup> To let i's utility reflect disappointment we can instead look at

$$u_i(z,\alpha_i) = \pi_i(z) - \theta_i \cdot [\mathbb{E}[\pi_i;\alpha_i] - (\pi_i(z) + k)]^+, \tag{4}$$

where  $k \geq 0$ . In words, i's utility equals money minus a pang of disappointment which is linked to his prior expectation. Note that k = 0 incorporates disappointment in the most straightforward way, while if k > 0 then

<sup>&</sup>lt;sup>28</sup>This suggests an alternative way to think of i's guilt towards j, namely that i is averse to j being disappointed.

i discounts "small" disappointments so that they have no effect on utility.<sup>29</sup> Below, we consider a case with k > 0 to make a technical point.

Utility (4) looks deceptively similar to (3) but is crucially different in that i's utility depends (in part) on i's plan, that is, the part of  $\alpha_i$  representing i's beliefs about the actions he is going to take (more on that in Section 6). Such "own-plan dependence," where i's beliefs about his choices impacts the utility of his choices, can lead to subtle complications as we now highlight (and see BC&D for more).

While (4) is applicable to any game form, and hence can shape strategic interaction generally, the clearest way to exhibit the essence of disappointment is to use a one-player game form with chance moves, like  $G_7$ . Assume that 0 < x < 1 while  $0 \le k \le \min\{x, 1 - x\}$ .

$$[G_7]$$

Can stay be an optimal plan for 1 in  $G_7$  (with utility given by (4))? This requires

$$\underbrace{x}_{\text{utility of } stay} \ge \underbrace{\frac{1}{2} \cdot 2 - \frac{1}{2} \cdot \theta_1 \cdot [x - (0 + k)]^+}_{\text{utility of } bet \text{ after planning } stay} \iff x \ge \frac{2 + \theta_1 \cdot k}{2 + \theta_1}. \quad (5)$$

Similarly, bet is an optimal plan if

$$\underbrace{\frac{1}{2} \cdot 2 - \frac{1}{2} \cdot \theta_1 \cdot [1 - (0 + k)]^+}_{\text{utility of } \text{bet after planning } \text{bet}} \ge \underbrace{x - \theta_1 \cdot [1 - (x + k)]^+}_{\text{utility of } \text{stay after planning } \text{bet}} \iff x \le \frac{2 + \theta_1 - \theta_1 \cdot k}{2 + 2 \cdot \theta_1}.$$
(6)

First, assume that k=0. Inspecting (5) and (6) one sees that if  $x\in [\frac{2}{2+\theta_1},\frac{2+\theta_1}{2+2\theta_1}]$  then either stay or bet can be an optimal plan. If  $x\in (\frac{2}{2+\theta_1},\frac{2+\theta_1}{2+2\theta_1})$  then 1 incurs a loss if he deviates from the plan. Such multiplicity of strictly optimal plans could never happen without own-plan dependent utility.<sup>30</sup> In

<sup>&</sup>lt;sup>29</sup>Disappointment aversion may violate first-order stochastic dominance. For example, if k in eq. (4) is 0 and  $\theta_i > 1$ , then i prefers a sure payoff x > 0 to the lottery that yields x and 2x with 50% chance. The axiomatization of Gul (1991) rules this out. Cerreia-Vioglio, Dillenberger & Ortoleva (2018) derive an explicit representation of preferences à la Gul (1991).

<sup>&</sup>lt;sup>30</sup>This statement is true if there is perfect recall; otherwise similar complications occur as, again, dynamically inconsistent preferences may appear, and the conditional expected utility of actions may depend on the planned probability of choosing "earlier" actions. See, e.g., Piccione & Rubinstein (1997), which is the lead article in a special issue devoted to imperfect recall.

the standard case, multiplicity of optimal plans is possible only if there is indifference.

An interesting variation arises if k > 0. Could it be that neither stay nor bet is an optimal plan? If so, then neither (5) nor (6) would hold, and we would get

$$\frac{2 + \theta_1 \cdot k}{2 + \theta_1} > x > \frac{2 + \theta_1 - \theta_1 \cdot k}{2 + 2 \cdot \theta_1}.$$
 (7)

To see that this is possible, pick a case that is easy to compute: assume that  $x = k = \frac{1}{2}$ , and study (7) as  $\theta_1$  increases. The leftmost term exceeds  $\frac{1}{2}$  for any  $\theta_1 \geq 0$ , while the rightmost term is lower than  $\frac{1}{2}$  for high enough  $\theta_1$  (it decreases from 1 to  $\frac{1}{4}$  as  $\theta_1$  goes from 0 to infinity). All in all, for a high enough value of  $\theta_1$ , (7) must hold.

We round up with two more remarks:

**Elation** This emotion, discussed by Bell (1985) and Loomes & Sugden (1986), is a sort of opposite of disappointment. It can be modeled by substituting  $[\cdot]^-$  for  $[\cdot]^+$  in (4) which then leads to p-games.<sup>31</sup>

Reference point Bell (1985) and Loomes & Sugden (1986) differ from us in the way they define rational choice/planning: We assume that realized payoff is compared to the *ex ante* (before choice) expected payoff, which depends on the agent's pre-determined plan. They instead assume that the term of comparison (reference point) depends on the actual (irreversible) choice. Since this relates to how Kőszegi & Rabin (2006, 2007, 2009) model belief-dependent loss aversion, we postpone the discussion to Section 5.

#### 3.3 Frustration & anger

Psychologists argue that people get frustrated when they are unexpectedly denied things they care about. That sounds like disappointment! However, while disappointment is mainly discussed in regards to pangs incurred and anticipated, frustration is more often discussed for how it influences decision

<sup>&</sup>lt;sup>31</sup>Elation is not discussed nearly as often as disappointment, and seems to be less often regarded as empirically relevant. In line with that, Gill & Prowse (2011) report results indicating "that winners are elated while losers are disappointed, and that disappointment is the stronger emotion" (p. 495).

making going forward. In particular, there is the "frustration-aggression hypothesis," originally proposed by Dollard *et al.* (1939) (see also, e.g., Averill 1982, Berkowitz 1978, 1989, Potegal, Spielberger & Stemmler 2010), whereby frustration breeds aggression towards others. We limit our discussion of frustration to its role in that context, which, we argue, suggests a difference in how to model frustration and disappointment.

Anger and aggression can have profound economic impact, though few economists studied the topic. Battigalli, Dufwenberg & Smith (2019) propose a broadly applicable model. They do not develop applications, but mention pricing, domestic violence, riots, recessions, contracting, arbitration, terrorism, road rage, support for populist politicians, and bank bail-outs as potentially interesting ones.<sup>32</sup> We sketch key features of the approach, and start with an example from the authors— $G_8$ —designed to make a technical point about frustration and how it compares with disappointment:

 $[G_8]$ 

Suppose that if 2 is frustrated she will consider 1 an attractive target of aggression. What would she do if 1 chooses forward? The answer may seem intuitively obvious, but consider what would happen if frustration were modeled as disappointment (more disappointment giving higher inclination to aggression). Building on eq. (4), there would be multiple optimal plans for 2, following the logic of (ii) in Section 3.2. If 2 plans to choose havoc, and if she believes 1 will choose down, then she would be disappointed after forward, hence choose havoc in order to hurt 1.

With outcome (2,2) available, this seems psychologically implausible. Battigalli, Dufwenberg & Smith resolve the issue by requiring players to focus on what has happened and what they can achieve in the future.

Maybe she will be frustrated and end up meting out a costly punishment, but that should be a reaction to, rather than a cause of, her frustration. This consideration leads to the following definition of i's frustration at history h:

$$F_i(h; \alpha_i) = \left[ \mathbb{E}[\pi_i; \alpha_i] - \max_{a_i \in A_i(h)} \mathbb{E}[\pi_i|(h, a_i); \alpha_i] \right]^+, \tag{8}$$

<sup>&</sup>lt;sup>32</sup>As the authors discuss, some of these topics have been analyzed by others empirically or using models that feature anger which however is not modeled using PGT. See, e.g., Rotemberg (2005, 2011) on pricing, Card & Dahl (2011) on family violence, and Passarelli & Tabellini (2017) on political unrest.

where  $\mathbb{E}[\pi_i|(h,a_i);\alpha_i]$  is the material payoff that i expects to get, according to his first-order beliefs  $\alpha_i$ , given history h and action  $a_i$ . Applied to  $G_8$ , let p be the probability 2 initially assigns to forward while q is the probability with which 2 plans to choose bliss (thus,  $\alpha_2$  is described by p and q). We get  $F_2(forward;\alpha_2) = [(1-p)\cdot 1 + p\cdot q\cdot 2 - 2]^+ = 0$ . Zero frustration breeds no aggression, so 2 will choose bliss.

While eq. (8) differs from the disappointment-part of (4), it is still a feature that brings own-plan dependence and belief-dependent p-game utilities. Having defined frustration, the next step is to model how that breeds anger and frustration. We avoid going into technical details—see the article for that—and here just highlight some key themes. Number one is that one must now theorize about blame. Consider  $G_9$  (where players payoffs are listed in alphabetical order, and Don is a dummy player):

 $[G_9]$ 

Battigalli, Dufwenberg & Smith assume that a frustrated player (which in  $G_9$  could only be Penny because (8) must equal 0 at the root) becomes inclined to hurt those deemed blameworthy. They develop three models based on different blame notions. We indicate how they play out in  $G_9$ :

Simple anger: All co-players are blamed independently of how they have behaved.<sup>33</sup> In  $G_9$ , if Penny's anger sensitivity  $\theta_P$  is high enough, she would choose d, going after Don whom she is most efficient at punishing.

Anger from blaming behavior: i's co-players are blamed to the extent that they could have averted i's frustration had they chosen differently. In  $G_9$ , with  $\theta_P$  high, Penny would choose b, going after Ben, since Don is no longer blameworthy (he had no choice!), and Penny is more efficient at beating up Ben than Abe.

Anger from blaming intentions: i's co-players are blamed to the extent that i believes they intended to cause i's frustration. In  $G_9$ , with  $\theta_P$  high,

<sup>&</sup>lt;sup>33</sup>Some psychologists argue that frustrated people tend to be unsophisticated and inclined to blame in such a way; see Marcus-Newhall *et al.* (2000) for a discussion. It seems to us that how and why people blame is an interesting empirical issue, which may depend on, e.g., how tired a person is, or on whether he or she has drunk a lot of beer.

Penny would choose a, going after Abe, since also Ben is no longer blameworthy (while he could have averted Penny's dismay, he had no rational way of correctly figuring out chance's actual choice, and thus can't have had bad intentions). This third category, because Penny cares about others' intentions, injects a second form of belief-dependence in players' utilities.

Finally, a comment about how these models apply to the Mini-ultimatum Game form,  $G_5$ . A comparison with reciprocity theory is of interest, since both approaches can help explain the prevalence of fair offers and rejections. In both cases (anger and reciprocity), 2 may rationally plan to choose reject (if  $\theta_2$  is high enough, and, in the case of anger, if 2's initial belief that 1 will choose fair is strong enough). However, whereas in Dufwenberg & Kirchsteiger's theory it is possible that 1 chooses greedy even if he expects 2 to choose reject (since 1 then views 2 as unkind, and so may want to retaliate), this could never happen in (any of the three versions of) Battigalli, Dufwenberg & Smith's theory. As hinted at in the previous paragraph, at the root a player cannot be frustrated and he must therefore maximize his expected material payoff.

#### 3.4 Regret

Despite Édith Piaf's assertion, regret can be a powerful feeling. To appreciate this, ask anyone who didn't sell stock while the coronavirus was ravaging China and Italy, but had not yet hit the US where stock prices remained close to all-time high. Zeelenberg & Pieters (2007) discuss other examples and synthesize much evidence from psychology.<sup>34</sup>

Research on regret starts with theoretical work by Bell (1982) and Loomes & Sugden (1982), who focus on pairwise choice. Quiggin (1994) proposes an extension for general choice sets. These authors restrict attention to single decision maker settings, but regret makes equal sense with strategic interaction. B&D, BC&D, and Dufwenberg & Lin (2019) formulate relevant definitions. We explain why (unlike in the case with disappointment) PGT is not needed for handling the decision theorists' settings, and why nevertheless PGT is crucial for analyzing general game forms.<sup>35</sup>

<sup>&</sup>lt;sup>34</sup>See also Zeelenberg (1999) and Connolly & Butler (2006).

<sup>&</sup>lt;sup>35</sup>A handful of papers proposed ways, not based on PGT, to incorporate regret in particular games; see, e.g., Engelbrecht-Wiggans (1989), Filiz-Ozbay & Ozbay (2007), and

Consider the following version of Quiggin's approach: Let  $\Omega$  and A be (finite) sets of states (chosen by chance, or nature) and actions of the decision maker (= 1). The payoff function  $\pi_1: \Omega \times A \to \mathbb{R}$  has a finite range  $C \subseteq \mathbb{R}$  of monetary consequences. Function  $v_1: C \to \mathbb{R}$  describe 1's "choiceless utility" (Loomes & Sugden's terminology) of consequences. However, after 1 chooses  $a \in A$ , chance's choice  $\omega \in \Omega$  is revealed and 1 now ruminates on what could have been. His regret-adjusted utility, which is what he wants to maximize, is a function  $u_1: \Omega \times A \to \mathbb{R}$  defined by

$$u_1(\omega, a) = v_1(\pi_1(\omega, a)) - f(\max_{a' \in A} v_1(\pi_1(\omega, a')) - v_1(\pi_1(\omega, a))),$$
 (9)

where  $f: \mathbb{R}_+ \to \mathbb{R}_+$  is strictly increasing. In words, i's overall utility involves pangs of regret that reflect i's evaluation of how much better off he could have been had he chosen differently. For our purposes it is useful to re-formulate this as a one-player game form with a chance-move, with perfect information at end-nodes: Chance makes a choice from  $\Omega$ . Player 1 is not informed of chance's choice, and chooses  $a \in A$ . Then end-node (terminal history)  $(\omega, a)$  is reached and revealed to 1, whose utility is computed using (9). Note that this is a traditional game, because 1's utility is uniquely determined at each end-node.

However, if one generalizes the above steps to apply to any game form, then one arrives at a p-game: To see this, fix an extensive game form, focus on player i, and try to compute his regret-adjusted utility at end-node z (and at the associated terminal information set; here is one more instance where terminal information will influence the analysis). To do that, one needs to figure out what would have happened had i chosen differently. This, in turn, depends both on what choices i's co-players actually made, and which ones they would have made at any history in the game tree that i could have made play reach had he chosen differently than he did. And that computation, of course, will reach a different answer dependent on which choices the coplayers are assumed to make. In contrast to the single-player example of the previous paragraph, i's regret-adjusted utility will not be uniquely defined. If i regret-adjusts based on his beliefs about what would have happened had he chosen differently, we would get a p-game. The belief-dependence of i's utility involves his own beliefs at end-nodes (and associated information sets) regarding co-players' choices.

Halpern & Pass (2014).

For example, consider  $G_4$ . Would 1 experience regret if he chose *stay*, and if so how much? The answer depends on p, the probability with which 1 believes that 2 would choose *take* had 1 chosen *reach*. Analogous remarks apply to, e.g.,  $G_5$ ,  $G_6$ , and  $G_8$ .

#### 3.5 Anticipatory feelings

So far we considered either the action tendencies caused by emotions, as in the frustration/aggression hypothesis, or how actions cause emotions (own or of others) with positive or negative valence and how players take this into account in their choice, as with guilt, disappointment or regret. Behavior of the second kind is explained by the anticipation of future feelings under different courses of action. Now we consider how uncertainty about the future can cause "anticipatory feelings" with negative or positive valence in the present (cf. Loewenstein, Hsee, Weber, & Welch 2001). Of course, the anticipation of such anticipatory feelings can drive behavior in earlier periods. Timing is essential to model anticipatory feelings. The simplest setting for a meaningful discussion is one with two periods  $t \in \{1, 2\}$  between three dates 0, 1, and 2. Each period t is between dates t - 1 and t. Action profile  $a^t$  is selected in period t. To make the problem interesting, player t—the decision maker under consideration—has to be active in period 1 and another player (typically, chance) has to be active in period 2.

Anxiety is an anticipatory feeling with negative valence caused by uncertainty about future material outcomes (e.g., health, or consumption). Huang (2020) argues that anxiety has major welfare consequences during pandemics. Drawing on earlier work by Kreps & Porteus (1978) on preferences for the temporal resolution of uncertainty, Caplin & Leahy (2001) put forward an axiomatic model of utility of "temporal lotteries" and consider specific functional forms. As one example, they analyze portfolio choice. Using our notation, they consider the following utility

$$u_i\left((a^1, a^2), \alpha_i\right) = -\left(\theta_i^V \mathbb{V}\left[\pi_i | a^1; \alpha_i\right] - \theta_i^E \mathbb{E}\left[\pi_i | a^1; \alpha_i\right]\right) + v_i^2\left(\pi_i\left(a^1, a^2\right)\right), \tag{10}$$

where  $\mathbb{V}$  is the variance operator,  $\theta_i^V$  and  $\theta_i^E \geq 0$  are sensitivity parameters, and  $v_i^2$  is the period-2 utility of the realized material outcome. The higher is the variance exhibited by i's beliefs about his material payoff, the lower is his utility. The theory helps explain the risk-free rate puzzle and the equity-premium puzzle: when buying safe assets an agent is "paying for his peace

of mind".

Caplin & Leahy (2001) also briefly mention how their general theory can be adapted to model suspense, i.e., the pleasure experienced immediately prior to the anticipated resolution of uncertainty. This theme is explored in depth by Ely, Frankel, & Kamenica (2015). Finally, Caplin & Leahy (2004) draw on their (2001) theory to study interaction between, e.g., an anxious patient and his caring doctor, who decides whether or not to reveal information affecting the patient's anticipatory feelings.

#### 3.6 Wrap-up remarks on emotions

We round up this section by collecting three remarks on distinct topics:

Valence and action-tendency Emotions have many characteristics, two important ones being valence, meaning the (material or psychological) costs or rewards associated with an emotion, and action-tendency, or how an emotion's occurrence incites new behavior. When modeling emotions using PGT one may want to choose which aspect to highlight, or abstract from. For example, Battigalli & Dufwenberg's (2007) models of guilt (cf. Section 3.1) are all about valence, abstracting away from action-tendency. This could well be restrictive; see, e.g., Silfver (2007) for a discussion of "repair behavior," which could be thought of in terms of an action-tendency of guilt. Similar remarks apply concerning the approaches to disappointment and regret presented in Sections 3.2 and 3.4. By contrast, Battigalli, Dufwenberg & Smith's models of frustration and anger are all about action-tendency, as frustration has no valence in their models. Again, this may be a restrictive abstraction. Frustration may, e.g., plausibly have similar valence as disappointment).

Is reciprocity an emotion? A referee raised the issue. Judging by similarity of mathematical styles (in Sections 2 and 3), perhaps the answer could be yes. However, scholars working on reciprocity rarely describe what they model as involving emotions, and reciprocity usually does not figure in the list of emotions. We have chosen to structure our presentation accordingly.

Elster's and Keltner & Lerner's lists While we have covered several emotions, and highlighted their connections with PGT, we have not been

exhaustive. Elster (1998) discusses anger, hatred, guilt, shame, pride, admiration, regret, rejoicing, disappointment, elation, fear, hope, joy, grief, envy, malice, indignation, jealousy, surprise, boredom, sexual desire, enjoyment, worry, and frustration. Keltner & Lerner (2010) offer another list (p. 330), which overlaps to a large degree but also adds contempt, disgust, embarrassment, contentment, enthusiasm, love, compassion, gratitude, awe, interest, amusement, and relief. We suspect that many of the additional sentiments listed here involve belief-dependent motivation that could be explored using PGT. However, rather than pursue these topics we propose that they hold promise for rewarding research to come.

#### 4 Image concerns

Introspection and empirical and experimental evidence suggest that people are willing to give up some material payoffs to improve the opinion of others about them. Evidence about deception can be explained by a trade-off between monetary payoff and a reduction of the perceived extent of cheating or lying (see the example below). Other models instead assume that agents try to signal that they have "good traits" such as being altruistic or fair (e.g., Bénabou & Tirole 2006; Andreoni & Bernheim 2009; Ellingsen & Johannesson 2008; Grossman & van der Weele 2017), which may explain behavior in the Dictator Game, or why people seldom give anonymously to charities, while they are happy to give non-anonymously (as shown by Glazer & Konrad 1996). Several other articles explore various forms of image concerns explaining, e.g., conformity, job-seeking effort, randomized survey-response, shame avoidance, peer evaluations, and pricing distortions.<sup>36</sup>

The aforementioned examples suggest two broad kinds of image about which people are concerned: others' (terminal) beliefs about (i) imperfectly observed bad/good actions, and (ii) imperfectly observed bad/good traits. Both are modeled by psychological utility functions.

<sup>&</sup>lt;sup>36</sup>See Bernheim (1994), Dufwenberg & Lundholm (2001), Blume, Lai & Lim (2019), Tadelis (2011), and Sebald & Vikander (2019). We note that some of the cited models of image concern do not make the PGT-connection explicit.

#### 4.1 Opinions about bad/good actions

Play in a game form is represented by a terminal history, or path  $z \in Z$  of actions taken by the players (including chance, when relevant). Suppose for simplicity that, according to some standard, paths in  $Z_i^B$  (resp.  $Z_i^G$ ) are such that player i behaved in a bad (resp. good) way. Some paths may be neutral, e.g., because i did not play. For example, in a deception game form  $Z_i^B$  could be the set of paths where i lies; in a Trust Minigame (e.g.,  $G_6$  above and  $G_{11}$  below) i is the trustee and  $Z_i^B$  (resp.  $Z_i^G$ ) contains the paths where he grabs (resp. shares).<sup>37</sup> Let j be an observer who obtains possibly imperfect information about the realized path z, and let  $p_{j,i}^B(z;\alpha_j)$  (resp.  $p_{j,i}^G(z;\alpha_j)$ ) denote the observer's ex post probability of bad (resp. good) deeds conditional on what he observed, given z and j's system of beliefs about actions  $\alpha_j$ . An image concern related to bad/good deeds can be captured by a simple functional form like

$$u_i(z,\alpha_j) = \pi_i(z) + \theta_i \left[ p_{i,i}^G(z;\alpha_j) - p_{j,i}^B(z;\alpha_j) \right]. \tag{11}$$

More generally, one can assume that intrinsic motivations—besides image concerns—also play a role (i (dis)likes good (bad) deeds as in Gneezy, Kajackaite & Sobel 2018 and Khalmetski & Sliwka 2019), or that i cares about the perceived distance from the standard rather than mere compliance. We expand on the second theme presenting the model of cheating by Dufwenberg & Dufwenberg (2018), which is a useful illustrative example.

Perceived cheating aversion A large recent literature explores humans' reluctance to lie or cheat using an experimental "die-roll paradigm" introduced by Fischbacher & Föllmi-Heusi (2013).<sup>38</sup> Dufwenberg & Dufwenberg (2018) propose a PGT-based account of this behavior. We draw on their work to illustrate how a concern with others' opinions regarding chosen actions can be modeled.

A subject is asked to roll a six-sided die in private and to report the outcome, but the report is non-verifiable and can be submitted with impunity. The subject is paid in proportion to the reported number, with one exception: reporting six yields a payout of zero. We will refer to a six as a "zero".

<sup>&</sup>lt;sup>37</sup>Note that paths record the behavior of every active player, hence we can accommodate norms such as behaving (or not) like the majority.

<sup>&</sup>lt;sup>38</sup>See Abeler, Nosenzo & Raymond (2019) for a survey.

Formally, chance (player 0) draws  $x \in \{0, ..., 5\}$  from a uniform distribution (x = 0 corresponding to rolling a six). Player 1 observes x and then chooses a report  $y \in \{0, ..., 5\}$  after which he is paid y.<sup>39</sup> Choice y, but not realization x, is observed by player 2, who is an "audience". In applications the audience might be a fellow citizen, but in the lab it could be the experimenter or an observer "imagined" by player 1. Player 2 has no (active) choice, but forms beliefs about x after observing y. The associated game form is  $G_{10}$ :

$$[G_{10}]$$

The analysis will not depend on 2's payoffs, which are therefore not specified. The dotted lines depict *information sets across end-nodes*. This is a feature rarely made explicit in traditional game-theoretic analysis, but here it will be critical. In  $G_{10}$ , these sets reflect player 2's end-of-play information.

Consider the following preference: Player 1 feels bad to the extent that player 2 believes that 1 cheats. Measure actual cheating at end-node (x, y) as  $[y - x]^+$ , i.e., cheating involves reporting a higher number than the roll and downward lies do not count as cheating. Player 2 cannot observe x, but draws inferences about x conditional on y. Let  $\alpha_2(x'|y) \in [0,1]$  be the probability 2 assigns to chance event x = x' given report y, with  $\sum_{x'} \alpha_2(x'|y) = 1$ , so 2's expectation of 1's cheating equals  $\sum_{x'} \alpha_2(x'|y)[y - x']^+$ . Player 1's utility of (x,y) given  $\alpha_2$  is

$$u_1((x,y),\alpha_2) = y - \theta_1 \cdot \sum_{x'} \alpha_2(x'|y)[y-x']^+,$$
 (12)

where  $\theta_1 \geq 0$  measures 1's sensitivity to 2's expectation of 1's cheating. Note that (12) is independent of x. This reflects the fact that 1 cares about his image, not about cheating *per se*. Also, 1 may feel bad even if he does not lie, if the audience believes that he cheats.

Appending utility function (12) to game form  $G_{10}$ , we obtain a p-game because  $\alpha_2(x'|y)$  is an endogenous belief, i.e., it has to be derived by strategic reasoning. Adopting the traditional equilibrium approach, the strategic analysis of this p-game is tractable and delivers testable predictions. In this game form, the relevant beliefs of player 1 about actions describe 1's plan (or behavior strategy), so  $\alpha_1(y|x)$  is the probability that  $\alpha_1$  assigns to y

<sup>&</sup>lt;sup>39</sup>That is, the monetary payoff function is  $\pi_1(x,y) = y$ .

after 1 observes x.<sup>40</sup> Dufwenberg & Dufwenberg solve for equilibria such that  $\alpha_1$  maximizes (12) given 2's beliefs, and  $\alpha_2(x'|y)$  is computed as a conditional probability using correct initial beliefs, that is,  $\alpha_1(x) = 1/6$  and  $\alpha_2(y|x) = \alpha_1(y|x)$ .<sup>41</sup> It can be shown that an equilibrium always exists. However, if 1's concern for his image is strong enough  $(\theta_1 > 2)$ , neither honesty  $(\alpha_1(x|x) = 1$  for all x) nor selfish choice  $(\alpha_1(5|x) = 1$  for all x) is an equilibrium. The striking implication: if  $\theta_1 > 2$  then equilibrium play involves partial lies (in expectation).

Walking through a sketch of the proof is helpful to get intuition for why this result holds: If honesty were expected by 2 then  $\alpha_2(x|x) = 1$  for all x, so cheating by 1 to y = 5 > x would raise no suspicion, hence be 1's best response, ruling out an honest equilibrium (for any value of  $\theta_1 \ge 0$ ). If selfish play  $(\alpha_1(5|x) = 1$  for all x) were expected then 2's expectation of 1's cheating would equal  $\sum_x \frac{1}{6}[5-x]^+ = 2.5$ ; if  $\theta_1 > 2$  player 1 could then increase his utility by deviating to y = 0 (so that perceived cheating = 0).

The analysis just conducted depends critically on the information across the end-nodes. To see this, consider what would happen if those informations sets were split into singletons. That is, assume that 2 is told about both x and y, i.e., which path (x,y) occurred. At (x,y), player 2 would form beliefs such that  $\alpha_2(x|y) = 1$ , implying that perceived and actual cheating coincide. If  $\theta_1 > 1$  then 1's choices would be honest  $(\alpha_1(x|x) = 1$  for all x); if  $\theta_1 < 1$  then 1's choices would be selfish  $(\alpha_1(5|x) = 1$  for all x). The partiallies prediction evaporates. This illustrates a feature, reminiscent also of the earlier comparison of  $G_3^*$  and  $G_3^{**}$ , that is unique to p-games. In traditional game theory, utilities are not affected by information across end-nodes, which therefore has no impact on the strategic analysis.<sup>42</sup>

 $<sup>^{40}</sup>$ Player 1's initial beliefs about chance moves — exogenously given by the uniform distribution — are irrelevant because he chooses after observing x. However, 1 takes into account that 2 knows the chance probabilities.

<sup>&</sup>lt;sup>41</sup>Formally, (i)  $\alpha_1(y|x) > 0 \Rightarrow y \in \arg\max_{y'}(y' - \theta \cdot \sum_{x'} \alpha_2(x'|y')[y' - x']^+)$  and (ii)  $\sum_x \alpha_1(y|x) > 0 \Rightarrow \alpha_2(x'|y) = \frac{\alpha_1(y|x')}{\sum_x \alpha_1(y|x)}$ . Our use of equilibrium analysis in this example does not mean that we endorse it in general; see the last part of Section 9.

<sup>&</sup>lt;sup>42</sup>This explains why in traditional game-theoretic analysis information sets over endnodes are usually not drawn, even if such information is objectively determined by the rules of interaction.

#### 4.2 Opinions about bad/good traits

The second kind of image concern starts from intrinsic motivation. People have heterogeneous intrinsic motivations to do good deeds and avoid bad ones, and are imperfectly informed about the motivations of others. This expands the domain of uncertainty: now we have to consider systems of first-order beliefs about both (paths of) actions and traits. Suppose, just for the sake of simplicity, that actions are perfectly monitored  $ex\ post$ . Then, after the realization of any path of play (terminal history) z, each player j holds an endogenous conditional belief  $\alpha_j(\cdot|z)$  about the traits of others  $\theta_{-j}$ . Intrinsic motivation of i is measured by parameter  $\theta_i^{\mathbf{I}} \geq 0$ , and i—besides liking material payoff and being intrinsically motivated—also cares about his **reputation**, that is, j's  $ex\ post$  estimate of  $\theta_i^{\mathbf{I}}$ . For example, i's psychological utility could be

$$u_{i}(z, \alpha_{j}, \theta_{i}) = \pi_{i}(z) + \theta_{i}^{\mathbf{I}} \left[ \mathbf{I}_{i}^{G}(z) - \mathbf{I}_{i}^{B}(z) \right] + \theta_{i}^{\mathbf{R}} \mathbb{E} \left[ \widetilde{\theta_{i}^{\mathbf{I}}} | z; \alpha_{j} \right], \qquad (13)$$

where  $\theta_i = (\theta_i^{\mathbf{I}}, \theta_i^{\mathbf{R}})$  is *i*'s trait vector, and  $[\mathbf{I}_i^G(z) - \mathbf{I}_i^B(z)]$  denotes the net intensity of *i*'s good deeds in path z.<sup>44</sup> More generally, *j*'s ex post belief about  $\theta_i^{\mathbf{I}}$  may be conditional on possibly imperfect information about the realized path of play. This allows for comparing non-anonymous and anonymous donations, or to consider the possibility that *i* has imperfect recall and is his own observer (j = i), as in the work Bénabou & Tirole (2002, 2006, 2011).<sup>45</sup>

Utility functions like (13) introduce a familiar element of signaling into the strategic analysis: even if i's intrinsic motivation to do good  $(\theta_i^{\mathbf{I}})$  is low, he may be willing to pay a material cost to make j believe that  $\theta_i^{\mathbf{I}}$  is high, hence that i is a "good guy". The simplest models of this kind are signaling games where only the sender is active and the receiver is a mere observer.<sup>46</sup>

<sup>&</sup>lt;sup>43</sup>Formally, we are considering beliefs in games with incomplete information; see Section 5 below (under the heading "Private sensitivities") for a discussion.

<sup>&</sup>lt;sup>44</sup>In the binary case considered at the beginning of this section,  $\mathbf{I}_{i}^{D}(\cdot)$  is the indicator function of  $Z_{i}^{D}$ , the set of paths where i made good (D = G) or bad (D = B) deeds.

<sup>&</sup>lt;sup>45</sup>See below. Note that Bénabou & Tirole (2006) consider a model similar to the one in the main text  $(j \neq i)$ , but also put forward a reinterpretation with imperfect recall where j is a future self of i.

<sup>&</sup>lt;sup>46</sup>Readers may fail to recognize that Bénabou & Tirole's (2006) model is a signaling game, because they choose not to frame it explicitly as such, making it seem more like a decision problem. But they indirectly hint (in footnote 17) at the fact that they are considering a refinement of signaling equilibria.

A noteworthy application of this approach concerns privacy. Depending on available technology and regulation, what we do may be monitored even when it does not affect the material payoff of anybody else. Many people seem to care about this and in Western countries there is a consensus that privacy should be protected. Gradwohl & Smorodinsky (2017) model this by considering functional forms such that (as in eq. 13), for each  $(z, \theta_i)$ ,  $u_i(z, \alpha_i, \theta_i)$ depends on the ex post belief of the "audience" j about  $\theta_i$ . In particular they assume that—other things being equal—the agent either wants j's posterior to be the same as the prior, 47 or dislikes being identified. 48 Focusing on the simple case where i is the only active agent and actions are observable by j (lack of privacy), they analyze the pooling and separating (Bayesian perfect) equilibria of the resulting signaling game. Pooling distorts actions from the first best that would obtain under perfect privacy. Separation inflicts a psychological utility loss due to identification.

Another application concerns identity. According to Hupkau & Maniquet (2018), one's own personal traits are part of an agent's identity, and he can suffer from the discrepancy between his true identity and others' perceptions of it.<sup>49</sup> This in turn is affected by actions *via* signaling, possibly causing inefficiencies. For example, high-type agents may refrain from requesting useful service from a provider to avoid being pooled with low-type agents. Bénabou & Tirole (2011) instead consider forgetful agents who "care about who they are". Their actions depend on such self-perception, which may be forgotten later on. Thus, actions are also identity signals for the future self, introducing an identity-investment concern in the choice of the current self.<sup>50</sup>

#### More motivations 5

The previous three sections focused on the three categories of motivation mentioned in the Introduction. We now complement that presentation by

<sup>&</sup>lt;sup>47</sup>That is,  $\alpha_{j,\Theta_i}^0 \in \arg\max_{\mu_j \in \Delta(\Theta_i)} u_i(z,\mu,\theta_i)$  for each  $(z,\theta_i)$ , where  $\alpha_{j,\Theta_i}^0 = \max_{\Theta_i} \alpha_j (\cdot | \varnothing)$  is the exogenous and known prior of j about  $\theta_i$ .

<sup>48</sup>That is,  $\delta_{\theta_i} \in \arg\max_{\mu_j \in \Delta(\Theta_i)} u_i(z,\mu_j,\theta_i)$  for each  $(z,\theta_i)$ , where  $\delta_{\theta_i}$  is the degenerate

<sup>(</sup>Dirac) marginal belief that assigns probability 1 to  $\theta_i$ .

<sup>&</sup>lt;sup>49</sup>In this case, the third element in eq. (13) is replaced by  $\ell\left(\left|\theta_i^{\mathbf{I}} - \mathbb{E}\left[\widetilde{\theta_i^{\mathbf{I}}}|z;\alpha_j\right]\right|\right)$ , where  $\ell: \mathbb{R}_+ \to \mathbb{R}_+$  is an increasing (e.g., quadratic) loss function.

<sup>&</sup>lt;sup>50</sup>Bénabou & Tirole's model is a p-game, because utility depends on an endogenuos belief, but—unlike Hupkau & Maniquet—they do not make the link to PGT explicit.

discussing additional forms of belief-dependent motivations, and broader related issues.

**Opposites** Sometimes a meaningful belief-dependent motivation takes an "opposite" form of another sentiment, switching a sign or replacing  $[\cdot]^+$  with  $[\cdot]^-$  in the utility formula. We already saw examples in Sections 3.2 and 3.5, where, respectively, elation was compared to disappointment, and suspense was related to anxiety.

Another example involves an opposite to guilt. Ruffle (1997) and Khalmetski, Ockenfels & Werner (2015) consider situations in which player i enjoys surprising j, in particular so that j gets a higher material payoff than j expected. See also Dhami, Wei & al-Nowaihi (2019). This can be modeled by substituting  $[\cdot]^-$  for  $[\cdot]^+$  in (3). Is surprising others this way another "emotion"? Maybe yes; obviously it is a kind of joy, and joy is often listed as an emotion.

The desire to surprise actually has venerable PGT-ancestry. GP&S explored the idea in their verbally presented opening example, although a different variety than the work cited above. GP&S's example does not require that the co-player is surprised in terms of material payoff.<sup>51</sup> Here is the quote (from p. 62), illustrating the sentiment and a feature idiosyncratic to p-games:

Think of a two-person game in which only player 1 moves. Player 1 has two options: she can send player 2 flowers, or she can send chocolates. She knows that 2 likes either gift, but she enjoys surprising him. Consequently, if she thinks player 2 is expecting flowers (or that he thinks flowers more likely than chocolates), she sends chocolates, and vice versa. No equilibrium in pure strategies exists. In the unique mixed strategy equilibrium, player 1 sends each gift with equal probability. Note that in a traditional finite game with only one active player, there is always a pure strategy Nash equilibrium. That this is untrue in psychological games demonstrates the impossibility of analyzing such situations merely by modifying the payoffs associated with various outcomes: any modification will yield a game with at least one pure strategy equilibrium.

<sup>&</sup>lt;sup>51</sup>Yet another example appears in Geanakoplos (1996), which reconsiders the classical "hangman's paradox" from philosophy, where the desire to surprise has a sadistic flavor.

Belief-dependent loss aversion When we discussed disappointment, in Section 3, we mentioned how that sentiments is closely related to ideas explored by Kőszegi & Rabin (2006, 2007, 2009) and by Shalev. The goal of these authors, however, is not to model disappointment, but rather to tie in with Kahneman & Tversky's (1979) work on prospect theory. Kőszegi & Rabin model prospect theory's central notion of a "reference level" as a decision maker's initially expected outcome. When he gets less than he expects he experiences loss, effectively much like in disappointment theory. Kőszegi & Rabin allow for losses in many dimensions, e.g., in n + 1 dimension if there are n goods as well as money. To capture that, we would have to augment the framework in Section 6, allowing i's outcome function  $\pi_i$  to be vector-valued.

The key features we highlighted in regards to disappointment in Section 3.2 have counterparts in the work of Kőszegi & Rabin. Most notably the feature of own-plan dependent utility is there, and it may lead to multiplicity of non-equivalent rational plans, as well as non-existence of degenerate rational plans. Beyond those technical similarities, details differ quite a lot. The exact way in which Kőszegi & Rabin define belief-dependent loss is different from the way that Bell (1985) and Loomes & Sugden (1986) (and we) define disappointment, and they also consider more notions of rational planning than we did when we discussed disappointment. The recent and penetrating survey on "Reference-Dependent Preferences" by O'Donoghue & Sprenger (2018) discusses all of these aspect in depth, so we refer to their text (and especially their Sections 5-7) for further details. Here we only mention one aspect. The notion of rational plan that we illustrated in Section 3.2 corresponds to Kőszegi & Rabin's concept of "personal equilibrium." They also consider (i) a refinement, "preferred personal equilibrium," that—in case of non-trivial multiplicity—selects the personal equilibrium most favorable to the initial self, and (ii) another concept, "choice-acclimating personal equilibrium," whereby the referent to which realized outcomes are compared is determined by the actual (irreversible) choice of the agent rather than his ex ante plan, i.e., it is the expected outcome conditional on the agent's action.

Consistently with Kőszegi & Rabin's explicit reference to different time frames, we can accommodate such distinctions in the PGT framework by explicitly introducing time periods, which in turn may comprise multiple stages (see Section 6 and BC&D): Endgame utility is the sum of the utilities of different periods, the referent for one-period gain-loss utility is determined by beginning-of-period beliefs. With this, personal equilibrium refers to rational planning in one-period decision problems, whereas choice-acclimating

equilibrium applies to two-period situations where the agent chooses in the first period and uncertainty realizes in the second period. Except for differences concerning the exact definition of the referent, the approach of Bell (1985) and Loomes & Sugden (1986) fits the choice-acclimating equilibrium concept.

As regards applied work, Kőszegi & Rabin discuss consumption, risk-preferences, and savings. O'Donoghue & Sprenger discuss other papers that deal with endowment effects, labor supply, job search, pricing, and mechanism design.

**Self-esteem** Self-esteem reflects an individual's overall subjective emotional evaluation of his own worth. It is "the positive or negative evaluations of the self" and "how we feel about it" (Smith & Mackie, 2007). We can model self-esteem by assuming that a valuable personal trait  $\theta_{0,i}$  of player i is imperfectly known by i. Such trait could be general intelligence, or ability. Player i's utility is increasing in his ex post estimate of  $\theta_{0,i}$  conditional on the path of play z,  $^{52}$  as in function

$$u_{i}(z, \alpha, \theta) = \pi_{i}(z, \theta) + v_{i}^{e}\left(\mathbb{E}\left[\widetilde{\theta_{0,i}}|z; \alpha_{i}\right]\right), \tag{14}$$

where the "ego-utility"  $v_i^e$  is increasing, and we allow material payoff  $\pi_i$  to depend on parameter vector  $\theta$  because traits such as ability typically affect material outcomes. For example, Mannahan (2019) shows that if  $\pi_i$  is observed ex post and  $v_i^e$  is concave, i may decide to handicap himself ensuring a bad outcome (e.g., by not sleeping before an exam) rather exposing himself to the risk of discovering that his ability is low.<sup>53</sup>

Also, better informed players may engage in signaling to affect *i*'s self-esteem: Does a teacher want to reveal to a student how bad his performance was? Better information may allow for a better allocation of the student's time (more study, less leisure), but it may also be detrimental: by decreasing the student's estimate of his ability it can bring it in a range where ego-utility is more concave and cause the self-handicapping effect described above.

A few economic studies of self-esteem model utility in line with our description here, although (unlike Mannahan) they do not make the PGT-

<sup>&</sup>lt;sup>52</sup>Path z may include a randomly chosen output, whose distribution depends on  $\theta_{0,i}$ .

<sup>&</sup>lt;sup>53</sup>There is a discussion in psychology of similar self-handicapping strategies, with implications regarding for example drug use. Berglas & Jones (1978) is a classic experimental study on this topic.

connection explicit. See Kőszegi (2006), Eil & Rao (2011), Möbius, Niederle, Neihaus & Rosenblat (2011), Sebald & Walzl (2015), and Kőszegi, Loewenstein & Murooka (2019).<sup>54</sup>

Higher-order belief-dependence The framework presented in Section 6 restricts the domain of a player's utility to depend on beliefs (own and others') up to only the first order.<sup>55</sup> This is enough to handle almost all forms of motivation that to date have been modeled using PGT.<sup>56</sup> The main exception is Battigalli & Dufwenberg's (2007) model of guilt-from-blame (but see also B&D, p. 14). We now indicate how that sentiment works in an example designed to provide a contrast with simple guilt (as presented in Section 3). Guilt-from-blame plugs a third-order belief into the domain of a player's utility, so we leave the framework of Section 6. We sketch the approach without going into all formal details:

First, for each end-node z in a game, measure how disappointed j is as  $[\mathbb{E}[\pi_j; \alpha_j] - \pi_j(z)]^+$  (compare (3) & (4)). Second, calculate how much of  $[\mathbb{E}[\pi_j; \alpha_j] - \pi_j(z)]^+$  could have been averted had i chosen differently. Third, calculate i's initial belief regarding  $[\mathbb{E}[\pi_j; \alpha_j] - \pi_j(z)]^+$ . Fourth, for each z, calculate j's belief regarding i's initial belief regarding  $[\mathbb{E}[\pi_j; \alpha_j] - \pi_j(z)]^+$ ; this is how much j would blame i if j knew he were at z. Finally, i suffers guilt-from-blame in proportion to j's blame, and i's utility trades off avoidance of that pang against i's material payoff.

Battigalli & Dufwenberg (2007; see Observation 1) prove that simple guilt and guilt-from-blame sometimes have similar implications. However, this is not true in general. To illustrate, consider  $G_{11}$ , a modified version of  $G_6$  in which even if 2 chooses *share* there is a  $\frac{1}{6}$  probability that 1 gets a material payoff of 0. Moreover, if 1 gets 0 then 1 is not informed of 2's choice. As in

<sup>&</sup>lt;sup>54</sup>The aforementioned work of Bénabou & Tirole may be interpreted as modeling self-esteem, but their approach relies on imperfect recall and self-signaling.

 $<sup>^{55}</sup>$ Player i may still have to consider his second-order beliefs, if his utility depends on j's first-order beliefs (as it did in our presentation of reciprocity, guilt, anger from blaming intentions, and image concerns). Since i does not know j's beliefs, he has to form beliefs about them to calcuate a best response.

<sup>&</sup>lt;sup>56</sup>This includes reciprocity, if formulated as in Section 2. (As we noted in a footnote there, others use a different formulation with utilities that depend on second-order beliefs.)

Section 3, 1 believes that there is probability q that 2 would choose share.

 $[G_{11}]$ 

Everything we said about simple guilt and (3) in Section 3.1 we could have formulated vis-à-vis  $G_{11}$  rather than  $G_6$ . We used  $G_6$  merely because it is more spare, but Charness & Dufwenberg (2006) (cited in Section 3.1) actually used  $G_{11}$  rather than  $G_6$ .<sup>57</sup>

If player 2 is instead motivated by guilt-from-blame then the implications are very different in  $G_6$  and  $G_{11}$ . In  $G_6$ , following trust, if player 2's second-order beliefs assign probability 1 to q=1, then for a high enough  $\theta_2$  player 2's best response is share. This is true just as it would be also under simple guilt. In  $G_{11}$ , however, following trust, if player 2's second-order beliefs assign probability 1 to q=1, then player 2's best response is grab regardless of how high  $\theta_2$  is! To appreciate why, note that if 2 believes that q=1 then 2 believes that 1 will not blame 2 if 2 chooses grab. Therefore, 2 can grab with impunity.<sup>58</sup>

Game form  $G_{11}$ , with guilt-from-blame appended to it, joins models of image concern such as  $G_{10}$  in illustrating the critical role information across end-nodes can play in p-games. Modify  $G_{11}$  such that 2's doubleton information set is broken up into two singletons. That is, if 1 gets 0 then 1 is informed of 2's choice.<sup>59</sup> The logic of the previous paragraph no longer applies, and in the modified version of  $G_{11}$  guilt-from-blame and simple guilt again work similarly.

**Social norms** Fehr & Schurtenberger (2018, p. 458) define social norms as "commonly known standards of behavior that are based on widely shared views of how individual group members ought to behave in a given situation." Similar ideas are discussed by Elster (1989), Bicchieri (2006), Andrighetto, Grieco & Tummolini (2015), and Cartwright (2019).

D'Adda, Dufwenberg, Passarelli & Tabellini (2020) develop a model for a restrictive context (a form of Dictator Game) where the central notions

<sup>&</sup>lt;sup>57</sup>Their reason is conceptual; from a contract-theoretic viewpoint  $G_{11}$  may be seen to incorporate an element of "moral hazard" which is absent in  $G_6$ . See Charness & Dufwenberg (2006, p. 1582).

<sup>&</sup>lt;sup>58</sup>The logic is similar to that we illustrated in regards to  $G_{10}$  in Section 4.1.

<sup>&</sup>lt;sup>59</sup>Tadelis compares behavior in experimental treatments that resemble  $G_{11}$  as well as the variation that we are describing here.

concern a player's conception of "the right thing to do" and a proclivity to do what others believe is the right thing to do, especially if there is consensus about this (which would then be an ideal case of a social norm). Departing from the social norm entails an element of disappointing the expectations of others, and the authors explore the idea that decision makers are averse to doing so. In this regard, the motivation resembles guilt, as modeled in Section 3.1. However, d'Adda et al. consider players' expectations regarding how one ought to behave, rather than regarding how one will actually behave. This marks a way that the approach is not formally captured by p-games, as we have described them in this paper.

Many scholars have written papers about social norms, but few proposed formal models, in particular models that can be generally applied.<sup>60</sup> There is work to do in this arena, and we suggest that it should involve (some possibly extended version of) PGT.<sup>61</sup>

**Punishing transgressors** Several motivations that we discussed incorporate some form of desire to punish those who, somehow, "misbehave." Negative reciprocity (Section 2) and anger (Section 3.3) have such features built in, and other notions may quite naturally be extended in that direction. For example, Sebald & Walzl (2015) explore the idea that player i may wish to be unkind (as in reciprocity theory) to j if j produced information (e.g., a performance review) that reduced i's self-esteem. Another example could be if someone is motivated to punish those who violate a social norm.

A recent paper by Molnar, Chaudhry & Loewenstein (2020) highlights an intriguing aspect of such punishment, one that would add a novel feature of belief-dependent motivation (addressable via suitable modifications of each of the motivations we just mentioned). Namely, the punisher may care about whether transgressors are able to figure out that they are being sanctioned. The PGT-connection would be that the punisher's utility depends on the transgressor's belief regarding the punishers choice (and maybe also his intention).

<sup>&</sup>lt;sup>60</sup>López-Pérez (2008) is an important exception. His model is not PGT-based however.

<sup>&</sup>lt;sup>61</sup>We do not expect the topic to be easy to address. There are many subtle issues. Is a norm a strategy or a strategy profile (or, possibly, a set thereof)? If people like to follow norms, what exactly is the nature of the preference involved? Is the cost of breaking a norm dependent on whether and how many others do so?

Emotion carriers In most of the models we discussed, the belief-dependent part of a player's utility was built up with reference to particular material payoffs. For example, following Battigalli & Dufwenberg (2007), player 2's guilt in  $G_4$  has the dimension of (expected) material payoff of player 1. And in Battigalli, Dufwenberg & Smith's model, player i's frustration has the dimension of (expected) material payoff of i. This is not a necessary feature of belief-based utility, and alternatives have been considered. Attanasi, Rimbaud, Villeval (2019) consider "situations where donors need intermediaries to transfer their donations to recipients and where donations can be embezzled before they reach the recipients". They discuss how intermediaries may experience guilt if they do not meet the owner's expectation, although the associated material cost would be incurred by the recipient rather than the donor. And Battigalli, Dufwenberg & Smith (in their Section 5) mention how in principle frustration may depend on regret of a previous decision, unexpected perceived unfairness, or negative shocks to self-esteem.

Unawareness Almost all game-theoretic analysis assumes that the game form is commonly known between the players. Casual observations of reality suggest that this assumption may be too strong. It seems plausible that before the "Fosbury Flop" and the "V-style" were popularized by, respectively, Dick Fosbury and Jan Boklöv, many high jumpers and ski jumpers were not aware of these techniques, or at least of their lucrative payoff consequences. Or, awareness of the possibility and expected nature of hi-jackings may have been altered by the 9/11 events.

There are formal models of unawareness—see Fagin & Halpern (1988) for a pioneering effort—and recent literature develops related techniques for modeling unawareness in games. Heifetz, Meier & Schipper (2006) is a key contribution here (and Burkhard Schipper provides an "Unawareness bibliography" with further references on his homepage at UC Davis). It is natural to imagine that belief-dependent motivation interacts with unawareness. For example, negative surprises that reveal previously unforeseen danger could instill fear; so, if 9/11 involved unawareness, then the occurrence of that event might have consequences for subsequent demand for air travel or supply of airport security.

Exploring unawareness using PGT seems potentially interesting, but we only know one paper that is devoted to the topic: Nielsen & Sebald (2017). We quote the informal example with which they open their paper (pp. 2-3)

as it illustrates nicely how unawareness may interact with a belief-dependent motivation (namely, guilt):<sup>62</sup>

Assume it is Bob's birthday, he is planning a party and would be very happy, if Ann could come. Unfortunately Bob's birthday coincides with the date of Ann's final exam at university. She can either decide to take the exam the morning after Bob's party or two weeks later at a second date. Ann is certain that Bob would feel let down, if she were to cancel his party without having a very good excuse. Quite intuitively, although Ann would really like to get over her exam as soon as possible, she might anticipate feeling guilty from letting down Bob if she canceled his party to take the exam the following morning. As a consequence, Ann might choose the second date to avoid letting Bob down. In contrast, consider now the following variant of the same example: Ann knows that Bob is unaware of the second date. In this situation Ann might choose to take the exam on the first date and not feel guilty. Since Bob is unaware of the second date and the final exam is a good excuse, he does not expect Ann to come. Ann knows this and, hence, does not feel guilty as Bob is not let down. In fact, if she were certain that Bob would never become aware of the second date, she probably had an emotional incentive to leave him unaware in order not to raise his expectations.

Motivated beliefs The 2016 summer issue of the Journal of Economic Perspectives contains an interesting symposium on "Motivated Beliefs," with an introduction by Epley & Gilovich (who credit George Loewenstein for taking "the leading role in stimulating and organizing the papers") and contributions by Bénabou & Tirole; Golman, Loewenstein, Moene & Zarri; and Gino, Norton & Weber.<sup>63</sup> The idea is this: Beliefs affect people's well-being. This, in turn, affects how they reason, control information, and gather & evaluate evidence. To some extent, it is argued, they may even choose their

<sup>&</sup>lt;sup>62</sup>We notice, however, that equilibrium concepts that do not capture strategic reasoning, such as self-confirming equilibrium, are consistent with some forms of unawareness. For example, players' whose utility depend on their action and a state variable aggregating others' actions may hold conjectures (confirmed in equilibrium) about such state variable without being aware of the actions of others. See Battigalli, Panebianco, & Pin (2018).

<sup>&</sup>lt;sup>63</sup>See also Loewenstein & Molnar (2018) who discuss related themes.

beliefs, although such choice may be unconscious and the ability to do so is hampered by reality-checks and various costs of having faulty beliefs. Epley & Gilovich mention how the topic has "a long history in psychological science" (p. 139). A particularly important reference would seem to be Kunda (1990), who wrote a highly influential paper on how motivation influences reasoning.

PGT is obviously useful for describing how beliefs affect well-being; such links are embodied in almost every example of belief-dependent utilities that we have exhibited. Second, relatively little work in the literature on motivated beliefs has been formal, and PGT may provide relevant tools for scholars who want to develop theory. Third, PGT is well equipped to deal with how belief-dependent motivation may impact how people control information, and how they gather evidence. These aspects concern choices that presumably can be straightforwardly described in carefully selected game forms. To see this more clearly, note that PGT models the (rational) choice of an agent as a process that takes as qiven his system of conditional beliefs, but the actual beliefs held on the realized path may well depend on his actions (as well as actions of others and exogenous shocks).<sup>64</sup> For example, an agent with imperfect recall may store and recall, possibly at a cost, the flow of information he receives, thus manipulating what he is able to remember and his beliefs. 65 As regards addressing other cognitive phenomena that the literature on motivated beliefs has discussed (e.g., modes of reasoning, evaluating evidence, and unconscious manipulation of beliefs) it seems less clear how PGT may provide useful tools. However, we are optimists and conjecture that PGT might prove useful for approaching those topics as well.

#### 6 Formal framework

The previous exposition of PGT relied on intuition, examples, and some formulas involving game forms and beliefs that could be understood in the context of specific models. But a better appreciation of PGT requires an

<sup>&</sup>lt;sup>64</sup>Indeed, we touch on examples of this sort, e.g. in Section 5 on self-esteem, and also where we discussed the impact of different information structures (Section 4.1 and the part on "higher-order belief-dependence" of this section).

<sup>&</sup>lt;sup>65</sup>Compare with Bénabou & Tirole (2002, p. 871) and their citation from Darwin (1898), where the great scientist describes how he manipulates memory of unpleasant facts to counteract unconscious removal.

understanding of its formal framework, which is necessarily abstract. Here we focus on a simple class of games that covers all the examples of this paper except  $G_3^*$  and  $G_3^{**}$  in the Introduction: we consider p-games obtained from *finite multistage* game forms with *monetary outcomes*, in which players may move simultaneously at some stage and perfectly observe past moves (including chance moves) when they have to make a choice. However, we allow for the possibility of imperfect terminal information, which—as highlighted in the previous sections—may matter for psychological reasons.  $^{67}$ 

The key feature of the analysis is the representation of players' beliefs about how the game form is played, and their beliefs about beliefs, as such beliefs affect the (psychological) utility of end-nodes and expected utility calculations at non-terminal nodes. We mostly assume common knowledge of the rules of the game form and of players' utility functions, i.e., *complete information*. Yet, as illustrated in Sections 4 and 5, incomplete information has to be addressed when we analyze specific motivations such as image concerns and self-esteem, where utility depends on terminal beliefs about unknown personal traits.

Our conceptual perspective mostly relies on B&D, rather than the seminal work of GP&S. The reason is that GP&S only encompasses utilities that depend on players' initial hierarchical beliefs, since at the time of their writing (i) a formal analysis of hierarchical conditional beliefs had yet to be developed, and (ii) the importance of letting utility depend on updated beliefs had not been underscored in applications. B&D instead could leverage on the recently developed theory of hierarchical conditional beliefs (Battigalli & Siniscalchi 1999) and a wealth of applications where updated beliefs enter the utility function. Motivated by conceptual arguments as well as applications, B&D substantially generalize GP&S in several ways. We will briefly point out the differences. Finally, our formalism relies on the recent methodological article by BC&D, which simplifies the analysis by putting only first-order

 $<sup>^{66}</sup>G_3^*$  and  $G_3^{**}$  involve game forms with a continuum of actions.  $G_9$  is covered only if modified (or interpreted) such that chance and Ben move simultaneously.

<sup>&</sup>lt;sup>67</sup>We further simplify in two ways: First, we do not explicitly describe players' non-terminal information when they are not active, which might be relevant for some anticipatory feelings (3.5). Our analysis works "as is" under the assumption that non-active players have the coarsest information consistent with perfect recall. Second, we assume that material consequences accrue at end-nodes only. See BC&D for a more general and explicit analysis of time, in which the game may last for one or more periods, which may have multiple stages, and consequences accrue after each period.

beliefs of *all* players in the domain of utility (so that expected utility depends only on second-order beliefs), but sharpens other aspects, such as the representation and role of players' plans.

**Game form** Formally, we start with a **game form**  $G = \langle I, \bar{H}, \iota, p_0, (\mathcal{P}_i, \pi_i)_{i \in I} \rangle$  with the following elements:

- I is the set of **players** not including **chance**, who is player 0; the set of personal players plus chance is  $I_0 = I \cup \{0\}$ .
- $\bar{H}$  is a finite set of possible sequences of action profiles, or **histories**  $h = (a^k)_{k=1}^{\ell}$  (for different values of  $\ell$ , and possibly including actions of chance) with a *tree* structure: every prefix of a sequence in  $\bar{H}$  (including the empty sequence  $\varnothing$ ) belongs to  $\bar{H}$  as well. Thus, histories in  $\bar{H}$  correspond to nodes of the game tree and  $\varnothing$  is the root. Set  $\bar{H}$  is partitioned into the set of **non-terminal** histories/nodes H and **terminal** histories (paths, end-nodes) Z.
- For each  $h \in H$ ,  $\iota(h) \subseteq I_0$  is the set of **active players**, who perfectly observe h. With this,  $H_i = \{h \in H : i \in \iota(h)\}$  denotes the set of histories where i is active, and the set of feasible action profiles is

$$A\left(h\right) = \left\{ \left(a_{i}\right)_{i \in \iota\left(h\right)} : \left(h, \left(a_{i}\right)_{i \in \iota\left(h\right)}\right) \in \bar{H} \right\} = \times_{i \in \iota\left(h\right)} A_{i}\left(h\right),$$

where  $A_i(h)$  denotes the set of feasible actions of  $i \in \iota(h)$ .

- $p_0$  is the **chance probability** function, which specifies a (discrete) probability density function  $p_0(\cdot|h) \in \Delta(A_0(h))$  for each  $h \in H_0$ .
- For each personal player  $i \in I$ ,
  - $\mathcal{P}_i$  is a partition of Z describing the terminal information of i that satisfies perfect recall (taking into account that active players perfectly observe non-terminal histories),  $\mathcal{P}_i(z)$  denotes the cell containing z;
  - $-\pi_i: Z \to \mathbb{R}$  is a material payoff function.

To illustrate, in  $G_{10}$  (Section 4.1),  $I = \{1, 2\}$ ,  $\bar{H} = \{\emptyset\} \cup \{0, ..., 5\} \cup Z$  with  $Z = \{0, ..., 5\}^2$ ,  $\iota(\emptyset) = \{0\}$ ,  $p_0(x|\emptyset) = \frac{1}{6}$  and  $\iota(x) = \{1\}$ ,  $\pi_1(x, y) = y$ ,  $\mathcal{P}_1(x, y) = \{(x, y)\}$ , and  $\mathcal{P}_2(x, y) = \{0, ..., 5\} \times \{y\}$  for every  $(x, y) \in Z$ .

Beliefs We model the first-order beliefs of (personal) player i as a system  $\alpha_i = (\alpha_i(\cdot|h))_{h \in H \cup \mathcal{P}_i}$  of conditional probabilities about paths of play  $z \in Z$ . We are not assuming that i observes h when he is not active at h ( $h \in H \setminus H_i$ ). In this case we interpret  $\alpha_i(\cdot|h)$  as a "virtual" conditional belief. We assume that: (i)  $\alpha_i$  is consistent with  $p_0$ , (ii) the chain rule holds, and (iii) i's beliefs about simultaneous or past and unobserved actions of other players do not depend on i's chosen action. The latter implies that, for each  $h \in H$ , i's conditional beliefs about the continuation can be obtained by multiplication from i's plan (behavior strategy)  $\alpha_{i,i} \in \times_{h \in H_i} \Delta(A_i(h))$  and i's conjecture  $\alpha_{i,-i} \in \times_{h \in H_{-i}} \Delta(A_{\iota(h)\setminus\{i\}}(h))$  about co-players. Note that i's plan is part of his first-order beliefs. For example, i's initially expected material payoff  $\mathbb{E}[\pi_i; \alpha_i]$  (which may affect his utility via disappointment or frustration) depends on both  $\alpha_{i,i}$  and  $\alpha_{i,-i}$ . As we further explain below, the interpretation is that i plans his contingent choices given his conjecture and thus ends up with an overall system of beliefs about paths.

Let  $\Delta_i^1$  denote i's space of first-order beliefs. We model **second-order** beliefs as systems  $\beta_i = (\beta_i(\cdot|h))_{h \in H \cup \mathcal{P}_i}$  of conditional probabilities about both paths of play  $z \in Z$  and co-players first-order beliefs  $\alpha_{-i} \in \times_{j \in I \setminus \{i\}} \Delta_j^1$  such that: (i) the marginal beliefs about paths form a first-order belief system in  $\Delta_i^1$  (hence they are also consistent with  $p_0$ ), (ii) the chain rule holds, and (iii) i's beliefs about  $\alpha_{-i}$  and simultaneous or past and unobserved actions of other players do not depend on i's chosen action. We let  $\Delta_i^2$  denote the set of second-order beliefs systems of i.

To summarize,  $\alpha_i \in \Delta_i^1$  denotes *i*'s (first-order) beliefs about sequences of actions, or paths, whereas  $\beta_i \in \Delta_i^2$  denotes *i*'s (second-order) overall beliefs about paths and co-players' (first-order) beliefs. In formulas with two-level hierarchies  $(\alpha_i, \beta_i)$ , we maintain the coherence assumption that  $\alpha_i$  is the marginal of  $\beta_i$ .

We point out two conceptually relevant differences with B&D: (i) There, we represented behavior (what players have first-order beliefs about) as a complete description of the actions that players would take at each history where they are active, that is, a (pure) strategy profile rather than a path of play. (ii) In B&D, we explicitly represented first-order beliefs as beliefs about the strategies of *others*. Our explicit interpretation in B&D was that each

 $<sup>^{68}</sup>$ For example, consider a variation of  $G_{10}$  where player 2 observes the report y and then bets on whether player 1 reported the truth or not. Then 2's terminal beliefs are the same as his beliefs before the bet.

player knows his (pure) plan and there is a necessary coincidence between each player's plan and the objective description of how he would behave whenever active, and that such coincidence is transparent to all players (see B&D, p. 11). Here instead we follow BC&D in modeling players' beliefs about paths, hence beliefs about the behavior of everybody.<sup>69</sup> Beliefs about own behavior are interpreted as (possibly non-deterministic) plans, which need not coincide with actual behavior. For example, if i is initially certain that j's plan is  $\alpha_{j,j}$  and then observes a deviation from  $\alpha_{j,j}$ , he may still believe that j' plan was indeed  $\alpha_{j,j}$  but that he took an unplanned action by mistake (a kind of "tremble" as in Selten 1975). The analysis of B&D instead rules this out: every observed action is necessarily interpreted as a planned choice (cf. our discussion of guilt aversion in Section 3.1). In sum, our framework is sufficiently expressive to model players' intentions, their perceptions of the intentions of others, and how such perceptions are affected by observing actions. This is important in standard games to elucidate the difference between, say, forward- and backward-induction reasoning.<sup>70</sup> It is even more important when players care intrinsically for the intentions of others, as with many forms of belief-dependent preferences.

Standard utility Before we describe the belief-based utilities that are characteristic of p-games, it may be helpful to recall how utilities are defined in traditional game theory. Namely, player i's utility has the general form  $u_i: Z \to \mathbb{R}$ . This does not imply that i is "selfish". Caring only about own material reward is a special case  $(u_i(z) = \pi_i(z))$  for all  $z \in Z$ , but i could alternatively be motivated by a host of "social preferences" including altruism, inequity aversion, maximin preferences, or warm glow.<sup>71</sup> However, the forms of motivation that we discussed are ruled out as they require a richer notion of utility.

**Psychological utility and p-games** As argued by BC&D, most forms of belief-dependent motivations for a given player i can be modeled by assuming that, for some terminal history z, i's utility for reaching z depends on the first-order beliefs profile  $(\alpha_j)_{j \in I}$ . Thus, we have utility functions with the general

 $<sup>^{69}</sup>$  The set of strategy profiles is exponentially more complex than the set of paths. Hence, beliefs about paths are simpler.

<sup>&</sup>lt;sup>70</sup>See Battigalli & De Vito (2018) and the references therein.

<sup>&</sup>lt;sup>71</sup>For prominent examples of specific functional forms, see, e.g., Fehr & Schmidt (1999), Bolton & Ockenfels (2000), or (the main text model of) Charness & Rabin (2002).

form  $u_i: Z \times \left(\times_{j \in I} \Delta_j^1\right) \to \mathbb{R}$ . These typically involve both the material payoffs and some features of own or others' initial, interim, or terminal first-order beliefs. For example, in the cheating model of Section 4.1, player 1's utility at terminal history (x,y) depends on his monetary payoff  $\pi_1(x,y) = y$  and on 2's terminal belief about die roll x given report y. In this case, utility depends on the terminal first-order beliefs of someone else. If instead i (besides liking money) dislikes disappointing j, then his utility for reaching z is decreasing j's disappointment  $\left[\mathbb{E}[\pi_j; \alpha_j] - \pi_j(z)\right]^+$ , which depends on j's payoff and his initial belief (cf. Section 3.1 on guilt aversion). In both cases, i's utility of terminal histories depends on payoffs and the (unknown) first-order beliefs of another player. This is like a standard state-dependent utility function. As noted by B&D, the maximization of its expected value can be analyzed with standard techniques leveraging on the dynamic consistency of subjective expected utility maximizers.

For other motivations like aversion to disappointment (Section 3.2), or belief-dependent loss aversion (Section 5), i's utility depends on his expectations (e.g., on the initially expected material payoff  $\mathbb{E}[\pi_i; \alpha_i]$ ), hence on his own plan  $\alpha_{i,i}$ . We showed in Section 3.2 that such forms of own-plan dependence yield dynamic inconsistency of preferences, which implies that some care is required in defining what it means to be subjectively "rational". Similar considerations apply to emotions like frustration and anger (Section 3.3) and to anticipatory feelings with negative or positive valence like anxiety, or suspense (Section 3.5). Essentially, i's plan  $\alpha_{i,i}$  must form an "intra-personal equilibrium" given his overall belief system  $(\alpha_i, \beta_i)$  (cf., e.g., Kőszegi 2010 and the relevant references therein). Next, we explain this in detail.

The combination of a game form and psychological utilities for all players gives a p-game. We consider p-games where the belief-dependence of utility is limited to first-order beliefs. The part of Section 5 on "higher-order belief-dependence" mentions exceptions.

**Subjective rationality** Fix a second-order belief  $\beta_i \in \Delta_i^2$  with marginal first-order belief  $\alpha_i \in \Delta_i^1$  including *i*'s plan  $\alpha_{i,i}$ . For every non-terminal or terminal history  $h' \in \bar{H}$ , we can determine the expectation of  $u_i$  conditional on h', written  $\mathbb{E}[u_i|h';\beta_i]$ . Now consider a history at which *i* is active, viz.  $h \in H_i$ . Each action  $a_i \in A_i(h)$  yields expected utility

$$\bar{u}_{i,h}\left(a_{i};\beta_{i}\right) = \sum_{a_{-i} \in \times_{j \in \iota(h) \setminus \{i\}} A_{j}(h)} \alpha_{i,-i}\left(a_{-i}|h\right) \mathbb{E}\left[u_{i}|\left(h,\left(a_{i},a_{-i}\right)\right);\beta_{i}\right].$$

Belief system  $(\alpha_i, \beta_i)$  satisfies **rational planning** if every action that i expects to take with positive probability is a local best reply, that is,

$$\alpha_{i,i}\left(a_{i}|h\right) > 0 \Rightarrow a_{i} \in \arg\max_{a'_{i} \in A_{i}(h)} \bar{u}_{i,h}\left(a'_{i};\beta_{i}\right)$$

for all  $h \in H_i$ ,  $a_i \in A_i(h)$ . Subjective rationality requires that player i plans rationally given his beliefs and carries out his plan when given the opportunity. Thus, from our perspective, consistency between plan and behavior is not a necessity, but rather a rationality condition. A player's plan of how to play cannot constrain his behavior in the game form, it can only guide it. Consistency between plan and behavior should be a rationality condition. This is the perspective adopted here and in more recent work on p-games.<sup>72</sup>

When  $u_i(z,\alpha)$  does not depend on  $\alpha_i$ , or—more generally—does not depend on i's plan  $\alpha_{i,i}$ , then rational planning is equivalent to the standard sequential rationality condition<sup>73</sup> and i's rational plan can be non-deterministic (not a pure strategy) if and only if i is always indifferent between the pure strategies in the "support" of  $\alpha_{i,i}$  (cf. Remark 1 in BC&D). If instead  $u_i(z,\alpha)$  depends on  $\alpha_{i,i}$ , first, it may be impossible to satisfy the standard sequential rationality condition, second, deterministic rational plans may not exist. The discussion of disappointment in game form  $G_7$  (Section 3.2) provides a simple illustrative example.

Let us here note that a feature of early work on p-games and, most specifically, on guilt was the tacit assumption that a player's behavior is necessarily consistent with his plan, independently of his being rational or not. B&D, in fact, explicitly assume that behavior is necessarily consistent with plan and that this is transparent to the players. We can illustrate this with reference back to the discussion of guilt in Section 3.1). In  $G_6^*$ , when 2 observes trust, he takes for granted that  $\hat{\pi}_1 = 10 \cdot q$  where q denotes the probability subjectively assigned by 1 to share (given trust), because 1 must have planned to trust. With this, the relevant psychological utility of 2 given (trust, grab) can be written as  $14 - \theta_2 \cdot 10 \cdot q$ . This makes a difference for the definition and implications of some solution concepts, like sequential equilibrium, as explained in detail by BC&D in their Section 7.6.

<sup>&</sup>lt;sup>72</sup>See B&D and also Battigalli & De Vito (2018) and Battigalli et al. (2020).

<sup>&</sup>lt;sup>73</sup>The strategy of *i* is ex ante optimal, and the continuation strategy is optimal starting from every  $h \in H_i$ .

Local utilities, incomplete information, and time Solution concepts for p-games can be defined and analyzed starting from the "local" expected utility functions  $\bar{u}_{i,h}: A_i(h) \times \Delta_i^2 \to \mathbb{R}$  ( $i \in I, h \in H_i$ ). To model some belief-dependent action tendencies such as the desire to reciprocate (un)kind behavior (un)kindly (Section 2), or the desire to vent one's own frustration by harming others (3.3), it is convenient to work directly with such history-dependent utility functions, without deriving them from utilities of terminal histories.

A realistic analysis of strategic thinking may have to account for uncertainty about personality traits, i.e., incomplete information. This can be achieved by parameterizing such traits with some vector  $\theta$  and letting players' first-order beliefs concern the unknown part of  $\theta$  as well as behavior. We argued that beliefs about personal traits are also essential to model some motivations such as image concerns and self-esteem (see Sections 4 and 5).

Finally, as in several game-theoretic models in standard economics and behavioral economics, it may be necessary to introduce the role of time explicitly, not only to model preferences over streams of outcomes, but also to model (i) preferences over streams of emotional states (e.g., Ely et al. 2015), (ii) changing belief-dependent referents (e.g., Kőszegi & Rabin 2007, 2009), and (iii) decay of emotional states, such as anger, and the related actions tendencies (e.g., Gneezy & Imas 2014). To do this, one has to distinguish between **stages** and **periods**. Actions are taken within stages and information accrues to players between stages. Outcomes occur within periods, which may comprise multiple stages, and decay or discounting concern different periods. For example, we mentioned in Section 5 how Kőszegi & Rabin's concepts of Personal Equilibrium and Choice-acclimating Personal Equilibrium can be accommodated within our framework by applying the aforementioned notion of rational planning to different time frames.

For a general analysis of the relationship between "global" and "local" utility functions, of incomplete information, and of the role of time see BC&D and the relevant references therein.

**Differences with GP&S** Our perspective and formal analysis differs from that of GP&S in several ways. Let us first address the least important one: unlike GP&S (and B&D) we stop at beliefs of the second order. To our knowledge, this is enough to encompass the overwhelming majority of applications. Thus, we focus on the case where only (first- and) second-order

beliefs matter for expected utility calculations. With this, GP&S consider only initial beliefs about the behavior and the initial beliefs of others. In particular, in game forms with simultaneous moves (where  $Z = A := A(\emptyset)$ ) GP&S consider utilities of the following form:  $\hat{u}_i(a, \beta_{i,-i}^{\varnothing})$ , where  $\beta_{i,-i}^{\varnothing} \in$  $\Delta(A_{-i} \times (\times_{j \neq i} \Delta(A_{-j})))$  denotes i's initial belief about the behavior and the (first-order) beliefs of co-players. We obtain such functional forms in the special case where only initial beliefs about others matter (see B&D for details). The approach of GP&S has three important limitations. First, it rules out models where utility depends on updated beliefs, such as models of sequential reciprocity (Section 2), image concerns (4), deception (4.1), and self-esteem (5). Second, it rules out own-plan-dependent utility as in models with belief-based reference-dependence (Sections 3.2, 3.3, and 5) and anticipatory feelings (3.5). Third, even if utility depends only on initial beliefs, GP&S' framework restricts the toolbox of strategic analysis to (extensions of) traditional equilibrium concepts whereby players have correct beliefs about the (initial) beliefs of others, which therefore never change as play unfolds, on or off the equilibrium path. Indeed, if this were not the case (as in appropriate versions of rationalizability, see BC&D), it would be necessary to address the issue of how players update their beliefs concerning what they care about, i.e., others' beliefs.

## 7 Experiments

In developing theory, our favored approach is to focus on interesting assumptions about forms of belief-dependent motivation, and to explore what they imply. Our main goal is not to explain data (although we are not saying that doing so cannot make sense). With this outlook, once a theory is formulated, a natural next step is to be inquisitive as regards its empirical support. Theories formulated using PGT can be tested for empirical relevance in lab experiments. Our main goal in this section is to give a fairly exhaustive review of aspects of experimental design that are of particular relevance to tests of PGT-based theory. We do not try to give an exhaustive account of results, although we briefly describe some prominent findings.

We cover a series of themes we deem important. For clarity, we give each a separate heading. The first themes are closely related. We list them next to each other. Later on, the themes are less closely related. We apologize that the transitions then may feel choppy.

Belief elicitation Models formulated using PGT suggest ways that particular beliefs impact preferences and play. To conduct lab tests it is often helpful to elicit those beliefs. The very first experiment specifically designed to test a PGT-based prediction was built around that insight. Dufwenberg & Gneezy (2000) considered versions of  $G_1$  (recall: player 2 chooses  $t \in \{0, ..., M\}$ ) as well as Trust Games where 1 could take an outside option (choose out) or choose in and let 2 choose in a subgame structured like  $G_1$  ("Lost Wallet Games"). They measured 1's first order-belief (FOB = expectation of t) by asking 1 to  $guess\ t$  (with rewards for accuracy). And they measured 2's second-order belief (SOB = the conditional expectation of 1's FOB) by asking 2 to  $guess\ 1$ 's  $guess\ (again\ with\ rewards\ for\ accuracy).$  The test for guilt checks whether for subjects in the position of player 2 there is positive correlation between t and those guess-guesses. Dufwenberg & Gneezy performed such tests and, by and large, found support for the theory.

There is a large follow-up literature testing for the empirical relevance of guilt in various game forms, and which often elicits beliefs. Most commonly, binary Trust Games like  $G_6$  are explored. See Cartwright (2019) for a survey.<sup>75</sup>

Some studies elicit beliefs in order to study other forms of motivation than guilt, again to a large degree reporting support. The pioneer to do this for reciprocity theory is Dhaene & Bouckaert (2010), who carefully elicit precisely the belief-data that the task demands (including particular conditional beliefs).<sup>76</sup> And a few recent studies testing aspects of Battigalli, Dufwenberg

 $<sup>7^4</sup>$ The description is precise as regards  $G_1$ . In the Lost Wallet Games, 2 was actually asked about the average guess of all the subjects in the role of 1 who chose *in*. This is crucial to eliciting the right belief, namely 2's belief conditional on 1 choosing *in*.

<sup>&</sup>lt;sup>75</sup>Many scholars report continued support for the theory, though this is not universal. See also Guerra & Zizzo (2004), Charness & Dufwenberg (2006, 2010, 2011), Bacharach, Guerra & Zizzo (2007), Vanberg (2008), Miettinen & Suetens (2008) Reuben, Sapienza & Zingales (2009), Ellingsen, Johannesson, Tjøtta & Torsvik (2010), Bellemare, Sebald & Strobel (2011), Chang, Smith, Dufwenberg & Sanfey (2011), Dufwenberg et al. (2011), Amdur & Schmick (2013), Beck, Kerschbamer, Qiu & Sutter (2013), Bracht & Regner (2013), Kawagoe & Narita (2014), Morrell (2014), Regner & Harth (2014), Andrighetto et al. (2015), Khalmestki et al. (2015), Hauge (2016), Khalmetski (2016), Woods & Servatka (2016), Balafoutas & Fornwanger, Balafoutas & Sutter (2017), Bellemare, Sebald & Suetens (2017), Attanasi, Battigalli, Manzoni & Nagel (2019), Attanasi, Rimbaud & Villeval, Danilov et al. (2019), Dhami et al. (2019), Di Bartolomeo, Dufwenberg, Papa & Passarelli (2019), and Inderst, Khalmetski & Ockenfels (2019).

<sup>&</sup>lt;sup>76</sup>See also Dufwenberg et al. (2011) and Attanasi, Battigalli et al. (2013, 2019). Results

& Smith's models of frustration and anger also do it.<sup>77</sup>

Some scholars we met, mainly of decision-theoretic bent, seemed skeptical on grounds of principle to the idea of belief-dependent motivation, and in the sequel to the idea of eliciting beliefs. This would be people who revere "revealed preference", who argue that beliefs are not real, or at least not observable. Beliefs are merely a theory feature, something that should be viewed only as part of a preference "representation". Conceivably, it is then argued that belief elicitation is pointless, as one then measures something that is not real. In our view, this position has little merit. It militates against introspection, against the rationale involved in humans' use of language, and against some of the experimental evidence we cited. Much like psychologists have by and large given up on their analogous "behaviorist" approach (favored in much of the twentieth century), many calls for revealed preference in economics seem obsolete to us. That said, however, one has to admit that there are many thorny methodological issues surrounding how to best measure subjects' beliefs.<sup>78</sup> Different PGT-related papers take different approaches and some (e.g., Cartwright) discuss pros & cons. See Schotter & Trevino (2014) for a (broader than just PGT) critical survey of the literature on belief elicitation in laboratory experimental economics.

Belief disclosure Charness & Dufwenberg (2006) point out that the guilt hypothesis just discussed is confounded by a form of "false consensus," if 2's choice (done for whatever reason) shapes her SOB such that she believes others believe she made that choice. This would imply that a subject's choice drives his SOB, rather than the other way around (as the guilt story has it). Ellingsen, Johannesson, Tjøtta & Torsvik (2008) propose a clever alternative design, which avoids that issue but which has another problem.

indicate that many, but not all, subjects conform with the theory.

<sup>&</sup>lt;sup>77</sup>See Aina, Battigalli & Gamba (2020) and Dufwenberg, Li & Smith (2018 a,b). Also Persson (2018) performs such a test, although he does so without eliciting beliefs. With the exception of Persson not documenting support for "simple anger," much of this evidence is supportive. See also Battigalli, Dufwenberg & Smith's (Section 5) discussion of several other older experiments which were not designed to test the theory and yet may be viewed in such a light ex post.

<sup>&</sup>lt;sup>78</sup>For example, should guesses be done before or after choices are made; refer to probabilities of a particular co-player's choices or frequencies of choices among a set of subjects one might be matched with; be incentivized or not, and if so how? These questions often have no obvious answers (for example, a quadratic scoring rule may provide precise incentives to reveal a particular expectation, but may also be harder for a subject to understand).

Rather than elicit 2's SOB they elicit 1's FOB, which they then disclosed to 2 before she made her choice. This procedure induces 2's SOB without the risk of false consensus. The drawback, however, is a potential loss of control. In Ellingsen et al.'s design 2 is informed that 1 was not informed that his elicited belief would be handed down to 2. This design feature is important, because if 1 knew then he would have had an incentive to lie (if he believed 2 would believe him). The problem is that when 2 learns that some design information is withheld from the players she may wonder if possibly there are other design aspects that are withheld from her. Perhaps that would affect her behavior.<sup>79</sup>

No elicitation It is not always necessary to elicit beliefs to meaningfully test PGT-based hypotheses. Sometimes patterns of behavior are idiosyncratic enough to a specific theory that clear conclusions can be drawn by observing only choice data. To illustrate, recall the perceived cheating aversion theory of Section 4.1. In their die-roll experiment, using a design matching  $G_{10}$ , Fischbacher & Föllmi-Heusi found reporting frequencies fell in between what would obtain with honest choices (16.7% for each y) and selfish reporting (100% y = 5). One does not need to elicit any beliefs to see that this is in line with the theoretical prediction we described.<sup>80</sup> A further striking insight concerns treatments that manipulate player 2's information. Recall our discussion of how (with perceived cheating aversion) such a change is predicted to undermine the partial lies prediction. Gneezy et al. (2018) ran such treatments, where player 2 were given information about both x and y, and report that 1's behavior indeed changed in the direction of all-or-nothing lies.

Other cases where belief-elicitation was not necessary include Charness & Dufwenberg's (2011) tests regarding "guilt-from-blame" (note especially their remark at the top of p. 1231); Dufwenberg et al.'s (2013) test of negative reciprocity in hold-up problems; and tests that involve a single active player and where the relevant beliefs are pinned down by chance moves—examples include tests of Kőszegi & Rabin's theory as pioneered by Ericson

<sup>&</sup>lt;sup>79</sup>This line of criticism made Ellingsen *et al.*'s approach controversial. Yet the technique has come to be frequently relied on. See, e.g., Attanasi, Battigalli *et al.* (2013, 2019), Khalmetski *et al.*, Bellemare *et al.* (2017), Dhami *et al.*, and Danilov *et al.* 

<sup>&</sup>lt;sup>80</sup>More precisely, 35% choose y = 5, 25% choose y = 4, and all other reports occur with positive frequency that declines with y. Dufwenberg & Dufwenberg explain how this pattern conforms particularly well with an equilibrium they call "sailing-to-the-ceiling.".

& Fuster (2011) and Smith (2019, but written contemporaneously) and Persson's (2018) test of Battigalli, Dufwenberg & Smith's theory.

There is also a large experimental literature which discusses and tests (and often finds support for) more informally formulated notions of reciprocity. Typically, no explicit connection to PGT is made. We will not discuss details, but see Fehr & Gächter for an early highly influential discussion.

Which theory wins? That's the wrong question! It is true that to the extent that we have so far cited experimental results, these have mostly been supportive. Someone may wonder: can really each theory, whether it concerns guilt or reciprocity or anger, be supported? And shouldn't we figure out which one is more relevant? Our answers are yes and a qualified no. We find it plausible that all (or most) of the many motivations we have discussed matter (including, e.g., all the emotions list cited at the end of Section 3.6). Why else would humans have come up with words to describe them?

As regards experiments, the style we focus on in this section is not trying to find a "best" theory for explaining data. Rather, we explore the empirical relevance of propositions that are mostly not mutually exclusive. Do we believe all of these motivations to always be relevant in all situations at the same time? No we certainly don't, and we would like to propose that an important research area of high (but seemingly largely untapped) potential is to figure out which situational cues trigger what kind of motivations.

Communication Charness & Dufwenberg (2006) argue that guilt can help explain why communication, and in particular promises, can foster trust & cooperation. Recall our discussion in Section 3.1. They designed an experiment to test that hypothesis, using methods similar to those of Dufwenberg & Gneezy described above. Vanberg (2008) argued that the results are confounded by another "commitment-based theory," i.e., that decision makers have a belief-independent preference not to break a promise they made. To test his theory, Vanberg came up with an ingenious design, based on a "switching feature". Any subject to whom a pre-play promise were issued was "switched" and replaced by another subject who would play with the person who issued the promise. If there were a switch, the promisor was told but the promisee was not. The key idea is that promisors would suffer expectations-based guilt independently of whether or not a switch occurred, whereas any cost of breaking a promise would apply only if no switch

took place. The commitment-based theory is *not* PGT-based. However, discussions of it typically involve comparisons with Charness & Dufwenberg's belief-based account, so it is important for PGT-scholars to know about Vanberg's work.

The studies cited in the previous paragraph report support for the story they set out to test, and lots of follow-up work has attempted to evaluate which one has stronger pull on data. It is fair to say that the some studies support one while others support the other, and that a debate is ongoing. For more information, see Section 4 in Cartwright (2019) and the discussion in Di Bartolomeo *et al.* (2019).

Exogeneity & causal inference Vanberg's approach is important also for the following methodological reason: Testing for belief-dependent preferences by comparing subjects who self-report different beliefs, as Charness & Dufwenberg did, has the drawback of not relying on exogenously created variation. Subjects are not randomly assigned to their (so-to-say, homegrown) beliefs. This weakens the force with which valid causal evidence can be drawn.

Similarly, if subjects can choose which message to send, then they are not randomly assigned to their messages. Vanberg overcame this last issue via his switching mechanism, creating exogenous variation is whether or not a subject had sent a promise to the player he eventually interacted with. Vanberg did not attempt to create exogenous variation in subjects SOB though, so his design is not ideal for reconsidering Charness & Dufwenberg's hypotheses.

Ederer & Stremitzer (2017) developed a design that involves exogenous variation in subjects' SOB's, and Di Bartolomeo *et al.* developed a design that features exogenous variation in both SOB's and promises. We refer to these studies for more information, while noting that exogenous variation and causal inference has become of high importance in this literature.

Avoidance For certain PGT-related testing purposes it may be useful to employ designs that allow a subject to avoid making another subject aware of a game being played. This would presumably be useful if one were to test ideas that directly involved unawareness, like in Nielsen & Sebald's "party example" of section 5, but it can also be useful for testing whether subjects care about image as described in section 4. For example, consider the design of Dana, Cain & Dawes (2006): Subject i were given a choice whether to

"exit" a \$10 Dictator Game form (like  $G_1$ , with M=10) and take \$9 instead, knowing that the exit option would leave the receiver j nothing and ensure that j never knew that a dictator game form could have been played. Dana et al.'s (2006) design provides a test whether i cares for his image. The idea is that by exiting i may enjoy a (rather) high payoff without suffering a bad image.<sup>81</sup>

Alternative designs can test similar hypotheses without leaving players unaware of strategic possibilities, by allowing a player to avoid revealing a choice. For example, the design of Andreoni & Bernheim (2009) examines "an extended version of the Dictator Game in which (i) chance sometimes intervenes, choosing an unfavorable outcome for the recipient, and (ii) the recipient cannot observe whether chance intervened" (p. 1609).<sup>82</sup> The idea is that by choosing directly the "unfavorable outcome" that chance might implement, a subject can avoid a bad image.

Method of play PGT provides a perspective to elucidate the costs and benefits of the methods of play used in experiments about sequential games and decision problems: the direct-response method, whereby subjects play sequentially (thus generating limited and unbalanced data about responses to different actions of early movers), and the strategy method, whereby subjects commit in advance (but covertly) on their contingent choices, i.e., they play the game in strategic form. <sup>83</sup> When and how should we expect the strategy method to distort behavior compared to what would occur in a sequential situation? This depends on the game form and the motivations that are supposed to drive behavior.

From a theoretical viewpoint, playing with the strategy method should be equivalent to playing with the direct response method when preferences are dynamically consistent. If instead, preferences are dynamically inconsistent, the ability to "tie the hands of one's later selves" by committing to a conditional response may significantly affect behavior in some situations. PGT underscores a wealth of psychological reasons for dynamic inconsistency, and

 $<sup>^{81}</sup>$ An exit choice gives \$-payoff combination (9,0) to i and j, whereby i reveals a preference for that outcome over each of the \$-payoff combinations (10,0), (9,1), and (5,5) which he could have obtained by not exiting. This contradicts models of distributional social preferences like Fehr & Schmidt, Bolton & Ockenfels, and Charness & Rabin.

<sup>&</sup>lt;sup>82</sup>For other related designs, see Dana, Weber & Kuang (2007), Broberg, Ellingsen & Johannesson (2007), and Lazear, Malmendier & Weber (2012).

<sup>&</sup>lt;sup>83</sup>See Brandts & Charness (2011) and the references therein.

suggests when they may be relevant. To illustrate, if subjects play an Ultimatum Game, anger and inequity aversion (a non-belief-dependent motivation) both provide potential explanations for rejections. According to the theory of frustration and anger of Battigalli, Dufwenberg & Smith, anger can only arise when a responder actually observes a greedy offer; instead, the mere contemplation of such a possibility with the strategy method is unlikely to engender enough frustration and anger to make him reject.<sup>84</sup> By contrast, if the only reason for rejections were inequity aversion, a consequentialist model of distributional preferences,<sup>85</sup> then the strategy method should be equivalent to the direct response method.

For example, in the Ultimatum mini-game  $G_5$ , the commitment to reject the greedy offer makes a difference for the distribution of monetary payoffs only if the proposer indeed makes the greedy offer; therefore, the ex ante preference between committing to accept or reject coincides with the post-offer preference between the payoff pairs (0,0) (reject) and (9,1) (accept, with \$1 for the responder). In the Trust mini-games  $G_6$  and  $G_{11}$ , second movers cannot be negatively surprised when they have to choose. Thus, even if they are prone to anger, such emotion cannot affect behavior. Guilt aversion and image concerns, instead, are likely drivers of behavior. But, as explained in Section 6, in such models psychological utility does not depend on the agent's plan, and the resulting belief-dependent preferences are dynamically consistent. Thus, the strategy method should be equivalent to the direct response method. The strategy method should be equivalent to the direct response method.

Other forms of data It may be useful to consider other kinds of data than choices and elicited beliefs to test PGT-based hypotheses. For example, brain imaging data (e.g., fMRI), emotion self-reports ("please rate how strongly you feel emotion X on a scale..".), electrodermal activity, or face-reader data may be useful. Chang, Smith, Dufwenberg & Sanfey (2011) pioneered the use of fMRI for PGT-related purposes, in a study taking the theory of simple guilt to the brain scanner. Chang  $et\ al.$ 's study also involved

<sup>&</sup>lt;sup>84</sup>See Remarks 1-2 in Battigalli, Dufwenberg & Smith.

<sup>&</sup>lt;sup>85</sup>See, e.g., Fehr & Schmidt or Bolton & Ockenfels.

<sup>&</sup>lt;sup>86</sup>Aina *et al.* report that the rejection rate in the Ultimatum mini-game is indeed higher with the direct response method, consistently with the hypothesis that anger is a driver of behavior.

<sup>&</sup>lt;sup>87</sup>For this reason, Attanasi, Battigalli *et al.* (2013, 2016, 2019) analyze a simultaneous-move version of the Trust mini-game.

emotion self-reports, in a way that was mindful of the possibility that pangs of guilt might be counterfactual and yet crucial (compare observation (x) in the guilt part of Section 3 above).<sup>88</sup> We do not know of any face-reader study which was conducted with an explicit PGT-connection in mind, but van Leeuwen, Noussair, Offerman, Suetens, van Veelen & van de Ven (2018) use the technology to explore anger and Battigalli, Dufwenberg & Smith cite their results when motivating their own theory.

## 8 Applications

We hope to inspire applied work. One may think of applications on two levels. The broader one would be to formulate within PGT a model of some belief-dependent motivation that can be used to analyze some class of game forms. We already discussed this topic, for models involving reciprocity, guilt, regret, and anger. Our focus in this section will instead be to applications of a more targeted form, exploring some particular economic setting or question. Our first and third examples of the introduction provide two specific examples, tipping and fashion choice being the economic phenomena scrutinized. Most applied work involves reciprocity, guilt, belief-dependent loss aversion, and image concerns. We listed relevant references in previous sections, but space constraints make it impossible to describe details. Instead, in this section, we attempt to provide some perspective by commenting on possible angles that applied research may follow.

First, one may find settings in which some form of belief-dependent motivation is plausibly operational, but which economists gave scant attention. Likely, economists may not have been used to exploring the sentiment in question, and therefore didn't look at the overall scenario. An example of such research could be Caplin & Leahy's (2004) exploration of how a caring doctor should disclose (or not) information to a patient who is about to undergo an operation and who is affected by anxiety.

Second, one may take a well-known classical model and explore whether predictions are robust to the incorporation of some particular form of beliefdependent motivation which might be plausible in that setting. The paper by

<sup>&</sup>lt;sup>88</sup>Chang *et al.* write (p. 569): "To confirm that participants were actually motivated by anticipated guilt, we elicited their counterfactual guilt for each trial following the scanning session. After displaying a recap of each trial, we asked participants how much guilt they would have felt had they returned a different amount of money."

Patel & Smith (2019) can exemplify. They consider a classic setting of public goods provision and ask how the equilibrium set is affected when players are influenced by guilt.

Third, one may explore whether comparative statics that hold under classical analysis are relevant also if some form of belief-dependent motivation is at play. Dufwenberg et al.'s (2013) study of "the hold-up problem" can exemplify. The analysis involves a form of reciprocity theory where players only reciprocate negative kindness, and involves the game form  $G_{12}$ :

 $[G_{12}]$ 

The parameter  $\omega$  takes different values depending on how details of a hold-up problem are described,<sup>89</sup> but under classical assumption that does not matter. If utility equals (own) money, then the backward induction solution is independent of  $\omega$ ! However, if player 1 is inclined to take revenge, applying the theory, he will wish to choose reject if  $\omega$  is low enough, so the hold-up problem is mitigated in these cases. The authors discuss how this may have implications for how to organize firms, a topic discussed in the academic discipline called *Strategy*.

Fourth, one may observe some pattern of behavior that looks like a puzzle from a classical viewpoint, and explore whether some form of belief-dependent motivation could plausibly be involved and resolve the puzzle. One example of such work could be the study by Conconi et al. (2017), which starts with the arguably puzzling observation that trade disputes are often initiated by incumbent political candidate shortly before elections (e.g., by Obama in 2012). Conconi et al. show that if the electorate is motivated by reciprocity then the observed pattern may be predicted. Another example is the demonstration by Herweg, Muller & Weinschenk (2010) that if agents in a contract-theoretic setting are influenced by belief-dependent loss aversion à la Kőszegi & Rabin, then under some circumstances optimal incentive

 $<sup>^{89}</sup>$ The authors use the following story to derive the special case of  $\omega=2$ : "An artist (player 1) has been asked by a presumptive buyer (player 2) to paint a 'beautiful portrait of 2.' 1 may disagree or agree. In the former case, 1 and 2 go separate ways. In the latter case, 1 spends \$2,000 worth of his/her time on the painting, and a contract says 2 should subsequently pay \$5,000 to 1. The value to 2 is \$8,000, but 2 may complain and claim (falsely) that the portrait is 'rather ugly' and attempt to renegotiate offering a newprice of \$1,000. Given the ambiguity of what constitutes beauty, 1 cannot enforce the \$5,000 payment and will have to accept or reject the new offer. 1 knows that no person other than 2 would pay to acquire the painting." In another story, where the fruits of the effort becomes human capital of player 2 (which 1 cannot deny 2 access to),  $\omega=10$ .

schemes will have a simple form, involving only a fixed wage and a bonus. This addresses the puzzling observation that many real-world contracts look that way, whereas traditional contract theory predicts more complexity and nuance in regards to how payment depends on performance.

Fifth, one may demonstrate how theoretical tools that are known to have certain properties in traditional games may somehow have different implications when a belief-dependent motivation is at play. Bernheim (1994) can demonstrate what we have in mind. Besides shedding light on an image concern and conformity, his study also reveals some new insights solutions on signaling games. Bernheim applies a particular refinement (viz, the "D1 criterion"), argues that it has proven "hostile to pooling in a variety of contexts" (p. 855), and emphasizes how special and telling it is that he nevertheless gets pooling/conformity in his model.

Sixth, we have emphasized several times how in traditional games predictions do not depend on information across end-nodes, but how this ceases to be true in some p-games. One may look for interesting economic settings where this matters. An example is provided by Dufwenberg & Nordblom's (2018) study of tax evasion and guilt, in which it is shown that whether tax returns are private or public influences information across end-nodes and strategic play.

Seventh, experiments may have revealed insights regarding the impact of some belief-dependent preferences. The documentation of effects may, however, have been confined to rather special lab games, and one may explore how the ideas apply in naturally occurring economic settings. For example, consider Charness & Dufwenberg's (2006) evidence that pre-play communication, via mechanisms that depend on guilt, may foster trust and cooperation. They derive that insight in a simple binary Trust Game form. One could consider exploring whether the pattern extends also to a setting with cartelists who meet informally to discuss price-fixing.

Finally, as we wrap up this section, it is natural to wonder where the meat will be going forward. What new specific topics should one look at? It is beyond the scope of our text to exhaustively answer that question. Our goal is to inspire efforts in that direction rather than conduct new studies ourselves. Yet we cannot resist pointing in two particular directions which we feel offer fertile grounds for timely and original contributions:

(i) Starting from the perspective of a particular emotion, we propose that there is scope for interesting work to be done involving *anger*. In Section 3.3, we mentioned many settings including road rage and support for populist candidates. Take the last of these. Suppose that (two or many) competing non-populist candidates (or parties), somehow, make promises that, for one reason or the other, they do not live up to once elected (or at least some folks feel that way). So the elected politicians end up frustrating many voters. Now comes Donald Trump (or Matteo Salvini) and says he'll "drain the swamp," which presumably would be a policy that is annoying for the incumbent. That (populist) candidate may then get the votes of those frustrated voters just because they want to hurt the incumbent! Such a prediction would seem consistent with the model of Battigalli, Dufwenberg & Smith, although one would have to develop the application to say with more confidence.

(ii) Starting from the perspective of a particular societal topic, we believe PGT has much to offer for analyzing behavior and economic outcomes during pandemics. As a matter of fact, we are not the first to argue this point. It is the core idea of a recent paper by Peter Huang (2020). He "advocates building rigorous, theoretical models to analyze how leaders and leadership communications in pandemics can reduce anger, anxiety, and frustration, prevent panic, inhibit complacency, and foster compliance with NPIs" (i.e., non-pharmaceutical interventions, like social distancing and self-quarantining), and goes on to argue that PGT provides relevant tools. We couldn't agree more, and strongly recommend his text for inspiration.

### 9 Comments on methodology

We now address a collection of methodological topics concerning how to model economic settings and how to generate predictions.

Fictitious auxiliary games Can and should we study the psychological phenomena we have described using traditional game theory instead of PGT? Sometimes p-games, most notably perhaps those involving image concerns as described in Section 4, can be turned into "strategically equivalent" traditional games by endowing the observer with a fictitious action space whereby he reports a belief, or estimate of  $\theta_i^{\mathbf{I}}$  and is rewarded with an incentive compatible scoring rule. The receiver's belief—or estimate—is then replaced by his action/report in the sender's utility function. For example, the (pseudo-) utility of observer j is

$$u_j\left(\theta_i^{\mathbf{I}}, a_j\right) = -\left(\theta_i^{\mathbf{I}} - a_j\right)^2 \tag{15}$$

and  $\mathbb{E}\left[\widetilde{\theta_i^{\mathbf{I}}}|\mathcal{P}_j\left(z\right);\alpha_j\right]$  is replaced by  $a_j$  in (13). This works because j maximizes the expected value of (15) by letting  $a_j$  equal the conditional expectation of  $\widetilde{\theta_i^{\mathbf{I}}}$ . As long as i believes in j's rationality, the strategic analysis of the p-game and such an associated traditional game are equivalent.

As in many fields of pure and applied math, transforming a problem into an "equivalent" one may give access to the application of known techniques and results. However, the possibility of such transformations has also engendered the claim that PGT is, after all, not needed: choosing different "weird" assumptions about utility one can go back to good, old, familiar game theory, making everybody feel at home. We are very critical of such attitudes. They confuse formalism with reality. The reality is described by the true game form (something that can be designed and controlled in the lab) and the true utility (which—in so far as it exists—one can try to elicit under appropriate auxiliary assumptions). If player j is passive in reality, coming up with a false representation of reality to claim representability with an old framework can be misleading.  $^{92}$ 

For example, the suggestion that traditional game theory can model everything (possibly after some transformation), even if correct (which we doubt), distracts from important phenomena that cannot be explained when the *true* game is *not* a traditional one, such as the dependence of behavior on the information another player gets when he is inactive. More generally, we propose that adopting indirect techniques referring to known different fields may prevent the conceptually and mathematically correct understanding of the object of analysis, without giving significant benefits over the use of a direct approach. One of the problems is that the theorist often does not realize that he is relying on assumptions that may be wrong/unappealing in some contexts. We illustrate this with a PGT example. GP&S give indirect

 $<sup>^{90}</sup>$ To mention non-obvious ones, results about forward-induction reasoning and rationalizability in a class of infinite dynamic games (Battigalli & Tebaldi 2019) can be applied to p-games with image concerns.

<sup>&</sup>lt;sup>91</sup>As one of our colleagues put it.

<sup>&</sup>lt;sup>92</sup>Furthermore, nobody has shown that all interesting forms of p-games can be turned into "equivalent" standard games. Considering claims made at seminars we have given, we suspect this is not for lack of trying. The only article we know of dealing with the topic is that of Kolpin (1992). He limits attention to the class of p-games considered by GP&S. In our view, while his exercise is pioneering and useful as an attempt at proof-of-concept, the specific assumptions he engages are too convoluted to be practical.

definitions of different kinds of equilibria in p-games. 93 They fix a profile of (initial) belief hierarchies, then look at equilibria of the resulting traditional game and the induced degenerate belief hierarchies whereby beliefs of all orders are correct, and finally define psychological equilibria as fixed points of the resulting correspondence from the space of profiles of belief hierarchies to itself. This is a clever first step, but it is also convoluted. Indeed, B&D show that a direct definition is both feasible and more transparent. Furthermore, when GP&S consider (trembling-hand) perfect equilibria, their indirect approach induces, in our view, a conceptual error: They "forget" that—in the spirit of Selten (1975)—a perfect equilibrium of the p-game should be the limit of equilibria of perturbed p-games whereby players are constrained to "tremble" with strictly positive, although vanishing, minimal probabilities. Using this conceptually correct definition, one can show that a perfect equilibrium always exists under continuity of utility in beliefs.<sup>94</sup> GP&S instead obtain with their indirect definition a different equilibrium concept and show by example that it may fail to exist (p. 73). What went wrong? Using their indirect approach that considers equilibria of auxiliary traditional games, GP&S did not realize that trembles in actions must induce corresponding perturbations in beliefs, hence in the utility of terminal nodes. Once this is taken into account, their counterexample vanishes.

Fictitious auxiliary p-games A convenience-in-modeling argument may cut in the other direction as well: Models with an intrinsic concern for belief-dependent reputation can sometimes offer a compact reduced-form approach to modeling repeated interaction in settings where players are not assumed to have any belief-dependent motivation. That is, a p-game can in such a case be useful modeling tool. For examples that take such an approach, see Morris (2001) and Ottaviani & Sørensen (2006).

**Solution concepts** PGT-analysis involves two key steps: (i) modeling belief-dependent utility, and (ii) applying a solution concept. Step (i) is unique to p-games, step (ii) is relevant also for traditional game theory. Our main goal has been to emphasize what is unique to p-games, so we have focused mostly on step (i). However, since we feel strongly about step (ii),

<sup>&</sup>lt;sup>93</sup>Of the restricted class they analyze.

<sup>&</sup>lt;sup>94</sup>Essentially, B&D prove the existence of perfect equilibrium, to obtain existence of sequential equilibrium, a slightly weaker concept.

let us explain our view.

Economists frequently take for granted that ad hoc notions of equilibrium (whereby players are assumed to have correct beliefs) meaningfully describe strategic interaction. Often such assumption are not well justified. In one-shot play settings, if players reason about each other's rationality and beliefs, inferences should concern steps of deletion of non-best-replies (possibly all the way to a form of rationalizability). If learning is allowed by looking at recurrent interaction, the appropriate solution is (some version of) self-confirming equilibrium, in which beliefs may be incorrect, although consistent with evidence. In neither case is the most commonly applied solution—sequential equilibrium—generally implied.<sup>95</sup> Only in rare cases is it justified to assume that a sequential equilibrium will be played, most notably when this solution concept yields the same prediction as rationalizability, or self-confirming equilibrium.<sup>96</sup>

Since a proper discussion would call for its own article, we have not gone there. Our approach has mainly been consistent with our favored view as we focused on steps of deletion of non-best-replies. But since previous scholarship (including ours) often referred to more traditional notions of equilibrium, we made a few related references when recalling such work.

We hope future work will take the appropriateness and relevance of solution concepts more seriously than has been done in the past.

Transformations of the game form A small literature studies transformations of the game form (with monetary payoffs)<sup>97</sup> that do not change the reduced normal form (see in particular Thompson 1952, Dalkey 1953, and Elms & Reny 1994). Considering only those that preserve perfect recall,

<sup>&</sup>lt;sup>95</sup>B&D extend Kreps & Wilson's (1982) classic notion of sequential equilibrium to p-games. See BC&D for relevant p-games definitions of all solution concepts mentioned above. See Battigalli, Corrao & Sanna (2019) and Jagau & Perea (2018) for epistemic foundations of (versions of) rationalizability.

 $<sup>^{96}</sup>$ A more special circumstance may apply in  $G_{10}$  as analyzed in Section 4.1, given the interpretation that player 2 is player 1's "imagined" audience (as hinted at). If 1 is, so-to-say, "his own audience," we have a single active player, so "forming equilibrium expectations should be easy" (Dufwenberg & Dufwenberg, p. 262).

<sup>&</sup>lt;sup>97</sup>The literature refers to "extensive-form games," which include a specification of players' payoffs at terminal nodes. We interpret such payoffs as monetary one, hence we have what we call "game forms".

they are: interchanging essentially simultaneous moves (INT), <sup>98</sup> coalescing sequential moves (COA; inverse: sequential agent splitting), and addition of a superfluous move (ADD). Any solution concept that just requires knowledge of the normal form, such as Nash equilibrium, or iterated admissibility, is obviously invariant to such transformations (and their inverses). Yet, such invariance cannot be a dogma. Note that the most used solution concept, sequential equilibrium, is only invariant to INT. Versions of rationalizability for extensive-form games are invariant to INT and COA, but not ADD. <sup>99</sup> Our general view is that it is interesting to study the invariance properties of independently motivated solution concepts, without viewing violations of some form of invariance as a flaw. We have just one *caveat*: If simultaneity of moves is represented indirectly, by letting players move in an arbitrary sequence and defining information sets so that the choices of early movers are not observed, then invariance w.r.t. INT is a must. <sup>100</sup>

How does this relate to the analysis of p-games? We have a two-tiered answer:

(i) In general p-games, players' utility depends on the temporal sequence of beliefs. This means that an accurate representation must explicitly account for time and distinguish between periods and stages within periods, as—sav—in models of bargaining protocols. Transformations involving moves across periods should be expected to have an impact on behavior. For example, as shown by Gnezy & Imas (2014), angry agents "cool off" as time goes by and are therefore less willing to harm others (see the discussion of fast vs. slow play in the working paper version of Battigalli, Dufwenberg & Smith). Consider a game form with time in which player i in period t (when he may be frustrated) decides to stop (S) or continue (C). If C is chosen, the play moves to period t+1 when i can harm j(H) or not (N). Since i in period t+1 may have cooled off, a COA transformation making i choose either S, or C&H, or C&N in period t, makes it more likely that i harms j. Thus, it is important to use rich and accurate representations of the rules of the game (see BC&D for a first step in this direction). The rules of games played in the field and in the lab, by necessity, imply that information accrues to both active and inactive players. Since information affects beliefs,

<sup>&</sup>lt;sup>98</sup>Note that the formalism initially adopted to describe games indirectly represented simultaneous moves in an arbitrary sequential order without flows of information to late movers

 $<sup>^{99}\</sup>mathrm{See}$  the discussion and characterization by Battigalli, Leonetti & Maccheroni (2020).

<sup>&</sup>lt;sup>100</sup>See the characterization of INT-invariance by Bonanno (1992).

we should formally represent the information of inactive players. Also, the rules of the game only specify information flows, not how much information is retained by players. To keep a sharp separation between game form and players' personal features, the former should specify information flows as, for example, in Myerson (1986), or in the literature on repeated games with imperfect monitoring (e.g., Mailath & Samuelson 2006). The use of traditional information partitions is acceptable under the presumption that players have perfect recall (and this is common knowledge).

(ii) Compared to traditional game theory, invariance with respect to a particular transformation depends also on the form of psychological utility functions, not only on the solution concept. For example, if utility depends only on the initial, or the terminal beliefs of others, sequential equilibrium is invariant w.r.t. INT (which, however, may not be applicable if simultaneity is represented directly as we do), whereas extensive-form rationalizability is invariant to both INT and COA. Yet, utility may depend on beliefs in ways that prevent any form of invariance.

Private sensitivities and incomplete information In preceding sections  $\theta_i$  often denoted player i's "sensitivity" with respect to some psychological concerns (e.g. guilt or reciprocity). It makes sense to assume that player j is not informed of  $\theta_i$ , and has to form beliefs about it in order to predict i's behavior. In this case the analysts must consider beliefs in p-games with incomplete information. We already discussed this topic, in Section 4, where, however, we focused on image concerns, i.e. how i might care about j's beliefs about  $\theta_i$ . However, even absent image concerns, incomplete information is very natural and important. If  $\theta_i$  measures i's sensitivity to feel, e.g., guilt, reciprocity, or anger, then j's beliefs about  $\theta_i$  may matter not because i cares about those beliefs but because they will impact j's choices. i

Most applied work using PGT otherwise either assumes complete information, or (legitimately) ignores the issue by focusing on the shape of the best reply correspondences, without using solution concepts. The analysis of psychological games with incomplete information is developed in Attanasi, Battigalli & Manzoni (2016), Bjorndahl, Halpern & Pass (2015), and BC&D, which is more general and systematic. Incomplete information is the lack of

 $<sup>^{101}</sup>$ For analyses of games where players are uncertain about each others' such sensitivities, see Attanasi, Battigalli *et al.* (2013, 2016, 2019) and Battigalli, Charness & Dufwenberg (2013, footnote 5) on guilt, Sohn & Wu (2019) on reciprocity, and Aina *et al.* on anger.

common knowledge of some features of the game form (feasible sets, information structure, material payoff functions) or of players' preferences. As explained in BC&D, if we rely on solution concepts with an epistemic foundation (versions of rationalizability) or a learning foundation (versions of self-confirming equilibrium), incompleteness of information can be addressed directly by means of relatively straightforward extensions of such solution concepts. 102 The situation is different if the analyst maintains the traditional equilibrium assumption that players have correct conjectures about the decision rules associating co-players' private information with their behavior. In this case, following the seminal contribution by Harsanyi (1967-68), the traditional approach proceeds in two steps. It first posits an *implicit* representation of players' possible interactive beliefs about the unknown parameter vector  $\theta$  by means of "types" encoding both private knowledge about  $\theta$ and hierarchical exogenous beliefs; this yields a "Bayesian game". Then one proceeds to analyze the Bayesian perfect (or sequential) equilibria of such a game, that is, profiles of decision rules whereby every type of every player carries out a sequential best reply to the co-players' decision rules given the beliefs of this type. The key to connecting Harsanyi's approach to PGT is that, in a Bayesian equilibrium, each type of each player is associated to a hierarchy of (exogenous and endogenous) beliefs, which—in a p-game—enters the utility functions. This allows us to obtain all the known "traditional" equilibrium concepts (those that postulate correct beliefs, as in GP&S and B&D) as special cases of Bayesian perfect equilibrium. Harsanyi's approach gives an additional degree of flexibility seldom noticed by non-specialists: a Bayesian game may feature multiple types of the same player with the same exogenous hierarchy of beliefs, even when there is complete information! This allows for equilibria where types with the same exogenous belief hierarchy have different plans and different beliefs about co-players. Since observed actions signal types, in equilibrium players may update their higher-order beliefs about the beliefs (and plans) of co-players, hence their intentions, which is instead prevented by the equilibrium concepts of GP&S and B&D. 103

**Bounded rationality** All the work on PGT that we know of assumes that players are rational. Clearly, however, the idea of bounded rationality makes as much sense in the context of p-games as in standard games. We propose

 $<sup>^{102}</sup>$ See Sections 7.1-2 of BC&D and the references therein.

 $<sup>^{103}</sup>$ See Section 7.3 of BC&D.

that exploring related topics might be interesting, although we mainly leave exploring it for future efforts. We see at least one interesting potential interaction between bounded rationality and PGT: While the latter models emotions as part of players' utility, it is known that some emotions such as anxiety (Rauh & Seccia 2006) and anger (Gneezy & Imas 2014) can hamper rational cognition, and this is factored in by early movers who can trigger such emotions. PGT in its current form is not equipped to model such effects.

### 10 Concluding remarks

Decisions are driven by a plethora of desires. Yet economists' approaches traditionally took a narrow view, focusing mainly on concern for own income (or consumption). When richer models were proposed, it was often taken as an advantage if the deviations from the tradition were limited. For example, much of the literature on "social preferences" considers it a success if data sets can be explained using utilities defined on distributions of material payoffs according to simple formulas.<sup>104</sup>

Being spare is not necessarily a virtue. If human psychology is rich and multi-faceted, one cannot know the effect of the involved sentiments unless one dives in and explores how and why that plays out in economic contexts. Many interesting desires that shape behavior in important ways take the form of belief-dependent motivation. This includes reciprocity, emotions, image concerns, belief-dependent loss aversion, and self-esteem. We have argued that the mathematical framework of psychological game theory (PGT) is useful and needed for modeling such sentiments, and we have tried to show why & how.

Working with PGT is exciting, and interest is growing. For example, the Journal of Economic Behavior and Organization recently published a special issue on PGT, comprising 15 articles (cited above); see Dufwenberg & Patel (2019) for the guest editors' introduction. We derive utility from our hope (=item #12 in Elster's list from Section 3.6) to inspire others to follow suit.

<sup>&</sup>lt;sup>104</sup>See, e.g., Fehr & Schmidt, Bolton & Ockenfels, Charness & Rabin for models, and Cooper & Kagel (2009) for a survey in that spirit.

# References

- [1] Abeler, Johannes, Daniele Nosenzo, and Collin Raymond. 2019. "Preferences for Truth-Telling". *Econometrica* 87: 1115-53.
- [2] Aina, Chiara, Pierpaolo Battigalli, and Astrid Gamba. 2020. "Frustration and Anger in the Ultimatum Game: An Experiment". Games and Economic Behavior, in print.
- [3] Akerlof, George. 1982. "Labour Contracts as a Partial Gift Exchange". Quarterly Journal of Economics 97: 543-69.
- [4] Aldashev, Gani, Georg Kirchsteiger, and Alexander Sebald. 2017. "Assignment Procedure Biases in Randomized Policy Experiments". *The Economic Journal* 127: 873-895.
- [5] Amdur, David, and Ethan Schmick. 2013. "Does the Direct-Response Method Induce Guilt Aversion in a Trust Game?". Economics Bulletin 33: 687-693.
- [6] Andreoni, James, and B. Douglas Bernheim. 2009. "Social Image and the 50-50 Norm: A Theoretical and Experimental Analysis of Audience Effects". *Econometrica* 77: 1607-1636.
- [7] Andrighetto, Giulia, Daniela Grieco, and Luca Tummolini. 2015. "Perceived Legitimacy of Normative Expectations Motivates Compliance with Social Norms when Nobody is Watching". Frontiers in Psychology 6, 1413.
- [8] Attanasi, Giuseppe, Pierpaolo Battigalli, and Elena Manzoni. 2016. "Incomplete Information Models of Guilt Aversion in the Trust Game". Management Science 62: 648-667.
- [9] Attanasi, Giuseppe, Pierpaolo Battigalli, Elena Manzoni, and Rosemarie Nagel. 2019. "Belief-Dependent Preferences and Reputation: Experimental Analysis of a Repeated Trust Game". *Journal of Economic Behavior and Organization* 167: 341-360.
- [10] Attanasi Giuseppe, Pierpaolo Battigalli, and Rosemarie Nagel. 2013. "Disclosure of Belief-Dependent Preferences in the Trust Game". Bocconi University IGIER w.p. 506.

- [11] Attanasi, G., and R. Nagel. 2008. "A Survey of Psychological Games: Theoretical Findings and Experimental Evidence." In A. Innocenti and P. Sbriglia (Eds.) Games, Rationality and Behavior. Essays on Behavioral Game Theory and Experiments, Palgrave McMillan, 204-232.
- [12] Attanasi, Giuseppe, Claire Rimbaud, and Marie-Claire Villeval. 2019. "Embezzlement and Guilt Aversion". Journal of Economic Behavior and Organization 167: 409-429.
- [13] Averill, James R. 1982. Anger and Aggression: An Essay on Emotion. New York: Springer.
- [14] Azar, Ofer H. 2019. "The Influence of Psychological Game Theory". Journal of Economic Behavior and Organization 167: 445-453.
- [15] Bacharach, Michael, Gerardo Guerra, and Daniel J. Zizzo. 2007. "The Self-Fulfilling property of Trust: an Experimental Study". Theory and Decision 63(4): 349-388.
- [16] Balafoutas, Loukas. 2011. "Public Beliefs and Corruption in a Repeated Psychological Game". Journal of Economic Behavior and Organization 78: 51-59.
- [17] Balafoutas, Loukas, and Helena Fornwagner. 2017. "The Limits of Guilt". Journal of the Economic Science Association 3: 137-148.
- [18] Balafoutas, Loukas, and Matthias Sutter. 2017. "On the Nature of Guilt Aversion: Insights from a New Methodology in the Dictator Game". Journal of Behavioral and Experimental Finance 13: 9-15.
- [19] Battigalli Pierpaolo, Gary Charness, and Martin Dufwenberg. 2013. "Deception: The Role of Guilt". Journal of Economic Behavior and Organization 93: 227-232.
- [20] Battigalli Pierpaolo, Roberto Corrao, and Martin Dufwenberg. 2019. "Incorporating Belief-Dependent Motivation in Games". Journal of Economic Behavior and Organization 167: 185-218.
- [21] Battigalli, Pierpaolo, Roberto Corrao, and Federico Sanna. 2020. "Epistemic Game Theory without Types Structures: An Application to Psychological Games". Games and Economic Behavior 120: 28-57.

- [22] Battigalli, Pierpaolo, and Nicodemo De Vito. 2018. "Beliefs, Plans, and Perceived Intentions in Dynamic Games". Bocconi University IGIER w.p. 629.
- [23] Battigalli, Pierpaolo, and Martin Dufwenberg. 2007. "Guilt in Games". American Economic Review 97(2): 170-176.
- [24] Battigalli, Pierpaolo, and Martin Dufwenberg. 2009. "Dynamic Psychological Games". *Journal of Economic Theory* 144: 1-35.
- [25] Battigalli, Pierpaolo, Martin Dufwenberg, and Alec Smith. 2019. "Frustration, Aggression and Anger in Leader-Follower Games". *Games and Economic Behavior* 117, 15-39.
- [26] Battigalli, Pierpaolo, Paolo Leonetti, and Fabio Maccheroni. 2020. "Behavioral Equivalence of Extensive Game Structures". Games and Economic Behavior 121: 533-547.
- [27] Battigalli, Pierpaolo, Fabrizio Panebianco, and Paolo Pin. 2018. "Learning and Self-Confirming Equilibrium in Network Games". Bocconi University IGIER w.p. 637.
- [28] Battigalli, Pierpaolo, and Marciano Siniscalchi. 1999. "Hierarchies of Conditional Beliefs and Interactive Epistemology in Dynamic Games". *Journal of Economic Theory* 88: 188-230.
- [29] Battigalli, Pierpaolo, and Pietro Tebaldi. 2019. "Interactive Epistemology in Simple Dynamic Games with a Continuum of Strategies". *Economic Theory* 68: 737-763.
- [30] Baumeister, Roy, Arlene Stillwell, and Todd Heatherton. 1994. "Guilt: An Interpersonal Approach". *Psychological Bulletin* 115: 243-267.
- [31] Beck Adrian, Rudolf Kerschbamer, Jianying Qiu, and Matthias Sutter. 2013. "Shaping Beliefs in Experimental Markets for Expert Services: Guilt Aversion and the Impact of Promises and Money-Burning Options". Games and Economic Behavior 81: 145-164.
- [32] Bell, David. 1982. "Regret in Decision Making under Uncertainty". Operations Research 30: 961-981.

- [33] Bell, David. 1985. "Disappointment in Decision Making under Uncertainty". Operations Research 33: 1-27.
- [34] Bellemare, Charles, Alexander Sebald, Martin Strobel. 2011. "Measuring Willingness to Pay to Avoid Guilt: Estimation using Equilibrium and Stated Belief Models". *Journal of Applied Economics*. 26: 437-453.
- [35] Bellemare, C., A. Sebald, and Sigrid Suetens. 2017. "A Note on Testing Guilt Aversion". *Games and Economic Behavior* 102: 233-239.
- [36] Bénabou, Roland, and Jean Tirole. 2002. "Self-Confidence and Personal Motivation". Quarterly Journal of Economics, 117: 871-915.
- [37] Bénabou, Roland, and Jean Tirole. 2006. "Incentives and Prosocial Behavior". American Economic Review 96: 1652-78.
- [38] Bénabou, Roland, and Jean Tirole. 2011. "Identity, Morals, and Taboos: Beliefs as Assets". Quarterly Journal of Economics 126: 805-855.
- [39] Bénabou, Roland, and Jean Tirole. 2016. "Mindful Economics: The Production, Consumption, and Value of Beliefs". *Journal of Economic Perspectives* 30: 141-64.
- [40] Berglas, Steven, and Edwin E. Jones. 1978. "Drug Choice as a Self-Handicapping Strategy in Response to Noncontingent Success". *Journal of Personality and Social Psychology* 36: 405-417.
- [41] Berkowitz, Leonard. 1978. "Whatever Happened to the Frustration-Aggression Hypothesis?". American Behavioral Scientist 21: 691-708.
- [42] Berkowitz, Leonard. 1989. "Frustration-Aggression Hypothesis: Examination and Reformulation". *Psychological Bulletin* 106: 59-73.
- [43] Bernheim, Douglas. 1994. "A Theory of Conformity". Journal of Political Economy 102: 841-877.
- [44] Bicchieri, Cristina. 2006. The Grammar of Society: The Nature and Dynamics of Social Norms. Cambridge University Press.
- [45] Bierbrauer, Felix, and Nick Netzer. 2016. "Mechanism Design and Intentions". *Journal of Economic Theory* 163: 557–603.

- [46] Bierbrauer, Felix, Axel Ockenfels, Andreas Pollak, and Désirée Rückert. 2017. "Robust Mechanism Design and Social Preferences". *Journal* of Public Economics 149: 59-80.
- [47] Bjorndahl, Adam, Joseph Halpern, and Rafael Pass. 2015. "Bayesian Games with Intentions". Proceedings of the 15th Conference on Theoretical Aspects of Rationality and Knowledge: 75–82.
- [48] Blume, Andreas, Ernest K. Lai, and Wooyoung Lim. 2019. "Eliciting Private Information with Noise: The Case of Randomized Response". Games and Economic Behavior 113: 356-380.
- [49] Bolton, Gary, and Axel Ockenfels. 2000. "ERC: A Theory of Equity, Reciprocity, and Competition". American Economic Review 90: 166-93.
- [50] Bonanno, Giacomo. 1992. "Set-Theoretic Equivalence of Extensive-Form Games". International Journal of Game Theory 20: 429-447.
- [51] Bracht, Jurgen, and Tobias Regner. 2013. "Moral Emotions and Partnership". Journal of Economic Psychology. 39: 313-326.
- [52] Brandts Jordi, and Gary Charness. 2011. "The Strategy versus the Direct-Response Method: A First Survey of Experimental Comparisons". Experimental Economics 14: 375-398.
- [53] Broberg, Tomas, Tore Ellingsen, and Magnus Johannesson. 2007. "Is Generosity Involuntary?". *Economics Letters* 94: 32-37.
- [54] Caplin, Andrew, and John Leahy. 2001. "Psychological Expected Utility Theory and Anticipatory Feelings". Quarterly Journal of Economics 116: 55-79.
- [55] Caplin, Andrew, and John Leahy. 2004. "The Supply of Information by a Concerned Expert". Economic Journal 114: 487-505.
- [56] Card, David, and Gordon Dahl. 2011. "Family Violence and Football: The Effect of Unexpected Emotional Cues on Violent Behavior". Quarterly Journal of Economics 126: 103-143.
- [57] Cardella, Eric. 2016. "Exploiting the Guilt Aversion of Others: Do Agents Do It and Is It Effective?". Theory and Decision 80: 523-560.

- [58] Caria, Stefano, Marcel Fafchamps. 2019. "Expectations, Network Centrality, and Public Good Contributions: Experimental Evidence from India". Journal of Economic Behavior and Organization 167: 391-408.
- [59] Cartwright, Edward. 2019. "A Survey of Belief-Based Guilt Aversion in Trust and Dictator Games". Journal of Economic Behavior and Organization 167: 430-444.
- [60] Çelen, Bogaçhan, Andrew Schotter, and Mariana Blanc. 2017. "On Blame and Reciprocity: Theory and Experiments". Journal of Economic Theory 169: 62-92.
- [61] Cerreia-Vioglio, Simone, David Dillenberger, and Pietro Ortoleva. 2018. "An Explicit Representation for Disappointment Aversion and Other Betweenness Preferences". Bocconi University IGIER w.p. 631.
- [62] Chang, Luke, and Alec Smith. 2015. "Social Emotions and Psychological Games". Current Opinion in Behavioral Sciences 5: 133-140.
- [63] Chang, Luke, Alec Smith, Martin Dufwenberg, and Alan Sanfey. 2011. "Triangulating the Neural, Psychological, and Economic Bases of Guilt Aversion". Neuron 70: 560-72.
- [64] Charness, Gary, and Martin Dufwenberg. 2006. "Promises and Partnership". Econometrica 74: 1579-1601.
- [65] Charness, Gary, and Martin Dufwenberg. 2010. "Bare Promises: An Experiment". *Economics Letter* 107: 281-283.
- [66] Charness, Gary, and Martin Dufwenberg. 2011. "Participation". American Economic Review 101: 1213-39.
- [67] Charness, Gary, and Matthew Rabin. 2002. "Understanding Social Preferences with Simple Tests". Quarterly Journal of Economics 117: 817-869.
- [68] Conconi, Paola, David R. DeRemer, Georg Kirchsteiger, Lorenzo Trimarchi, and Maurizio Zanardi. 2017. "Suspiciously Timed Trade Disputes". Journal of International Economics 105: 57-75.

- [69] Connolly, Terry, and David Butler. 2006. "Regret in Economic and Psychological Theories of Choice". Journal of Behavioral Decision Making 19: 139-154.
- [70] Cooper, David J., and John H. Kagel. 2016. "Other Regarding Preferences: A Survey of Experimental Results". In *The Handbook of Experimental Economics*. Vol. 2. Princeton: Princeton University Press.
- [71] Cox, James, Daniel Friedman, and Vjollca Sadiraj. 2008. "Revealed Altruism". *Econometrica* 76: 31-69.
- [72] d'Adda, Giovanna, Martin Dufwenberg, Francesco Passarelli, Guido Tabellini. 2020. "Social Norms and Private Values". Unpublished.
- [73] Dalkey, Norman. 1953. "Equivalence of Information Patterns and Essentially Determinate Games". Contributions to the Theory of Games II, H. Kuhn and A. Tucker (Eds.). Princeton University Press, 217-243.
- [74] Dana, Jason, Daylian Cain, and Robyn Dawes. 2006. "What You Don't Know Won't Hurt Me: Costly (but Quiet) Exit in Dictator Games". Organizational Behavior and Human Decision Processes 100: 193-201.
- [75] Dana, Jason, Roberto Weber, and Jason Kuang. 2007. "Exploiting Moral Wiggle Room: Experiments Demonstrating an Illusory Preference for Fairness". *Economic Theory* 33: 67-80.
- [76] Danilov, Anastasia, Kyril Khalmetski, and Dirk Sliwka. 2019. "Norms and Guilt". CESifo w.p. 6999. CESifo Group Munich.
- [77] Darwin, Francis. 1898, *The Life and Letters of Charles Darwin*. Edited by his son Francis Darwin. New York: D. Appleton and Co..
- [78] Dhami, Sanjit, Mengxing Wei, and Ali al-Nowaihi. 2019. "Public Goods Games and Psychological Utility: Theory and Evidence". *Journal of Economic Behavior and Organization* 167: 361-390.
- [79] Dhaene, Geert, and Jan Bouckaert. 2010. "Sequential Reciprocity in Two-Player, Two-Stage Games: An Experimental Analysis". Games and Economic Behavior 70: 289-303.

- [80] Di Bartolomeo, Giovanni, Martin Dufwenberg, Stefano Papa, and Francesco Passarelli. 2019. "Promises, Expectations and Causation". Games and Economic Behavior 113: 137-46.
- [81] Dollard, John, Leonard W. Doob, Neal E. Miller, O. H., Mowrer, and Robert R. Sears. 1939. Frustration and Aggression. New Haven: Yale University Press.
- [82] Dufwenberg, Martin. 2002. "Marital Investment, Time Consistency and Emotions". Journal of Economic Behavior and Organization 48: 57-69.
- [83] Dufwenberg, Martin. 2008. "Psychological Games". In *The New Pal-grave Dictionary of Economics* edited by S.N. Durlauf and L.E. Blume. Volume 6: 714-718. Palgrave Macmillan.
- [84] Dufwenberg, Martin, and Martin A. Dufwenberg. 2018. "Lies in Disguise A Theoretical Analysis of Cheating". Journal of Economic Theory 175: 248-264.
- [85] Dufwenberg, Martin, Simon Gächter, and Heike Hennig-Schmidt. 2011. "The Framing of Games and the Psychology of Play". Games and Economic Behavior 73: 459-478.
- [86] Dufwenberg, Martin and Uri Gneezy. 2000. "Measuring Beliefs in an Experimental Lost Wallet Game". Games and Economic Behavior 30: 163-182.
- [87] Dufwenberg, Martin, and Georg Kirchsteiger. 2000. "Reciprocity and Wage Undercutting". European Economic Review 44: 1069-1078.
- [88] Dufwenberg, Martin, and Georg Kirchsteiger. 2004. "A Theory of Sequential Reciprocity". *Games and Economic Behavior* 47: 268-298.
- [89] Dufwenberg, Martin and Georg Kirchsteiger. 2019. "Modelling Kindness". *Journal of Economic Behavior and Organization* 167: 228-234.
- [90] Dufwenberg, Martin, Flora Li, and Alec Smith. 2018. "Promises and Punishment". Unpublished.
- [91] Dufwenberg, Martin, Flora Li, and Alec Smith. 2018. "Threats". Unpublished.

- [92] Dufwenberg, Martin, and Senran Lin. 2019. "Regret Games". Unpublished.
- [93] Dufwenberg, Martin, and Michael Lundholm. 2001. "Social Norms and Moral Hazard". *Economic Journal* 111: 506-525.
- [94] Dufwenberg, Martin, and Katarina Nordblom. 2018. "Tax Evasion with a Conscience". Unpublished.
- [95] Dufwenberg, Martin, and Amrish Patel. 2017. "Reciprocity Networks and the Participation Problem". *Games and Economic Behavior* 101: 260-272.
- [96] Dufwenberg, Martin, and Amrish Patel. 2019. "Introduction to Special Issue on Psychological Game Theory," Journal of Economic Behavior and Organization 67, 181-84.
- [97] Dufwenberg, Martin, and David Rietzke. 2016. "Banking on Reciprocity: Deposit Insurance and Insolvency". Unpublished.
- [98] Dufwenberg, Martin, Alec Smith, and Matt Van Essen. 2013. "Hold-up: With a Vengeance". *Economic Inquiry* 51: 896-908.
- [99] Ederer, Florian, and Alexander Stremitzer. 2017. "Promises and Expectations". Games and Economic Behavior 106: 161-178.
- [100] Eil, David, and Justin Rao. 2011. "Asymmetric Processing of Objective Information about Yourself". American Economic Journal: Microeconomics 3: 114-138.
- [101] Ellingsen, Tore, Magnus Johannesson, Sigve Tjøtta, Gaute Torsvik. 2010. "Testing Guilt Aversion". Games and Economic Behavior 68: 95-107.
- [102] Elms, Susan, and Phil Reny. 1994. "On the Strategic Equivalence of Extensive Form Games". *Journal of Economic Theory* 62: 1-23.
- [103] Ellingsen, Tore, and Magnus Johannesson. 2008. "Pride and Prejudice: The Human Side of Incentive Theory". American Economic Review 98: 990-1008.

- [104] Elster, Jon. 1989. "Social Norms and Economic Theory". The Journal of Economic Perspectives 3: 99-117.
- [105] Elster, Jon. 1996. "Rationality and the Emotions". *Economic Journal* 106: 1386-1397.
- [106] Elster, Jon. 1998. "Emotions and Economic Theory". Journal of Economic Literature 36: 47-74.
- [107] Ely, Jeffrey, Alexander Frankel, and Emir Kamenica. 2015. "Suspense and Surprise". *Journal of Political Economy* 123, 215-260.
- [108] Engelbrecht-Wiggans, Richard. 1989. "The Effect of Regret on Optimal Bidding in Auctions". *Management Science* 35: 685-92.
- [109] Engelbrecht-Wiggans, Richard, and Elena Katok. 2008. "Regret and Feedback Information in First-Price Sealed-Bid Auctions". Management Science 54: 808-819.
- [110] Epley, Nicholas, and Thomas Gilovich. 2016. "The Mechanics of Motivated Reasoning". *Journal of Economic Perspectives* 30: 133-40.
- [111] Ericson, Keith Marzilli, and Andreas Fuster. 2011. "Expectations as Endowments: Evidence on Reference-Dependent Preferences from Exchange and Valuation Experiments". Quarterly Journal of Economics 126: 1879-1907.
- [112] Fagin, Ronald, and Joseph Halpern. 1988. "Belief, Awareness and Limited Reasoning". *Artificial Intelligence* 34: 39-76.
- [113] Falk, Armin, and Urs Fischbacher. 2006. "A Theory of Reciprocity". Games and Economic Behavior 54: 293-315.
- [114] Fehr, Ernst, and Simon Gächter. 2000. "Fairness and Retaliation: The Economics of Reciprocity". *Journal of Economic Perspectives* 14: 159-181.
- [115] Fehr, Ernst, and Ivo Schurtenberger. 2018. "Normative Foundations of Human Cooperation". *Nature* 2: 458-468.
- [116] Fehr, Ernst, and Klaus Schmidt. 1999. "A Theory of Fairness, Competition, and Cooperation". Quarterly Journal of Economics 114: 817-868.

- [117] Filiz-Ozbay, Emel, and Erkut Ozbay. 2007. "Auctions with Anticipated Regret: Theory and Experiment". *American Economic Review* 97: 1407-1418.
- [118] Fischbacher, Urs, and Franziska Föllmi-Heusi. 2013. "Lies in Disguise

   An Experimental Study on Cheating". Journal of the European Economic Association 11: 525-547.
- [119] Geanakoplos, John. 1996. "The Hangman Paradox and the Newcomb's Paradox as Psychological Games". Cowles Foundation Discussion Paper No. 1128.
- [120] Geanakoplos, John, David Pearce, and Ennio Stacchetti. 1989. "Psychological Games and Sequential Rationality". *Games and Economic Behavior* 1: 60-80.
- [121] Gilboa, Itzhak, and David Schmeidler. 1988. "Information Dependent Games: Can Common Sense be Common Knowledge?". *Economics Letters* 27: 215–221.
- [122] Gill, David, and Victoria Prowse. 2012. "A Structural Analysis of Disappointment Aversion in a Real Effort Competition". *American Economic Review* 102: 469–503.
- [123] Gino, Francesca, Michael I. Norton, and Roberto A. Weber. 2016. "Motivated Bayesians: Feeling Moral While Acting Egoistically". *Journal of Economic Perspectives* 30: 189-212.
- [124] Glazer Amihai, and Kai A. Konrad. 1996. "A Signaling Explanation for Charity". *American Economic Review* 86: 1019-1028.
- [125] Gneezy, Uri, and Alex Imas. 2014. "Materazzi Effect and the Strategic Use of Anger in Competitive Interactions". *Proceedings of the National Academy of Sciences* 111: 1334-1337.
- [126] Gneezy, Uri, Agnel Kajackaite, and Joel Sobel. 2018. "Lying Aversion and the Size of the Lie". American Economic Review 108: 419-453.
- [127] Golman, Russell, George Loewenstein, Karl Ove Moene and Luca Zarri. 2016. "The Preference for Belief Consonance". *Journal of Economic Perspectives* 30: 165-88.

- [128] Goranson, Richard, and Leonard Berkowitz. 1966. "Reciprocity and Responsibility Reactions to Prior Help". *Journal of Personality and Social Psychology* 3: 227-232.
- [129] Gradwohl, Ronen, and Rann Smorodinsky. 2017. "Perception Games and Privacy," Games and Economic Behavior 104, 293-308.
- [130] Grossman, Zachary and Joel J. van der Weele. 2017. "Self-Image and Willful Ignorance in Social Decisions". Journal of the European Economic Association 15: 173-217
- [131] Guerra, Gerardo, and Daniel J. Zizzo. 2004. "Trust Responsiveness and Beliefs". *Journal of Economic Behavior and Organization* 55: 25–30.
- [132] Gul, Faruk. 1991. "A Theory of Disappointment Aversion". *Econometrica* 59: 667-686.
- [133] Gul, Faruk, and Wolfgang Pesendorfer. 2016. "Interdependent Preference Models As a Theory of Intentions". *Journal of Economic Theory* 165: 179-208.
- [134] Hahn, Volker. 2009. "Reciprocity and Voting". Games and Economic Behavior 67: 467-480.
- [135] Halpern, Joseph, and Rafael Pass. 2012. "Iterated Regret Minimization: A New Solution Concept". Games and Economic Behavior 74: 194-207.
- [136] Harsanyi, John. 1967-68. "Games of Incomplete Information Played by Bayesian Players. Parts I, II, III". Management Science 14: 159-182, 320-334, 486-502.
- [137] Hauge, K. 2016. "Generosity and Guilt: The Role of Beliefs and Moral Standards of Others". *Journal of Economic Psychology* 54, 35-43.
- [138] Heifetz, Aviad, Martin Meier, and Burkhard Schipper. 2006. "Interactive Unawareness". *Journal of Economic Theory* 130: 78-94.
- [139] Herweg, Fabian, Muller, Daniel, and Weinschenk, Philipp. 2010. "Binary Payment Schemes: Moral Hazard and Loss aversion". *American Economic Review* 100, 2451-2477.

- [140] Huang, Peter. 2000. "Pandemic Emotions, Public Health, Financial Economics, and Law, and Leadership". University of Colorado Law Legal Studies Research Paper No. 20-14.
- [141] Huang, Peter, Ho-Mou. 1994. "More Order without More Law: A Theory of Social Norms and Organizational Cultures". *Journal of Law, Economics, and Organization* 12: 390-406.
- [142] Hupkau, Claudia, and François Maniquet. 2018. "Identity, Non-Take-Up and Welfare Conditionality". Journal of Economic Behavior and Organization 147: 13-27.
- [143] Inderst, Roman, Kyril Khalmetski, and Axel Ockenfels. 2019. "Sharing Guilt: How Better Access to Information May Backfire". *Management Science*.65: 2947-3448.
- [144] Isoni, Andrea, and Robert Sugden. 2019. "Reciprocity and the Paradox of Trust in Psychological Game Theory". *Journal of Economic Behavior and Organization* 167: 219-227.
- [145] Jagau, Stephen, and Andrés Perea. 2018. "Common Belief in Rationality in Psychological Games". Epicenter w.p. 10.
- [146] Jang, Dooseok, Amrish Patel, and Martin Dufwenberg. 2018. "Agreements with Reciprocity: Co-Financing and MOUs". Games and Economic Behavior 111: 85-99.
- [147] Jiang, Lianjie, and Jiabin Wu. 2019. "Belief-Updating Rule and Sequential Reciprocity". *Games and Economic Behavior* 113: 770-780.
- [148] Kahneman, Daniel, and Amos Tversky. 1979. "Prospect Theory: An Analysis of Decision Under Risk". *Econometrica* 47: 263-291.
- [149] Kartik, Navin. 2019. "Strategic Communication with Lying Costs". Review of Economic Studies 76: 1359-1395.
- [150] Kawagoe, Toshiji, and Yosuke Narita. 2014. "Guilt Aversion Revisited: an Experimental Test of a New Model". Journal of Economic Behavior and Organization 102: 1-9.

- [151] Keltner, Dacher, Jennifer Lerner. 2010. "Emotions." In Handbook of Social Psychology, Vol. 1, edited by Susan Fiske, Daniel Gilbert and Gardner Lindzey: pp. 317-52. New York: Wiley.
- [152] Khalmetski, Kyril. 2016. "Testing Guilt Aversion with an Exogenous Shift in Beliefs". Games and Economic Behavior 97: 110-119.
- [153] Khalmetski, Kiryl. 2019. "The Hidden Value of Lying: Evasion of Guilt in Expert Advice". *Journal of Economic Behavior and Organization* 167: 296-310.
- [154] Khalmetski, Kiryl, Axel Ockenfels, and Peter Werner. 2015. "Surprising Gifts: Theory and Laboratory Evidence". *Journal of Economic Theory* 159: 163-208.
- [155] Khalmetski, Kiryl and Dirk Sliwka. 2019. "Disguising Lies Image Concerns and Partial Lying in Cheating Games". American Economic Journal: Microeconomics 11: 79-110.
- [156] Kolpin, Van. 1992. "Equilibrium Refinements in Psychological Games". Games and Economic Behavior 4: 218-231.
- [157] Kőszegi, Botond. 2006. "Ego Utility, Overconfidence, and Task Choice". *Journal of the European Economic Association*. 4: 673-707.
- [158] Kőszegi, Botond. 2010. "Utility from Anticipation and Personal Equilibrium". *Economic Theory* 44: 415-444.
- [159] Kőszegi, Botond, George Loewenstein, and Takeshi Murooka. 2019. "Fragile Self-Esteem". Unpublished.
- [160] Kőszegi, Botond, and Matthew Rabin. 2006. "A Model of Reference-Dependent Preferences". Quarterly Journal of Economics 121: 1133-1166.
- [161] Kőszegi, Botond, and Matthew Rabin. 2007. "Reference-Dependent Risk Attitudes". American Economic Review 97: 1047-1073.
- [162] Kőszegi, Botond, and Matthew Rabin. 2009. "Reference-Dependent Consumption Plans". American Economic Review 99: 909-936.

- [163] Kozlovskaya, Maria, and Antonio Nicolò. 2019. "Public Good Provision Mechanisms and Reciprocity". Journal of Economic Behavior and Organization 167: 235-244.
- [164] Kreps, David, and Evan Porteus. 1978. "Temporal Resolution of Uncertainty and Dynamic Choice Theory". *Econometrica* 46: 185-200.
- [165] Kreps, David, and Robert Wilson. 1982. "Sequential Equilibria". Econometrica 50: 863-894.
- [166] Kunda, Ziva. 1990. "The Case for Motivated Reasoning". Psychological Bulletin 108: 480-498.
- [167] Lazear, Edward, Ulrike Malmendier, and Roberto Weber. 2012. "Sorting in Experiments with Application to Social Preferences". *American Economic Journal: Applied Economics* 4: 136-63.
- [168] Le Quement, Mark, and Amrish Patel. 2018. "Communication as Gift-Exchange". Unpublished.
- [169] Lerner, Jennifer, and Dacher Keltner. 2000. "Beyond Valence: Toward a Model of Emotion-specific Influences on Judgement and Choice". Cognition and Emotion 14: 473-93.
- [170] Lerner, Jennifer, and Dacher Keltner. 2001. "Fehr, Anger, and Risk". Journal of Personality and Social Psychology 81: 146-159.
- [171] Lerner, Jennifer, Ye Li, Piercarlo Valdesolo, and Karim Kassam. 2015. "Emotion and Decision Making". *Annual Review of Psychology* 66: 799-823.
- [172] Levine, David K. 1998. "Modeling Altruism and Spitefulness in Game Experiments". Review of Economic Dynamics 1: 593-622.
- [173] Livio, Luca, and Alessandro De Chiara. 2019. "Friends or Foes? Optimal Incentives for Reciprocal Agents". *Journal of Economic Behavior and Organization* 167: 245-278.
- [174] Loewenstein, George, Christopher Hsee, Elke Weber, and Ned Welch. 2001. "Risk as Feelings". *Psychological Bulletin* 127: 267-286.

- [175] Loewenstein, George, and Andras Molnar. 2018. "The Renaissance of Belief-based Utility in Economics". *Nature Human Behavior* 2: 166-167.
- [176] Loomes, Graham and Robert Sugden. 1982. "Regret Theory: An Alternative Theory of Rational Choice under Uncertainty". *Economic Journal* 92: 805-824.
- [177] Loomes, Graham and Robert Sugden. 1986. "Disappointment and Dynamic Consistency in Choice under Uncertainty". Review of Economic Studies 53: 271-282.
- [178] López-Pérez, Raúl. 2008. "Aversion to Norm-Breaking: A Model". Games and Economic Behavior 64: 237-267.
- [179] Mailath, George, and Larry Samuelson. 2006. Repeated Games and Reputations: Long-Run Relationships. Oxford University Press.
- [180] Mannahan, Rachel. 2019. "Self-Esteem and Rational Self-Handicapping". Unpublished.
- [181] Marcus-Newhall, A., Pedersen, W.C., Carlson, M., Miller, N. (2000). "Displaced Aggression is Alive and Well: A Meta-analytic Review". Journal of Personality and Social Psychology 78, 670–689.
- [182] Mauss, Marcel. 1954. The Gift: Forms and Functions of Exchange in Archaic Societies. Glencoe, Illinois: The Free Press.
- [183] Miettinen, Topi, and Sigrid Suetens. 2008. "Communication and Guilt in a Prisoner's Dilemma". *Journal of Conflic Resolution* 52: 945-960.
- [184] Möbius, Markus M., Muriel Niederle, Paul Niehaus, and Tanya S. Rosenblat. 2010. "Managing Self-Confidence: Theory and Experimental Evidence". NBER w.p. 17014.
- [185] Molnar, Andras, Shereen Chaudhry & George Loewenstein. 2020. "It's Not About the Money. It's About Sending a Message! Unpacking the Components of Revenge". Unpublished.
- [186] Morrell, Alexander. 2014. "The Short Arm of Guilt: Guilt Aversion Plays Out More Across a Short Social Distance". W.p. series of the Max Planck Institute for Research on Collective Goods 2014-2019.

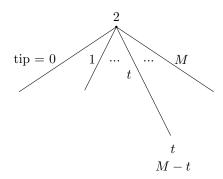
- [187] Morris Stephen. 2001. "Political Correctness". Journal of Political Economy 109: 231–265.
- [188] Myerson, Roger. 1986. "Multistage Games with Communication". Econometrica 54: 323-358.
- [189] Netzer, Nick, and Armin Schmutzler. 2014. "Explaining Gift-exchange The Limits of Good Intentions". Journal of the European Economic Association 12: 1586-1616.
- [190] Nielsen, Carsten, and Alexander Sebald. 2017. "Simple Unawareness in Dynamic Psychological Games" The B.E. Journal of Theoretical Economics 17: 1-29.
- [191] Nyborg, Karin. 2018. "Reciprocal Climate Negotiators". Journal of Environmental Economics and Management 92: 707-725
- [192] O'Donoghue, Ted, and Matthew Rabin. 1999. "Doing it Now or Later". American Economic Review 89: 103-124.
- [193] O'Donoghue, Ted and Charles Sprenger. 2018. "Reference-Dependent Preferences". In Douglas Bernheim Stefano DellaVigna David Laibson (eds.): *Handbook of Behavioral Economics*, Vol. 1, Elsevier.
- [194] Ottaviani, Marco, and Peter Sørensen. 2006. "Reputational Cheap Talk". The RAND Journal of Economics 37: 155–175.
- [195] Passarelli, Francesco, and Guido Tabellini. 2017. "Emotions and Political Unrest". *Journal of Political Economy* 125: 903-946.
- [196] Patel, Amrish, and Alec Smith. 2019. "Guilt and Participation". Journal of Economic Behavior and Organization 167: 279-295.
- [197] Persson, Emil. 2018. "Testing the Impact of Frustration and Anger When Responsibility is Low". Journal of Economic Behavior and Organization 145: 435-448.
- [198] Piccione, Michele, and Ariel Rubinstein. 1997. "On the Interpretation of Decision Problems with Imperfect Recall". *Games and Economic Behavior* 20: 3-24.

- [199] Potegal, Michael, Charles Spielberger, and Gerhard Stemmler. 2010. International Handbook of Anger: Constituent and Concomitant Biological, Psychological, and Social Processes. New York: Springer.
- [200] Quiggin, J. 1994. "Regret Theory with General Choice Sets". *Journal of Risk and Uncertainty* 8: 153-65.
- [201] Rabin, Matthew. 1993. "Incorporating Fairness into Game Theory and Economics". American Economic Review 83: 1281-1302.
- [202] Rauh, Michael, and Giulio Seccia, 2006. "Anxiety and Performance: A Learning-By-Doing Model". International Economic Review 47: 583-609.
- [203] Regner, Tobias, and N.S. Harth. 2014. "Testing Belief-Dependent Models". W.p. Max Planck Institute of Economics.
- [204] Rotemberg, Julio. 2005. "Customer Anger at Price Increases, Changes in the Frequency of Price Adjustment and Monetary Policy". *Journal of Monetary Economics* 52: 829-852.
- [205] Rotemberg, Julio. 2011. "Fair Pricing". Journal of the European Economic Association 9: 952-981.
- [206] Ruffle, Bradley. 1999. "Gift Giving with Emotions". Journal of Economic Behavior and Organization 39: 399-420.
- [207] Schotter, Andrew, and Isabel Trevino. 2014. "Belief Elicitation in the Laboratory". Annual Review of Economics 6: 103-128.
- [208] Sebald, Alexander. 2010. "Attribution and Reciprocity". Games and Economic Behavior 68: 339-352.
- [209] Sebald, Alexander, and Nick Vikander. 2019. "Optimal Firm Behavior with Consumer Social Image Concerns and Asymmetric Information". Journal of Economic Behavior and Organization 167: 311-330.
- [210] Sebald, Alexander, and Markus Walzl. 2015. "Optimal Contracts Based on Subjective Evaluations and Reciprocity". Journal of Economic Psychology 47: 62-76.

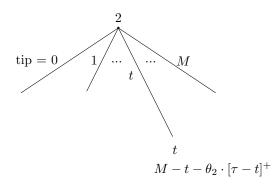
- [211] Selten, Reinhard. 1975. "Reexamination of the Perfectness Concept for Equilibrium Points in Extensive Games". *International Journal of Game Theory* 4: 25-55.
- [212] Shalev, Jonathan. 2000. "Loss Aversion Equilibrium". International Journal of Game Theory 29: 269-287.
- [213] Silfver, Mia. 2007. "Coping with Guilt and Shame: A Narrative Approach". *Journal of Moral Education* 36: 169-183.
- [214] Skinner, B.F. 1948. Walden Two. Englewood Cliffs, NJ: Prentice Hall.
- [215] Smith, Eliot R., and Diane M. Mackie. 2007. Social Psychology (Third ed.). Hove: Psychology Press.
- [216] Sohn, Jin and Wenhao Wu. 2019. "Reciprocity with Uncertainty about Others". Unpublished.
- [217] Sobel, Joel. 2005. "Interdependent Preferences and Reciprocity". *Journal of Economic Literature* 43: 396-440.
- [218] Tadelis, Stephen. 2011. "The Power of Shame and the Rationality of Trust". Unpublished.
- [219] Tangney, June Price. 1995. "Recent Advances in the Empirical Study of Shame and Guilt". American Behavioral Scientist 38: 1132-1145.
- [220] Thompson, Frederick. 1952. "Equivalence of Games in Extensive Form". Research memorandum, U.S. Air Force, Rand Corporation.
- [221] Trivers, Robert. 1971. "The Evolution of Reciprocal Altruism". Quarterly Review of Biology 46: 35-57.
- [222] Vanberg, Christoph. 2008. "Why Do People Keep Their Promises? An Experimental Test of Two Explanations". *Econometrica* 76, 1467-1480.
- [223] van Damme, Eric, et al. 2014. "How Werner Güth's Ultimatum Game Shaped our Understanding of Social Behavior". Journal of Economic Behavior and Organization 108: 292-318.

- [224] van Leeuwen, Boris, Charles Noussair, Theo Offerman, Sigrid Suetens, Matthijs van Veelen, and Jeroen van de Ven. 2018. "Predictably Angry – Facial Cues Provide a Credible Signal of Destructive Behavior". Management Science 64: 3352-3364.
- [225] Woods Daniel, and Maros Servatka. 2016. "Testing Psychological Forward Induction and the Updating of Beliefs in the Lost Wallet Game". Journal of Economic Psychology 56: 116-125.
- [226] Zeelenberg, Marcel. 1999. "Anticipated Regret, Expected Feedback and Behavioral Decision Making". *Journal of Behavioral Decision Making* 12: 93-106.
- [227] Zeelenberg, Marcel & Rik Pieters. 2007. "A Theory of Regret Regulation 1.0". Journal of Consumer Psychology 17: 3-18.

 $G_1$ 



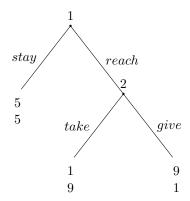
 $G_1^*$ 



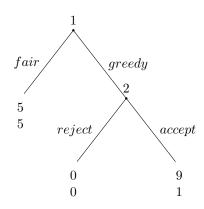
 $G_2$ 

	opera		boxing	
opera		1/2		0
-	1		0	
boxing		0		1
	0		1/2	

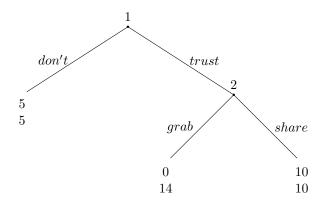
 $G_4$ 



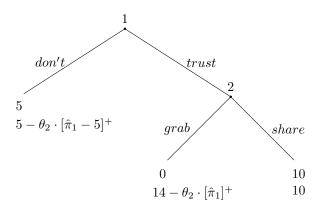
 $G_5$ 



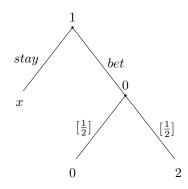
 $G_6$ 



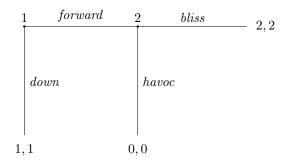




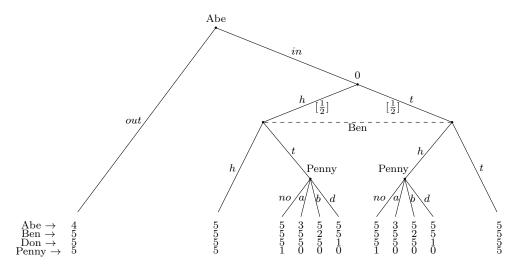
## $G_7$



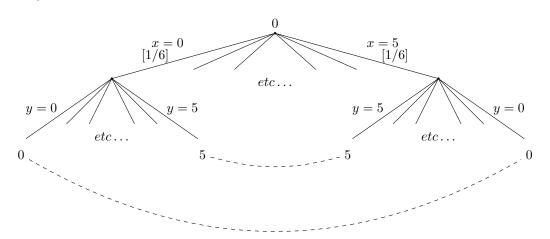
## $G_8$



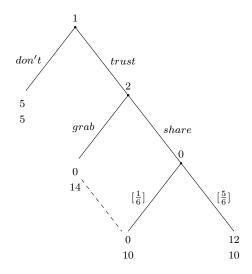




## $G_{10}$



 $G_{11}$ 



 $G_{12}$ 

